**317959**

**Mobile Opportunistic Traffic Offloading**

**D3.3.1 – Design and evaluation of enabling techniques for mobile data traffic offloading (release a, public)**

SEVENTH FRAMEWORK
PROGRAMME

| | |
|---|---|
| Grant Agreement No. | 317959 |
| Project acronym | *MOTO* |
| Project title | Mobile Opportunistic Traffic Offloading |
| Advantage | |
| | |
| Deliverable number | D3.3.1 |
| Deliverable name | Design and evaluation of enabling techniques for mobile data traffic offloading (release a) |
| Version | V 3.0 |
| | |
| Work package | WP3 – Offloading foundations and enablers |
| Lead beneficiary | CNR |
| Authors | Vania Conan (TCS), Filippo Rebecchi (TCS), Raffaele Bruno (CNR), Andrea Passarella (CNR), Elisabetta Biondi (CNR), Chiara Boldrini (CNR), Antonino Masaracchia (CNR), Giovanni Mainetto (CNR), Marcelo Dias de Amorim (UPMC), Filippo Rebecchi (UPMC), Engin Zeydan (AVEA), Ahmet Serdar Tan (AVEA), Eva Pierattelli (INTECS), Daniele Azzarelli (INTECS). |
| | |
| Nature | R – Report |
| Dissemination level | PU – Public |
| Delivery date | 31/10/2014 (M24) |

# Table of Contents

# List of Figures

# Executive summary

This document is the third deliverable of WP3, and reports on the activities and results obtained in the Tasks 3.2 and 3.3 during the second year. Activities on T3.1 are presented in a separate document, i.e. D3.2 [2], which is the logical output of T3.1 (now finished). Activities in WP3 have progressed along the methodology discussed already in D3.1 [1]. As far as capacity assessment is concerned (T3.2) we have both analysed the performance of individual building blocks in isolation, and their performance when they are combined in a complete offload networking solutions. Results of these activities include: (i) analysing convergence issues in opportunistic networks; (ii) providing end-to-end delay guarantees in opportunistic networks with duty cycling; (iii) assessing the performance of LTE through modelling; and (iv) assessing the performance of complete offload networks (also using infrastructure WiFi components in addition to cellular and opportunistic) in presence of both synchronised and non-synchronised content requests. Moreover, the document also reports results from T3.3 about scheduling. We have analysed scheduling from multiple dimensions. We have analysed both intra-technology and inter-technology scheduling issues. From the first standpoint, we have considered joint scheduling of multicast and D2D transmissions to optimise offloading. As far as inter-technology scheduling is concerned we have developed a general optimisation framework based on TOPSIS, in order to optimise allocation of users to the various possible technologies based on different QoS performance indices and criteria. Last but not least, we have analysed how to schedule various architectural components of an LTE network (i.e., pico and macrocells) in order to reduce the LTE energy consumption without compromising the efficiency in terms of throughput perceived by the users. In addition to presenting these results in detail, in Section 1 we remind the general strategy of activities in WP3 and how these results are aligned with it, and how they are synergic with the work undertaken in the rest of the project. At the end of the document, we discuss how WP3 activities are progressing based on these results.

# 1  Introduction

This deliverable reports on the main activities during the second year of the project in T3.2 and T3.3 of WP3. Task 3.1, focusing on the study of spatio-temporal contact patterns, was finished on M18, and the main outputs have been reported in D3.2 [2]. Activities described in this document, therefore, focused on two main broad topics. The first one is about assessing the capacity of offload networking solutions (with main focus on supporting terminal-to-terminal communication). The second one is about scheduling solutions, both inside single network technologies and across different technologies. Before summarising the key results achieved on these topics, let us recall the main objectives of the WP (related to T3.2 and T3.3), and how we are addressing them through the results presented in this document.

## 1.1  Problem statement: Objectives of the WP and approach in addressing them

As described in the DoW, the objectives of WP3 related to T3.2 and T3.3 are as follows:

1. To quantify capacity improvements that can be achieved when offloading traffic across different wireless infrastructures and/or using terminal-based offloading, in both single- and multi-operator environments.

2. To characterize the impact on offloading efficiency of factors such as user mobility patterns, heterogeneity of network deployments, traffic loads, QoS application requirements, and variable terminal densities due to distributed duty cycling techniques.

3. To develop inter-technology scheduling algorithms allowing a more efficient synergy – in presence of offloading techniques – between multiple wireless infrastructures and opportunistic networks, which a special focus on high-load conditions.

The first two objectives are addressed by T3.2, while the third one is addressed by T3.3. As far as the first objective is concerned, the rationale of the work undertaken is as follows:

1.1.  We have identified the main architectural blocks to be analysed as far as capacity improvements are concerned from the architecture definition in WP2. At the high level, these are (i) wireless broadband infrastructures and (ii) opportunistic networks.

1.2.  We analyse capacity limits of wireless infrastructures, primarily focusing on LTE cellular networks. Although this is not the main focus of the WP, in some cases we also study modifications of LTE components that can overcome some of these limits.

1.3.  We characterise the capacity that opportunistic networks can bring about. To this end, we consider a wide range of opportunistic networking protocols and users' mobility patterns, and assess capacity (in terms of throughput and/or end-to-end delay) as a function of these key elements.

1.4.  We consider both architectural blocks together, i.e. we study the capacity of an integrated heterogeneous network composed of both a wireless broadband infrastructure and opportunistic network. We study the actual capacity gain that can be achieved when these networking environments are put together.

With respect to the second objective, the rationale of the work undertaken is as follows:

2.1.  We have derived configurations for evaluation of the capacity of networks with offloading from the factors identified in the objective, primarily: (i) mobility patterns; (ii) heterogeneity of networks and of users' mobility; (iii) end-to-end delay requirements; (iv) contact patterns modifiers related to energy efficiency (duty cycling)

2.2.  We analyse how these factors impact on the capacity limits of wireless infrastructures, and on the capacity improvements brought by opportunistic networks.

Finally, with respect to the third objective, the rationale of the work undertaken is as follows:

3.1. We have identified relevant scheduling problems for the different MOTO scenarios. As discussed in the first period review meeting, this has lead to re-focusing some of the activities, that now take into consideration also intra-technology scheduling. The resulting lines of activities are as follows: (i) intra-technology scheduling in LTE to jointly exploit multicast and D2D communications; (ii) intra-technology scheduling in LTE to improve energy efficiency in the access network; (iii) inter-technology scheduling to optimize allocation of users to multiple wireless technologies available at the same time.

3.2. For each of the three lines of research, we identify specific research problems and address them. We propose algorithmic solutions to improve the efficiency of MOTO for each of them.

## 1.2 Enabling techniques for mobile data traffic offloading: A summary

Figure 1 provides a graphical representation of the main activities undertaken during the reporting period in WP3, and particularly in Tasks 3.2 and 3.3. Note that for each activity we also highlight the main methodological approach followed, consisting either of analytical modelling, definition of algorithms to improve capacity, or simulation-based analysis (or combinations thereof).



**Figure 1. Schematic representation of WP3 activities.**

The following subsections provide a summary of these activities, and highlight the key results achieved. Moreover we also present as a brief summary of the improvements for each activity with respect to the status at the end of the previous reporting period.

### 1.2.1 Task 3.2: Capacity limits and improvements in networks with offloading

Task 3.2 is devoted to better understand the capacity limits of LTE, and assess the capacity improvements that can be achieved through offloading. From a *methodological* standpoint, we take primarily an approach based on analytical modelling and simulation analysis. Simulation is used to explore the performance of specific systems or networking solutions. Analysis is also used for this purpose (in a pretty standard way with respect to the field of performance evaluation). In addition, it is also used to provide compact mathematical tools that can be used by network operators to plan how to dimension their network in presence of offloading. With respect to the specific *subjects of investigation*, we go step-by-step. At the high level, an offload network is made up of two main components, i.e. a wireless infrastructure part (primarily, LTE), and a mobile part (primarily, an opportunistic network). Therefore, several topics of research deal with assessing the performance and the capacity of these two building blocks in isolation.

This is important, as there are still several open points in understanding the capacity limits of these types of network alone.

With respect to **LTE networks**, while a number of studies have been carried out to characterise its capacity at the physical layer, little effort has been devoted in analysing the capacity perceived by the users, i.e. a number of network layers and functional blocks away from the basic physical layer. In T3.1 we are contributing to fill this gap, and this document presents some results on this.

---

*Specifically, in Section 3 we present the key results achieved with respect to this point, i.e.:*

- *we provide an initial model of the throughput experience by users of LTE networks, when we factor in not only parameters of the physical layer, but also the key mechanisms of retransmissions and data reliability*

- *we exploit learning mechanisms to extend LTE algorithms for automatic configurations of transmission parameters, making them adaptive to changing network conditions.*

---

With respect to **opportunistic networks**, research is still ongoing to model the end-to-end throughput (or, equivalently, the end-to-end-delay, as explained below) of opportunistic networks. As summarised in the following of this section, this document presents results on convergence properties of opportunistic networks, and on modelling of end-to-end delay in presence of energy saving mechanisms such as duty cycling (we will come back on the energy saving dimension afterwards in this section). Specifically, in Section 2 we first focus on the issue of convergence of forwarding protocols in opportunistic networks, where a protocol is convergent if it yields expected finite end-to-end delay. While at a first sight this may seem a very theoretical problem, it has significant practical implications. Diverging, in practise, means loosing packets, and not be able to provide and end-to-end delay guarantee. Unfortunately, analysis of real mobility traces has shown that protocols may indeed diverge, depending on the contact patterns between nodes. Then, we study the end-to-end delay of forwarding protocols, in case of exponential contact patterns, when duty cycling is used at nodes to conserve energy. This piece of work directly exploits results previously presented in D3.1 (about how to model end-to-end delay in with exponential contact patterns) and D3.2 (about the effect of duty cycling on temporal contact patterns).

---

*The key results achieved with respect to this point, are:*

- *we provide practical tools to select appropriate opportunistic networking protocols given a stochastic description of the contact patterns between nodes, such that convergence can be guaranteed. Moreover, we extend what presented in D3.1, by considering a much more vast family of routing protocols, and comparing them;*

- *we provide probabilistic end-to-end delay guarantees when a certain duty cycling is used*

- *we provide practical tools to select the optimal trade-off between the delay that can be guaranteed and the maximum energy saving that can be achieved while still achieving that delay.*

---

Advancing what presented in D3.1, in this document we also provide some results on the **capacity of offload networks, i.e. considering the two main building blocks working together**. Note that both level of analysis (considering individual blocks, and the two blocks together) are considered important outcomes of T3.2. Isolated analysis allows us to understand the performance of single technologies, and therefore provide tools to compare possible alternatives when configuring the two main building blocks. Joint analysis allows us to understand the interplay and the complementarity between the two blocks. We provide two main contributions along those lines, in Section 4. Remember that in D3.1 we have presented Push&Track (and its evolution, Droid), as one of the reference solutions in MOTO to implement offloading through opportunistic networks. In this document we extend Droid to use all the different networks that may be available in the MOTO context, i.e. cellular, WiFi and opportunistic. Then, in Section 4.2 we start investigating a relatively untapped problem in the offloading networking panorama, i.e. the offloading performance when the same content is requested in a non-synchronised way by users (note that results in

D3.1 and in Section 4.1 focus on simultaneous requests, as the vast majority of the literature on offloading networks).

---

*The key results achieved with respect to this point, are:*

- *through simulations based on real mobility trace and real WiFi deployments we compare the offloading performance when WiFi, cellular and opportunistic networks are used partly or all together at the same time (which clearly results in the best offloading performance).*

- *in case of offloading with non-synchronised requests we find that even with unfavourable configurations of the opportunistic networking protocols (i.e., those where resources of mobile devices are used at the minimum possible level), offloading can be very effective, saving up to 90% of the total traffic from flowing on the cellular network.*

---

### 1.2.2 Task 3.3: Scheduling issues in networks with offloading

Task 3.3 focuses on scheduling policies at various levels. We have identified *three main threads of activities* within the scheduling topic. As discussed at the first review meeting, this was the outcome of a restructuring of the activities originally planned for T3.3, which should have focused exclusively on inter-technology scheduling. While this topic is still investigated in T3.3, we have started two more lines of research that were considered equally important.

First, we study how to **schedule multicast vs. terminal-to-terminal transmissions in an offload network** (Section 5). This is an example of solutions that can be seen both as intra-technology scheduling (e.g., if LTE-D2D is used) and as inter-technology scheduling (when opportunistic networking solutions are used for terminal-to-terminal communication). This piece of work is quite interesting, as it start clarifying one of the basic questions of the overall offloading approach, i.e. if multicast wouldn't be enough.

---

*Our results show that, even in presence of synchronised content requests, multicasting alone is not as efficient as using multicast and D2D communication in an optimised way.*

---

A second thread of activity is related to **scheduling picocells** in the emerging scenario of dense networks, **where LTE macro cells are complemented by picocells for increase capacity at the edge of the cell**. While this in principle an interesting and effective approach, it may significantly increase energy consumption in the core network, as multiple eNBs need to be powered. We therefore asked ourselves if (i) macrocells can be run in a more efficient way in presence of picocells, by reducing their transmit power, and (ii) if picocells can be selectively switched off, concentrating traffic on fewer, better utilised, picocells. Initial results presented in Section 6 show that the answer is positive in both cases.

---

*The key results achieved with respect to this point, are:*

- *macrocells can be operated at reduced transmit power, without compromising the throughput perceived by the users thanks to the additional capacity provided by picocells (and the net power consumption is lower when the solution with picocells is adopted)*

- *picocells can be switched off, without impacting on the throughput perceived by users, which remains at the same level.*

---

Note that this activity is synergic with the work on end-to-end delay in presence of duty cycling, presented in Section 2.2. Specifically, the two works together provide initial results on how to operate offload networks in an energy efficient way, considering both the energy consumed in the core, and the energy consumed on the users' mobile devices.

Finally, in Section 7 we present results on the third activity in Task 3.3, related to **inter-technology scheduling**. We are using a standard optimisation framework (TOPSIS) to optimally allocate users to technologies, when multiple technologies are available at the same time. The framework is customised based on a number of QoS metrics, and modified to find the optimal allocation considering overall capacity performance goals (instead of maximising individual nodes benefit). Specifically, for now we consider on

the case where operators run both a cellular and a WiFi network, but we briefly discuss how we are extending this work to put also terminal-to-terminal protocols into the picture.

> ***Our results show that by using the TOPSIS framework, together with appropriate global throughput maximisation functions, we can effectively schedule users across multiple wireless technologies, avoiding cellular network overload.***

### 1.2.3 Progress with respect to Y1 activities

All in all, these results significantly progress the status of the work with respect to how it started at the end of the first year:

- The overall methodology of the work in T3.2 was confirmed, and activities have extended what presented in D3.1 in the first place. This is the case of (i) results about convergence of opportunistic networks; (ii) results about the end-to-end delay in opportunistic networks; (iii) results about the capacity of LTE networks in terms of throughput; (iv) results about the performance of Droid in joint cellular, WiFi and opportunistic networks

- In addition, we have started to analyse more comprehensively LTE networks by developing a user-oriented throughput model, and to analyse more in depth the capacity advantage of offloading in a complete offloading solution with non-synchronised content requests. Finally, we have started activities on scheduling producing several interesting results.

## 1.3  WP3 activities in the overall framework of the project

Due to the nature of how WP3 was planned, activities and results provide a library of solutions that can be composed together when needed, in order to assess the capacity of offload network either in operation or at the design stage. To give one concrete example of how such results can be combined and used, let us consider the case of an operator wishing to plan its network in case of opportunistic offloading with energy saving at the users' devices. Results presented in D3.2 tell how the process of inter-contact times between nodes is modified according to the use of duty cycling. Then, criteria presented in Section 2.1 of this deliverable can be used to select which specific protocol should be used, to avoid divergence problems. Depending on the application to be supported, the operator can identify the requirements in terms of end-to-end delay that need to be guaranteed. Therefore, models in Section 2.2 can be used to tune energy saving to meet these constraints. Finally, the operator would use a Push&Track (or Droid) system to implement the offloading solution. Based on the guaranteed end-to-end delay, it will be clear which fraction of the traffic can be expected by the operator to be served through the opportunistic network, thus finally obtaining an estimate of the additional capacity that would be gained through the so-configured offloading process.

In addition, the WP3 work presented in this deliverable is well aligned with the overall flow of activities of the project, and specifically with WP2 (architecture), WP4 (protocols), WP5 (performance evaluation). Specifically, all capacity building blocks and proposed algorithms are totally inline with the MOTO reference architecture defined in D2.2.1 [3]. Specifically, in the architecture we have clearly separated blocks dealing with how to manage and use wireless infrastructure (LTE and WiFi), and opportunistic networks. The separation of the capacity work in T3.2 is aligned with this architectural separation, and results can thus be fed back to the modules corresponding to these architectural blocks. Moreover, the integrated studies of offloading network follow the general design features of Push&Track, which implements all the basic architectural building blocks identified in WP2. Finally, activities on intra-technology scheduling naturally fit into the corresponding architectural elements, while inter-technology scheduling solutions are amenable to be implemented in the control blocks of the architecture dealing with orchestration between multiple wireless technologies. With respect to WP4, WP3 provides identifies initial solutions and alternatives, to be then analysed more precisely in the framework of the specification of networking protocols in WP4. An example is the work undertaken in WP4 about resource limitations in opportunistic network. This is based on the basic algorithmic tools of Push&Track, whose performance are assessed in WP3, and on the

selection of appropriate configurations for opportunistic protocols, available after the analysis in T3.2. Finally, with respect to WP5, in WP3 we evaluate either analytically or by simulation individual solutions, to better identify which technical solutions to integrate in the testbeds or in the integrated simulation platform. This is achieved often through simplified simulation models and scenarios, where we abstract (with respect to WP5) some characteristics, to be able to have quicker and "more agile" results about the performance of selected offloading building blocks.

In the rest of this document we expand on the concepts described above, and present the detailed results produced in WP3 during the second year of the project. Finally, in Section 8 we give a preview of open issues, and how we are addressing them.

**As a note to the reader**, each section starts with a summary of the content presented herein. Then, subsections present these activities in some more details. When appropriate, we omit technical details that will make the presentation too long. In these cases, Appendices are provided where all details are available. **Therefore, the document is structured so that it can be read at multiple levels of details**. The first parts of the sections are sufficient to understand the content in terms of the approach taken and the main results achieved. The rest of the sections go in more details presenting results and how they have been achieved. Appendices contain the rest of the details.

## 2 Capacity analysis: Assessing capacity of opportunistic networks

This section presents two main contributions, focused on the assessment of the capacity of opportunistic networks. Both contributions build upon, and significantly extend, results presented in previous documents, specifically in D3.1 (Initial results on offloading foundations and enablers) [1] and D3.2 (Spatiotemporal characterization of contact patterns in dynamic networks) [2].

In Section 2.1 we present results on convergence of opportunistic networking protocols. Remember that in our overall strategy, this is one of the necessary steps to characterise the capacity of the opportunistic network, because it allows us to identify configurations of the network (i.e., patterns of contacts between users) for which specific protocols do not converge. In these cases, the capacity gained by offloading would be zero, as protocols would yield infinite expected delay (more practically, messages will be lost and never delivered to the destination). Initial results on this topic were presented in Section 4.1 of [1]. In Section 2.1 we summarise the additional results we have obtained (complete details are available in Appendix A, which is a reprint of [11]). Specifically, we have fully characterised the convergence properties of both randomised (or social-oblivious) and social-aware routing protocols in case of Pareto distributed inter-contact times. Social-oblivious protocols do not use any contextual information about the behaviour (and thus resulting contact patterns) of users, while social-aware protocols are built to exploit such knowledge. The set of protocols considered cover the vast majority of forwarding algorithms defined in the opportunistic networking literature, while considering Pareto inter-contact times is inline with well-established results in the literature about the analysis of real mobility traces. Our results allow us to draw a pretty interesting and useful set of conclusions. Within each class of routing protocol (social-oblivious, social-aware) we are able to identify best solutions, i.e. those that guarantee convergence in most of the cases. Comparing best solution of each class, we find that there is not a unique winner, but the best overall choice actually depends on the pattern of contacts between nodes. It is particularly interesting to find that using social-oblivious protocols may yield convergence in some cases where social-aware would not.

In Section 2.2 we deal with end-to-end delay models in opportunistic networks in presence of duty cycling, for the case of exponential inter-contact time patterns (complete details are provided in Appendix A, which is a reprint of [8]). Specifically, we provide stochastic guarantees on the end-to-end delay as a function of the duty cycle period used by nodes. In other words, our results allow an operator to set the optimal duty cycling so that the end-to-end delay is below a given threshold with a given probability. This model provides an analytical tool to set the trade-off between the additional capacity that can be provided and the corresponding energy cost (in terms of battery of users' mobile devices). This result derives from two previous results achieved in the project. On the one hand, we exploit the results on how duty cycling modifies the patterns of *useful* contact between nodes in case of exponential contact patterns, where useful contacts are those that can be used to forward messages, i.e. those that occur when both nodes are active (not sleeping). These results have been presented in Section 4 of [2]. On the other hand, we exploit the models of end-to-end delay in case of exponential inter-contacts, presented in Section 4.2 of [1].

We refer the reader to Section 8 of this document for a discussion on how we will complete these activities, according to the overall logical structure of the WP activities described in Section 1.

### 2.1 Convergence of opportunistic networking protocols

Modelling the performance of social-oblivious and social-aware forwarding protocols for opportunistic networks is still an open research issue. Knowing the distribution of intermeeting times and the rules applied by the forwarding algorithm used in the network, one could - in principle - model the distribution of the delay experienced by messages and compute its expectation. In practice, modeling analytically the delay of the various forwarding protocols for general distributions of inter meeting times is very hard, and models exist only for some specific cases, typically assuming exponential intermeeting times [58][59][37][60][49][29]. A related modelling challenge is to assess the convergence of routing protocols, i.e. whether a specific protocol yields finite or infinite expected delay. Assessing convergence allows us to

understand whether a particular protocol can be safely used or not given a pattern of intermeeting times and how to configure it so that it converges, if possible. Although less informative than a complete delay model, convergence models can be derived for a large class of routing protocols releasing the exponential intermeeting time assumption.

The convergence of the expected delay is not guaranteed in all cases in which the expectation of the inter-meeting times may diverge. In fact, being the delay the result of the composition of the time intervals between node encounters, depending on the convergence of inter- meeting times, the expectation of the delay itself might diverge. This can happen, for example, when intermeeting times feature a Pareto (also known as power law) distribution, as first highlighted in [20]. The problem with Pareto distributions is that their expectation is finite only for certain values of their exponent $\alpha$. More specifically, the expectation is finite if $\alpha > 1$, while for $\alpha \leq 1$ it diverges to infinity. The first to postulate the existence of Pareto intermeeting times in real mobility scenarios (i.e., analyzing real traces of human mobility) were Chaintreau et al. in their seminal work in [20]. The relevance of Pareto intermeeting times in opportunistic networks is both theoretical and empirical. Cai and Eun [18] have mathematically derived that heavy-tailed intermeeting times can emerge depending on the relationship between the size of the boundary of the considered scenario and the relevant timescale of the network, showing that, at least in principle, Pareto intermeeting times are something that one may be faced with when studying opportunistic networks. Empirical evidence for the presence of Pareto intermeeting times was first suggested by [20], but it has been later criticised, arguing that the tail of the distribution is in fact exponential (e.g., [33]). Typically, these results are derived focusing on the aggregate inter-contact time distribution, while convergence depends on pairwise distributions. As proved in [47], the aggregate and pairwise distributions can be in general very different, and therefore analysis of pairwise inter-contact times are necessary, which are however mostly missing in the literature. To address this issue, we have performed a pairwise hypothesis testing on three popular publicly available contact datasets (Cambridge, Infocom'05, and RollerNet) and we have found that the Pareto hypothesis for intermeeting times cannot be rejected for 80%, 97%, and 85.5% of pairs, respectively. We believe that these results provide a strong case for Pareto intermeeting times in opportunistic networks and substantially motivate analyses like the one presented in this paper.

Under the Pareto intermeeting times assumption, in this work we derive the stability region (i.e., the Pareto exponent values of pairwise intermeeting times for which finite expected delay is achieved) of a broad class of social-oblivious and social-aware forwarding protocols (single- and multi-copy, single- and multi-hop). The starting point of our paper is the work by Chaintreau et al. [8], where such conditions have been studied for the two-hop scheme under the assumption of homogeneous mobility (i.e., i.i.d. intermeeting times across all pairs). However, measurement studies [21] [20] have shown that real networks are intrinsically heterogeneous. Thus, in this work, we investigate whether heterogeneity in contact patterns helps the convergence of the expected delay of a general class of social-oblivious and social-aware forwarding protocols, and whether convergence conditions can be improved using multi-copy strategies and/or multi-hop paths.

Overall, the key findings we obtain from this analysis are as follows:

• For *social-oblivious strategies*, if convergence can be achieved, *two hops are enough* for achieving it.

• Using n *hops* can help *social-aware schemes*, and make them converge in some cases when all other social- aware or social-oblivious schemes diverge.

• In both the social-oblivious and the social-aware case, we find that *multi-copy strategies* can achieve a finite expected delay even when single-copy strategies cannot.

• Comparing *social-oblivious and social-aware multi-copy solutions*, we are able to prove mathematically that there is no clear winner between the two, since either one can achieve convergence when the other one fails, depending on the underlying mobility scenario.

A concise presentation of these findings is provided in the following subsections. All details are available in Appendix A.

## 2.1.1 Social-oblivious protocols

The analysis of convergence for social-oblivious protocols was presented already in [1], and is briefly summarized here for the reader's convenience.

To accurately represent the different variants in this class, we identify three main groups, differing in the number of hops allowed between source and destination, the number of copies generated, and whether the source and relay nodes keep track of the evolution of the forwarding process or not. First, forwarding strategies can be single-copy or multi-copy. In the former case, at any point in time there can be at most one copy of each message circulating in the network. In the latter, multiple copies can travel in parallel, thus in principle multiplying the opportunities to reach the destination (we assume that all copies are generated by the source node). Second, forwarding protocols can be classified based on the number of hops that they allow messages to traverse, or, in other words, based on a TTL computed on the number of hops. When the number of allowed hops is finite, the last relay can only deliver the message to the destination directly. Third, the amount of knowledge that each agent in the forwarding process can rely on (or is willing to collect and store) is an additional element for classifying forwarding strategies. Focusing on the source node, there can be social-oblivious strategies in which the source node does not keep track at all of how the forwarding process progresses. In this case, considering the configuration in which the source node can generate up to m copies of the message, the m copies might end up being all distributed to the exact same relay, thus eliminating the potential benefits of multi-copy forwarding. A memoryful source, instead, is able to guarantee to use distinct relays. A similar problem holds for intermediate relays. Memoryless relays can forward the message to the same next hop more than once, because they are not at all aware of what happened in the past. On the other hand, memoryful relays possess this knowledge, and are able to refuse the custody of messages that they have already relayed.

The following conditions are found for convergence of these protocols

| | 1 hop | | 2 hops | | $n$-hop | |
|---|---|---|---|---|---|---|
| | 1 copy | $m$ copies | 1 copy | $m$ copies | 1 copy | $m$ copies |
| memoryless | $\alpha_{sd} > 2$ | - | [C1,C2] | [C1,C2] | [C1,C2] | [C1,C2] |
| memoryful source | - | - | - | [C3,C4] | - | [C1,C2] |
| memoryful relays | - | - | - | - | [C1,C2] | [C1,C2] |

**Table 1. Convergence conditions.**

Specifically, conditions C1 through C4 in the table are defined as follows

**C1** $\sum_{j \in \mathcal{P}_s} \alpha_{sj} > 1 + |\mathcal{P}_s|$, where $\mathcal{P}_s$ denotes the set of all nodes that can be encountered by node $s$;

**C2** $\alpha_{jd} > 2, \forall j \in \mathcal{P}_s - \{d\}$.

$C3 = m \leq m^*$

$C4 = \sum_{j=N-m}^{N-1} \alpha'_j > 1 + m$

where *s* and *d* denote the source and destination nodes, respectively, *m* denotes the number of copies generated by the source, *m\** is defined as follows

$$m^* = \begin{cases} 0 & if \sum_{j \in \mathcal{P}_s} \alpha_{sj} \leq N \\ \arg\max_m \{m + \sum_{i=m}^{N-1} \alpha_i^* > 1 + N\} & o.w. \end{cases}$$

and $\alpha_i^*$ denotes the *i*-th largest $\alpha_{sj}$ with $j \in P_s$.

It is useful to focus on one specific case, in order to clarify how these conditions can be interpreted. Let us consider the 2-hop 1-copy memoryless scheme, which converges iff conditions C1 and C2 are met. The physical meaning of the conditions is quite intuitive. Recall that in the 2-hop 1-copy scheme the source hands over the only copy of the message to the first encountered node, which then has to relay it directly to the destination. Condition C1 guarantees that the first phase occurs with a finite expected time.

Specifically, the source node encounters the first possible relay with a time that is distributed according to a Pareto law with shape $\sum_{j \in P_s} \alpha_{sj} - |P_s|$. Therefore, the first phase "converges" if the average value of this time is finite, which leads to condition C1. Condition C2 guarantees that whatever relay is chosen by s, it encounter the destination within a finite expected time (note that the time for such relay to meet the destination is the residual of their intermeeting time, as the process of encounter between nodes is asynchronous, and therefore node s meets the relay at a random point in time with respect to the meetings between the relay and the destination).

Note that conditions C3 and C4 are needed only in case of multi-copy forwarding. The value *m\** is a threshold on the number of copies, such that if the source generates up to *m\** copies, all of them are handed over to m\* distinct relays with finite expected delay, while if *m* exceeds *m\** the additional copies cannot be handed over with finite expected delay. Condition C3 thus imposes that the source can actually relay m distinct copies of the message, while condition C4 guarantees that the destination meets at least one of the used relays with finite expected delay.

## 2.1.2 Social-aware protocols

### 2.1.2.1 Definition of social-aware forwarding protocols

Due to the variety of social-aware schemes available in the literature, here we only consider an abstract social-aware protocol that measures how good a relay is for a given destination in terms of its *fitness*. The fitness $fit^d_i$ is assumed to be a function of how often node i meets the destination d, thus $fit^d_i$ can be taken as proportional to the rate of encounter between node i and the destination. $1/E[M_{id}]$.

Under this abstract and general social-aware strategy, upon encounter, a node i can hand over the message to another node j only if its fitness is lower than the fitness of the peer, i.e., if $fit^d_j > fit^d_i$ holds (in the following we drop superscript d). The fitness function considered here uses only information on contacts between nodes, which have a direct dependence on the intermeeting time distribution. This lets us clearly show what is the impact of the contact dynamics on the performance of opportunistic forwarding protocols. How such simple fitness function can be extended to more complex forwarding strategies has been discussed in [12]. In the following, we denote with $R_i$ the set of possible relays for node i, i.e., the set of nodes whose fitness is greater than that of node i. Therefore, with social-aware forwarding, nodes can hand over a message only to nodes with higher fitness. This means that the set of potential relays shrinks as the message is handed over from hop to hop towards the final destination. This is pictorially represented in Figure 2.



**Figure 2. Schematic representation of social-aware forwarding.**

### 2.1.2.2 Convergence conditions for social-aware protocols

Based on this definition of social-aware forwarding protocols, we have been able to analyse different families of protocols, again distinguishing between single- and multi-hop protocols, and between single- and multi-hop protocols. In all cases we consider memoryful protocols, as it does not make much sense to

assume memoryless protocols when they have already to store information about contact patterns between nodes.

The corresponding convergence conditions are provided in Table 2 (note that the definition of conditions is slightly modified with respect to [11] for the sake of clarity, in order to match conditions for social-oblivious protocols in Table 1).

| | 1 hop | | 2 hops | | $n$-hop | |
|---|---|---|---|---|---|---|
| | 1 copy | $m$ copies | 1 copy | $m$ copies | 1 copy | $m$ copies |
| social-oblivious | $\alpha_{sd} > 2$ | - | [C1,C2] | [C3,C4] | [C1,C2] | [C1,C2] |
| social-aware | $\alpha_{sd} > 2$ | - | [C5,C6] | [C9] | [C7,C8] | [C7,C8] |

**Table 2. Convergence conditions for both social-oblivious and social-aware schemes**

We hereafter only sketch the intuitive meaning of one such condition, by focusing on the 1-copy 2-hop scheme (and refer the reader to Appendix A for the details). This protocol converges if and only if conditions [C5,C6] are met, which are defined as follows

**C5** $\quad \sum_{j \in \mathcal{R}_s} \alpha_{sj} > 1 + |\mathcal{R}_s|$

**C6** $\quad \alpha_{jd} > 2, \forall j \in \mathcal{R}_s - \{d\}.$

where Rs denotes the set of nodes encountered by the source with higher fitness (than itself) towards the destination. Conditions C5 and C6 maps exactly conditions C1 and C2 of the social-oblivious case. Specifically, they state that the expected delay in the 2-hop 1-copy case if finite iff the source encounters one of the possible relays (i.e., one node in Rs) in a finite amount of time (condition C5), and each of these nodes encounters the destination in a finite amount of time (condition C6). Conditions C7,C8 and C9 can be derived for the other protocols using a similar line of reasoning (see Appendix A).

As in the social- oblivious case, multi-hop schemes do not benefit from the use of multiple copies, and in fact the 1-copy n-hop scheme and the m-copy n-hop scheme share the same convergence conditions. Similarly, the difference between 2-hop schemes mirrors that between the corresponding social-oblivious versions. Thus, the 1-copy 2-hop scheme is effective when $\alpha_{jd} > 2$ for all j $\in$ Rs, since it allows us to save resources by sending a single copy. However, when conditions [C5,C6] do not hold, the only chance to achieve convergence is to exploit multiple copies.

If we focus on single-copy schemes, it is interesting to note that, differently from the social-oblivious case in which using additional hops did not provide any advantage, 1-copy social-aware schemes may benefit from multiple hops. In fact, for the 1-copy 2-hop scheme we need to impose that all intermediate relays j meet the destination with $\alpha_{jd} > 2$, (conditions C6) which is a quite strong condition. On the other hand, if we use multiple hops (1-copy n-hop case), conditions C7 and C8 are required, which are milder than C5. Their definition requires several steps, and therefore we don't report them here for the sake of simplicity (see Appendix A for the details). Basically, the only constraint for the 1-copy n-hop case is that there must be at least one node z (the one with the highest fitness) meeting the destination with $\alpha_{zd} > 2$.

Finally, we compare the m-copy 2-hop case with the 1-copy n-hop case (which is equivalent to the m-copy n-hop scheme). There is no clear winner here, as each scheme can provide convergence when the other one cannot. For example, consider the case in which the source node is not able to send more than one copy within a finite amount of time. In this case, the m-copy 2-hop scheme becomes effectively a 1-copy 2-hop scheme, which fails to achieve convergence if some intermediate hop j does not have exponent αjd greater than 2 (condition C6). Instead, exploiting multiple hops pays off in this case, as it allows us to rely on more intermediate relays, which may not meet the destination within a finite expected time but can bring the message "closer" to nodes that do meet d with $\alpha_{jd} > 2$. Vice versa, when the source node can hand over multiple copies (m>1) within a finite delay, the cooperative delivery of the multiple copies can overcome the presence of intermediate relays for which conditions C8-C9 do not hold. For example, when there is not even one relay j with $\alpha_{jd} > 2$, then the m-copy 2-hop case is the only possible choice.

### 2.1.3 Comparing social-oblivious and social-aware schemes

In the following we take the champions of each class and we investigate whether there is a clear winner between social-oblivious and social-aware strategies when it comes to the convergence of their expected delay.

Let us first consider the case $\alpha_{sd} > 2$. With this configuration the Direct Transmission scheme is the best choice from the convergence standpoint. In fact, with social-oblivious schemes using more than one hop, "bad" relays can be selected even starting from a source that is already able to reach the destination with a finite expected residual intermeeting time. This does not happen with social-aware strategies. In fact, assume that the source is the only node with $\alpha_{sd} = 2 + \varepsilon$, while all other nodes meet the destination with $\alpha_{jd} = 1 + \varepsilon$, with $\varepsilon$ being a very small quantity. In the social-aware case, Rs contains only the destination, as all other nodes are clearly worse than the source node as relay. This shows the adaptability of social-aware schemes: the additional knowledge that they exploit makes them able to resort to simpler approaches (in this case, Rs = {d} is equivalent to the Direct Transmission) when they realize that additional resources in terms of number of copies or number of hops would not help the forwarding process. This implies that one can safely use the m-copy 2-hop or the 1-copy n-hop social-aware protocols because in the worst case they will do no harm (they will downgrade to simpler strategies, without exploiting wrong paths), while in the best case they are able to improve the convergence of the forwarding process.

When $\alpha_{sd} \leq 2$ and $\alpha_{jd} > 2$ for all nodes j in the relay set (i.e., $j \in$ Rs − {d} for the social-aware case and $j \in$ Ps − {d} for the social-oblivious case), the strategy of choice is the 1-copy 2-hop for both the social-oblivious and social-aware category. However, the 1-copy 2-hop social-aware scheme is overall more advantageous than its social-oblivious counterpart. More specifically, when the source node is the worst relay for the destination (i.e., $\min_i\{\alpha_{id}\} = \alpha_{sd}$), the social-oblivious and the social-aware approaches are equivalent (given that Ps = Rs). In all other cases, instead, Rs ⊂ Ps, thus, for the set of nodes in Ps − Rs, social-aware forwarding does not impose any constraint, while social-oblivious forwarding needs to impose constraints, thus resulting in stricter conditions for convergence.

Let us now focus on the remaining cases, namely i) when $\alpha_{sd} \leq 2$ and not all intermediate relays have exponent greater than 2, and ii) when $\alpha_{jd} \leq 2$ for all nodes j. In the first case, the social-aware m-copy 2-hop, the social-aware 1-copy n-hop, and the social-oblivious m-copy 2-hop can achieve convergence. In the second case, the only options for convergence are the social- aware m-copy 2-hop and the social-oblivious m-copy 2- hop. We first highlight the differences between the n-hop approach and the 2-hop approach by discussing when the social-aware 1-copy n-hop outperforms the other two strategies in terms of convergence (which can only happen in case i), then we focus on the social-aware and social-oblivious m-copy 2-hop strategies, thus covering both case i and ii.

Assume that there exists at least one node z that meets the destination with $\alpha_{zd} > 2$. The m-copy 2-hop strategies send multiple copies to a set of relays, which in turn can only deliver the message to the destination directly. This implies that intermediate relays must have collectively the capability of reaching the destination, for all subsets with size m of possible relays. Here, only meetings with the destination are relevant, and if all relays but z have very low exponent for encounters with the destination, convergence may not be achieved. Differently from the 2-hop strategies, the social-aware n- hop scheme do not rely exclusively on the capabilities of meeting with d, but it is able to generate a *path* towards the destination in which intermediate nodes may not be good relays for d but good relays towards nodes with high fitness (in the extreme case, only $\alpha_{zd} > 2$ can hold). Thus, in the n-hop case, as long as the message can leave intermediate relays within a finite expected time, this could be enough for convergence. When all three strategies achieve convergence, the one to be preferred can be chosen based on resource consumption considerations. With the m-copy 2-hop strategies there can be up to 2m transmissions, while with the 1-copy n-hop scheme there are n. Hence, when n < 2m, the single-copy scheme should be preferred.

Let us finally compare the social-oblivious and the social-aware m-copy 2-hop schemes. Since they seem to cover similar mobility scenarios (as discussed in the previous section) and to be based on similar

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

mechanisms (the mini and maxi quantities, whose relation with m determines the convergence), it may be difficult to intuitively evaluate which one performs better in terms of convergence. As shown in Appendix A, it may happen that either the social-oblivious m-copy 2-hop scheme achieves convergence when the social-aware m- copy 2-hop scheme does not, or vice versa, depending on the underlying mobility process. An example of the first case is when there are a lot of nodes that meet the source with high $\alpha_{sj}$; if those relays have very low $\alpha_{jd}$, they will not be used by the social-aware scheme, and this may hinder convergence of the second hop. It is easy to construct a corresponding example for the other case.

## 2.2   End-to-end delay in opportunistic networks in presence of duty cycling

A possible roadblock in using opportunistic networking for offloading is the fact that direct communications may consume significant energy. To address this, nodes are typically operated in duty cycling mode, by letting their WiFi (or Bluetooth) interfaces ON only for a fraction of time. The joint effect of duty cycling and mobility is that, even if the network is dense, the resulting patterns in terms of communication opportunities is similar to that of conventional opportunistic networks, as devices are able to directly communicate with each other only when they come in one-hop radio range *and* both interfaces are ON.

The net effect of implementing a duty cycling scheme is thus the fact that some contacts between nodes are missed because the nodes are in power saving mode. Hence, *detected* intercontact times, defined as the time between two consecutive contact events during which a communication can take place for a pair of nodes, are longer than intercontact times determined only by mobility, when a duty-cycling policy is in place. This heavily affects the delay experienced by messages, since the main contribution to message delay is in fact due to the intercontact times. In [9] (presented in D3.2 [2]) we have focused on exponentially distributed intercontact times and we have studied how these are modified by duty cycling, obtaining that intercontact times remain exponentially distributed but their rate is scaled by the inverse of the duty cycle. Building upon this result, we have then investigated how the first moments of the end-to-end delay vary with the duty cycle for a number of opportunistic forwarding schemes. In addition, we have found that energy saving and end-to-end delay both scale linearly with the duty cycling. Therefore, for a single message delivery, the same energy saved through duty cycling is spent because the network must stay alive longer. Thus, the main advantage of duty cycling is enabling the network to carry more messages by being alive longer (rather than improving the energy spent for each single delivery).

Our work in [9] assumed that the value of the duty cycle was given and studied its effects on important performance metrics such as the delay, the network lifetime, and the number of messages successfully delivered to their destination. More in general, the duty cycling can be seen as a parameter that can be configured, typically, based on some target performance metrics. To this aim, in this part of the work we develop a mathematical model that allows us to tune the duty cycle in order to meet a given target performance, expressed as a probabilistic guarantee (denoted as *p*) on the delay experienced by messages. Considering probabilistic, instead of hard, guarantees, allows us to cover a very broad range of application scenarios also beyond best-effort cases – all but those requiring real-time streaming. Specifically, we study the case of exponential, hyper-exponential and hypo-exponential delays (please recall that any distribution falls into one of these three cases, at least approximately), deriving the optimal duty cycle for each of them. For the simple case of exponential delays we are able to provide an exact solution. For the other two cases, we derive an approximated solution and the conditions under which this approximation introduces a small fixed error ε (which is always below 0.14) on the target probability *p*. Specifically, in the worst case, the approximated duty cycle introduces an error on the target probability *p* of about 0.1 (hyper-exponential case) and 0.14 (hypo-exponential case), while in the other cases the error is well below these thresholds.

### 2.2.1 Optimal duty cycling settings

In this section we discuss how to derive the optimal duty cycle $\Delta_{opt}$ such that the delay of a tagged message remains, with a certain probability p, under a target fixed threshold z or, in mathematical notation, $\Delta_{opt} = \min\{\Delta : P\{D_\Delta < z\} \geq p\}$. Since the delay increases with $\Delta$, the latter is equivalent to finding the solution to the following equation:

$$\Delta_{opt} = \{\Delta : P\{D_\Delta < z\} = p\} \tag{1}$$

In order to find the solution to this Equation, the distribution of the delay $D_\Delta$ should be known. To this end, we can exploit the results presented in [9] (see D3.2 [2]) as follows. From [9] we know that when inter-contact times are exponential the detected inter-contact times are also exponential. Therefore, we can use existing models to derive the moments of the end-to-end delay. For example, we can use the model presented in [12] and reported in D3.1 [1]. Based on the first and second moments, we can use well-known distribution approximation techniques for deriving $D_\Delta$. Specifically, if the resulting coefficient of variation is 1, we can approximate $D_\Delta$ with an exponential distribution. When it is greater than 1 with a hyper-exponential distribution, while when it is lower than 1 with a hypo-exponential distribution. Based on these approximations, we can find approximate solutions for the optimal duty cycling by solving the Equation above.

The case when $D_\Delta$ can be approximated with an exponential distribution is very easy. In this case, $D_\Delta$ is distributed exponentially with a rate $\lambda\Delta$ where $\lambda$ is the rate of the delay when no duty cycling is used. Therefore, the optimal duty cycling can be easily obtained as

$$\Delta = -\log(1 - p) / \lambda z$$

For the hyper- and hypo-exponential cases, the solution is not as straightforward, because by substituting the expression of $P\{D_\Delta < z\}$ in Equation 1, the resulting expression cannot be inverted to find $\Delta$. However, it is possible to use approximate expressions, that prove to be within a very reasonable margin of error (below 0.14 in all cases, see Appendix A).

Based on these analytical results, we can therefore set the tradeoff between energy saving and throughput (i.e., end-to-end delay) as needed. Figure 3 illustrates this tradeoff when the end-to-end delay can be approximated with an exponential distribution. Similar results can be obtained in the other cases.



**Figure 3. Example of delay/energy tradeoff.**

Figure 3 shows, for various possible values of $\Delta$, the delay that can be guaranteed (z) with a given probability (p). Clearly (i) for a given $\Delta$, the higher the delay, the higher the probability with which it can be guaranteed and (ii) for a given probability p the delay that can be guaranteed with lower duty cycling is higher than the delay that can be guaranteed with lower duty cycling (i.e., curves move "from right to left when $\Delta$ increases). Even more importantly from a network configuration standpoint, if we want to guarantee a certain delay z with a certain probability p, we can identify the corresponding point (z,p) in the graph. Remember that $\Delta$ is the fraction of time during which nodes are active, and therefore we aim at the minimum possible value of $\Delta$ compatible with the point (z,p). Therefore, the optimal duty cycle is the one corresponding to curve with the minimum $\Delta$ that remains "on the left" of the point (z,p) in the graph.

Graphs like those in Figure 3 (that can be derived through our analytical results) thus provide simple tools for network operators to set the maximum energy saving that can be achieved, given a constraint in terms

of extra capacity that they need to obtain from the opportunistic network, or, complementary, the maximum capacity that they can obtain if they have a constraint in terms of maximum energy consumption that users can tolerate.

# 3 Capacity analysis: Assessing capacity of LTE

To achieve high throughput performance, in addition to an advanced physical layer design LTE exploits a combination of sophisticated radio resource management functionalities, such as Channel Quality Indicator (CQI) reporting, link rate adaptation through Adaptive Modulation and Coding (AMC), and Hybrid Automatic Retransmission Request (HARQ) [1]. Specifically, a base station (eNB) can simultaneously serve multiple users on orthogonal subcarriers that are grouped into frequency resource blocks (RBs). Then, each user (UE) periodically measures channel state information that is fed back to the eNB in terms of CQI reports. Typically, only aggregate CQI values are reported to reduce channel feedback information. CQI measurements are used by eNBs for scheduling and link rate adaptation on the downlink [19]. For instance, the modulation and coding scheme (MCS) is typically selected in order to maximise the data rate to the scheduled UE subject to a constraint on the error probability. How CQI values should be computed by the UE using channel state information (e.g., SINR measurements) is implementation dependent. Unfortunately, past research has shown that it is difficult to derive accurate link performance predictors under realistic channel assumptions. Automatic retransmission protocols with channel coding (HARQ) are also exploited to mitigate errors at the physical layer. More precisely, HARQ procedures use the classical stop-and-wait algorithm, in which the eNB decides to perform a retransmission based on the exchange of ACK/NACK messages with the UE. Then, UEs try to decode the packet by combining the retransmitted copies.

Since user, cell and radio link throughputs are among the most important performance indicators that the operators adopt to asses the QoS in an LTE system [64], an extensive literature exists that investigates LTE throughput performance based on analytical models [23][64][38], simulation tools [19] or field tests [25][17]. However, it is evident that a complex interplay exists among the various mechanisms that operate at the MAC layer to improve communication reliability and to increase data rates. This makes accurate LTE throughput analysis notably difficult. Thus, most studies limit the analysis only to the radio link throughput or consider single MAC functions in isolation [31]. Furthermore, simplified error models are typically considered that only allow deriving upper bounds for the LTE throughput [64].

The contribution of this section is twofold. The first contribution is the development an initial model of the **user-perceived MAC-layer throughput on the downlink channel**. Our model is valid for homogeneous cells [55] and Rayleigh-distributed fading. Our model simultaneously caters for CQI feedback schemes that use spectral efficiency to generate CQI, as well as AMC and HARQ protocols. Furthermore, we include in the analysis an accurate link layer abstraction model that uses the Mean Mutual Information per coded Bit (MMIB) metric to derive the physical error probability [13]. The throughput estimates of our model are accurate, as validated using the ns-3 simulator extended with the LENA module for LTE [73]. As we will discuss in Section 8 this model is a first step towards a complete mathematical characterization of the capacity of a complete offloading solution, providing a detailed model of the cellular network performance. In particular, we are currently working towards a modelling framework that combines the analysis developed in this Deliverable with existing models for data disseminations in opportunistic networks. The second contribution of this section is a proposal for increasing the LTE downlink capacity by enabling **automatic tuning of AMC parameters**. Specifically, we have developed a new AMC scheme that exploits a reinforcement learning (RL) algorithm to adjust at run-time the MCS selection rules based on the knowledge of the effect of previous AMC decisions. The salient features of our proposed solution are: i) the low-dimensional space that the learner has to explore, and ii) the use of direct link throughput measurements to guide the decision process. Simulation results obtained using ns3 demonstrate the robustness of our AMC scheme that is capable of discovering the best MCS even if the CQI feedback provides a poor prediction of the channel performance. Therefore, this result provides a possible capacity enhancing tool for the LTE part of an offloaded network, to be used in addition (and orthogonally to) offloading solutions.

## 3.1 LTE MAC-layer Throughput

### 3.1.1 LTE MAC Model

We now briefly describe relevant details of the LTE downlink, with special attention to frame structure, CQI feedback mechanisms and HARQ protocols. We also introduce the system model and notation, and we discuss the main assumptions that underlay our analysis.

In LTE, each DL frame is 10 ms long and it consists of ten transmission time intervals (TTIs). Furthermore, each TTI consists of two 0.5 ms slots. Each slot contains seven OFDM symbols. In the frequency domain, the system bandwidth, $W$, is divided into several orthogonal subcarriers. Each subcarrier has a bandwidth of 15 kHz. A set of twelve consecutive subcarriers over the duration of one slot is called a physical Resource Block (RB). Let $q$ denote the total number of RBs available over the system bandwidth.

Since the RB bandwidth is only 180 kHz, it is reasonable to assume that the channel response is frequency-flat across all the twelve subcarriers of the RB[1]. Then, let us denote with $\gamma_{i,k}$ the SNR of the $i^{th}$ RB of the $k^{th}$ UE. Clearly, the statistics of the SNR depend on the channel model and the multi-antenna diversity mode of operation. As commonly adopted in other LTE models, e.g. [23], in this study we assume that the fading from the eNB to the UEs is *Rayleigh distributed*. This implies that *the SNR of each RB is an exponential* random variable (RV) [31]. Furthermore, we also assume an *homogeneous cell model* [55], i.e. the SNR is independent for different users and RBs. This also means that the SNRs of all RBs are *uncorrelated* in frequency and space, and $\gamma_{i,k}$ can be regarded as independent and identically distributed (i.i.d.) RV.

Popular methods (e.g., EESM and MIESM) that are typically used in LTE to compute CQI values rely on the concept of *effective* SNR. Basically, the UEs map the SNRs of multiple subcarriers/RBs into a single value by applying complex non-linear transformations. Then, the effective SINR is used to estimate the BLER experienced by a user and to determine the appropriate MCS, i.e. the MCS that allows the UE to decode the transport block with an error rate probability not exceeding 10%. However, the statistics of the effective SNR generated by EESM and MIESM techniques are not known in closed-form. Thus, they must be approximated or computed numerically, which makes performance analysis difficult [56][23]. An alternative approach proposed in [68] implement AMC capabilities is based on the *spectral efficiency*. Specifically, let us denote the with $\eta_{i,k}$ the spectral efficiency of the $i^{th}$ RB of the $k^{th}$ UE. Then, it holds that [53]

$$\eta_{i,k} = \log_2\left(1 + \frac{\gamma_{i,k}}{\Gamma}\right) \qquad (4.1)$$

where $\Gamma = -\ln(5/\beta)/1.5$ and $\beta$ is BLER upper bound. Now a static mapping can be determined between the spectral efficiency and the CQI index, as well as between the CQI index and the MCS value [68]. More formally, let us denote with $C_{i,k}$ the CQI index for the $i^{th}$ RB of the $k^{th}$ UE. Typically the value of CQI can range between 1 and $L$. Then, $C_{i,k} = j$ ( $j = 1,\ldots,L$ ) if $S_j \leq \eta_{i,k} \leq S_{j+1}$, with $S_0 = 0$ and $S_L = \infty$. In other words the CQI value is a quantised version of the spectral efficiency[2]. Furthermore, in the 3GPP-LTE standard the available MCS indexes are 32 but a 4-bit CQI allows selecting only 15 MCS. Thus, in practical LTE systems only a subset of available MCS is typically used.}. Closely related to the MCS selection is also the transport block (TB) size determination. More precisely, let $n_k$ the number of RBs allocated to the $k^{th}$ UE during a frame. Then, the number $B$ of bits that can be delivered in those RBs, which is called transport

---

[1] This assumption will not hold for highly dispersive channels with long delay spread.

[2] Note that in the 3GPP-LTE standard, L=15 and the $S_j$ thresholds are specified in Table 7.2.3-1 of [68].

block, is a function of the MCS index[3]. Furthermore, if $B > Z$ (with $Z = 6144$ bits in 3GPP-LTE) the transport block is *segmented* into a number $C$ of *code blocks* (CBs) that are independently encoded. Note that the CB size highly impacts the actual BLER performance for a given MCS [13]. Figure 4 exemplifies the transport block segmentation.



**Figure 4: Transport block segmentation.**

Regarding the HARQ protocol, LTE employs two types of HARQ schemes. In HARQ type-I, each encoded data frame is retransmitted until the frame passes the CRC test or the maximum number of retransmissions is reached. Erroneous frames are simply discarded. In contrast, in HARQ type-II, each transmission contains incremental redundancy (IR) about the data frame. Thus, consecutive transmissions can be combined at the receiver to improve error correction. Although our model is valid for all HARQ types, in the following we only consider HARQ type-II that is the most widely used in LTE. Note that in LTE systems *retransmissions typically use the same MCS index as the initial transmission*. It is also important to point out that the transmission of HARQ feedbacks (i.e. ACK/NACK messages) is not instantaneous but each received packet experiences a processing delay. According to the LTE standard, the processing delay at the receiver is about 3ms. Thus, assuming the same delay to process data transmissions and ACK/NACK messages, the HARQ round trip time, say $\tau_{ARQ}$, is 7 TTIs, as shown in Figure 5. For this reason, an eNB must support up to 8 parallel HARQ processes for each UE to enable *uninterrupted* communications. In this way, an eNB can continue to transmit new TBs while the UEs are decoding already received TBs.



**Figure 5: HARQ processes and timing in FDD-LTE DL.**

### 3.1.2 Throughput Analysis

In this section we develop the mathematical model of the MAC-layer throughput for the LTE downlink. First of all, let us assume that $n$ UEs are randomly distributed in the cell, and let $d_k$ be the distance of the $k^{th}$ UE from the eNB. As discussed in Section 3.1.1 we develop our analysis in the case $\gamma_{i,k} \sim Exp(\lambda_k)$, where

---

[3] See Table 7.1.7.2.1-1 of [68] for the static mapping between TB size, MCS and number of RBs allocated to the UE.

the rate parameter $\lambda_k$ of the exponential distribution depends on the UE position. Under this assumption the statistics of the spectral efficiency for each RB can be expressed in a closed-form as given by the following Theorem (unless otherwise stated, all proofs are reported in [16]).

*Theorem 4.1:* If $\gamma_{i,k} \sim Exp(\lambda_k)$ then the cumulative distribution function (CDF) of the spectral efficiency $\eta_{i,k}$ in equation~(4.1) is computed as:

$$F_\eta\left(x;i,k\right) = \begin{cases} 1 - e^{-\lambda_k \Gamma\left(2^x - 1\right)} & \text{if} \quad x \geq 0 \\ 0 & \text{if} \quad x < 0 \end{cases} \tag{4.2}$$

LTE specifies different types of CQI reporting: *wideband* and *subband*. Specifically, the wideband CQI represents the SNR observed by the UE over the whole channel bandwidth, while the subband CQI represents the SNR observed by the UE over a collection of adjacent RBs. Note that a vector of CQI values should be transmitted to the eNB when using the latter feedback scheme. Thus, the subband-level feedback scheme ensures a finer reporting granularity but it also generates a higher overhead. In this study, we focus on the wideband feedback scheme and we assume that the CQI reported by the $k^{th}$ UE, say $\hat{C}_k$ is the *arithmetic mean* of the CQI values computed over all RBs[4]. Then, we use the spectral efficiency to generate the CQI values from the SNR measures of all RBs. The statistics of the wideband CQI are mathematically derived below.

*Claim 4.1:* The probability mass function (PMF) of the CQI value for the $i^{th}$ RB assigned to the $k^{th}$ UE is given by

$$g_{i,k}\left[j\right] = F_\eta\left(S_{j+1};i,k\right) - F_\eta\left(S_j;i,k\right) \tag{4.3}$$

*Claim 4.2:* The probability mass function (PMF) of the CQI value for the $i^{th}$ RB assigned to the $k^{th}$ UE is given by

$$g_k\left[j\right] = \sum_{l=qj}^{q(j+1)-1} g_{i,k}^{(q)}\left[j\right] \ , \tag{4.4}$$

where $g_{i,k}^{(q)}\left[j\right]$ is the *q*-fold convolution of $g_{i,k}\left[j\right]$.

As described in Section 3.1.1 a static mapping is typically established between the CQI value received at the eNB and the MCS for the downlink transmissions. For simplicity of notation we indicate with $m(j)$ the MCS that the eNB uses when the wideband CQI reported by a UE is equal to $j$.

Before proceeding with the throughput analysis, we need to introduce the physical layer error model. In this study we adopt the general approach initially proposed in [68] to accurately approximate the BLER curves of OFDMA-based wireless systems, and later specialised for the LTE case in [43]. Specifically, we assume that the *mutual information per coded bit* (MIB) of MCS $m$, as defined in [30], can be accurately approximated by a combination of Bessel functions of the SNR $\gamma$ as follows

$$I_m\left(\gamma\right) \approx \sum_{h=1}^{H} \alpha_h J\left(\psi_h \sqrt{\gamma}\right) \ , \tag{4.5}$$

where $H$, $\alpha_h$ and $\psi_h$ parameters are empirically calibrated as a function of the MCS index. Subsequently, the *mean* MIB (MMIB) value for each UE is computed by averaging the corresponding mutual information

---

[4] Note that an alternative solution would be to report the worst CQI value over all (or a subset of) RBs as in [64].

of all RBs allocated to that UE. Specifically, let $\Omega(k)$ be the set of RBs that are allocated to the $k^{th}$ UE by the scheduler. Then, the MMIB value over the vector of SNR values for each RB assigned to the $k^{th}$ UE when $m$ is the adopted MCS is simply given by

$$\hat{I}_{m,k} = \frac{1}{\omega(k)} \sum_{i \in \Omega(k)} I_m(\gamma_{i,k}) \ , \tag{4.6}$$

where $\omega(k)$ is the cardinality of the $\Omega(k)$. The non-linear nature of (4.5) makes an exact analysis difficult. Thus, previous studies limit the computational complexity of deriving MMIB values in multi-user scenarios by considering a quantised version of the $I_m(\gamma)$ function (4.5) in order to *discretise the MIB metric [43]*. More precisely, let us define a set $V_m = \{\mu_m[0], \mu_m[1], \ldots, \mu_m[v_m]\}$ for each MCS $m$ such that

$$\mu_m[v] = I_m(Q_{m,v}) \ , \tag{4.7}$$

where $(Q_{m,v+1} - Q_{m,v}) = \delta\gamma$ is the quantisation step size, and $Q_{m,0}$ is the minimum usable SNR for MCS $m$. Now, let us denote with $H_{i,m,k}$ the discrete MIB value for the $i^{th}$ RB scheduled to the $k^{th}$ UE when $m$ is the adopted MCS. Similarly to the approach adopted for CQI mapping, we assume that $H_{i,m,k} = \mu_m[v]$ ($v = 0, \ldots, v_m$) if $Q_{m,v} \le \gamma_{i,k} \le Q_{m,v+1}$. In other words the discrete MIB value is associated to a *range* of SNRs. It is straightforward to derive the statistics of the discretised MIB metric as follows.

*Claim 4.3:* The probability mass function (PMF) of the CQI value for the $i^{th}$ RB assigned to the $k^{th}$ UE is given by

$$h_{i,m,k}[v] = \int_{Q_{m,v}}^{Q_{m,v+1}} f_\gamma(x; i, k) dx \ , \tag{4.8}$$

where $h_{i,m,k}[v] = \Pr\{H_{i,m,k} = \mu_m[v]\}$.

Similarly, we introduce a discrete MMIB metric, say $\hat{H}_{m,k}$, computed over the set of RBs allocated to the $k^{th}$ UE when $m$ is the adopted MCS. In particular, $\hat{H}_{m,k}$ can be obtained as the mean of the $H_{i,m,k}$ values over the set $\Omega(k)$. Thus, the statistics of the discretised MMIB value are derived using the same technique of Claim 4.2.

*Claim 4.4:* In an homogeneous cell the PMF of $\hat{H}_{m,k}$ is given by

$$h_{m,k}[v] \approx \sum_{l \in \Phi_v} h_{i,m,k}^{(\varpi(k))}[l] \ , \tag{4.9}$$

where $h_{i,m,k}^{(\varpi(k))}[l]$ is the $\omega(k)$-fold convolution of $h_{i,m,k}[l]$. The definition of the $\Phi_v$ set is quite involved and is given in [16].

Once the MMIB value is given, a direct MMIB to BLER mapping can be used to obtain the *code block error rate*, without necessarily defining an effective SINR. Following the approach proposed in [68], the empirical BLER curve for MCS $m$ can be approximated with a Gaussian cumulative model as follows

$$CBLER_m(y, e) = \frac{1}{2}\left[1 - erf\left(\frac{y - b_e}{c_e}\right)\right] \ , \tag{4.10}$$

where $y$ is the MMIB value, while $b_e$ and $c_e$ are parameters used to fit the Gaussian distribution to the empirical BLER curve[5]. These parameters depend on the Effective Code Rate (ECR), i.e. the ratio between the number of downlink information bits (including CRC bits) and the number of coded bits. Intuitively, the ECR value is a result of the selected TB size, MCS, and $\Omega(k)$. Then, the overall error probability for a transport block transmitted as a combination of $C$ code blocks, each one associated with a MMIB and ECR value, can be computed as

$$TBLER_m(y,e) = 1 - \frac{1}{2}\prod_{i=1}^{C}\left(1 - CBLER_m(y_i,e_i)\right) \quad . \tag{4.11}$$

However equation (4.11) does not take into account the impact of an IR-HARQ mechanism that combines retransmissions to improve error correction. To generalise equation (4.11) for a system with incremental redundancy we adopt the same approach as in [71]. In particular, we introduce an *equivalent* MMIB metric as the average of the mutual information values per HARQ block received on the total number of retransmissions. More precisely, let us assume that the original transport block has been retransmitted $r$ times. Then, let $\left(\hat{I}_{m,k}^{(0)}, \hat{I}_{m,k}^{(1)}, \ldots, \hat{I}_{m,k}^{(r)}\right)$ be the vector of MMIB values for each of these transmissions. The equivalent MMIB for the $r^{th}$ retransmission can be computed as follows

$$\hat{I}_{m,k,r} = \frac{1}{r+1}\sum_{i=0}^{r}I_{m,k}^{(i)} \quad . \tag{4.12}$$

Then, the PMF of the equivalent MMIB value for the $r^{th}$ retransmission is $h_{m,k}^{(r)}[v] = \Pr\left\{\hat{I}_{m,k,r} = \mu_m[v]\right\}$. This PMF can be obtained using the same technique as in Claim 4.4 and it is not reported here for the sake of brevity. Similarly, we compute the effective ECR after $r$ retransmissions, say $e^{(r)}$, by dividing the number of information bit of the original transmission with the sum of the number of coded bits of each retransmission. Finally, by applying the law of total probability the *average* TB error probability at the $r^{th}$ retransmission for the $k^{th}$ UE when $m$ is the adopted MCS can be computed as

$$P_e(m,k,r) = \sum_{v=0}^{v_m}TBLER_m\left(\mu_m[v], e^{(r)}\right) \cdot h_{m,k}^{(r)}[v]dy \quad . \tag{4.13}$$

To conclude the MAC-layer throughput analysis we have to model the operations of the packet scheduler at the eNB, which is responsible for allocating RBs to UEs every TTI. In this study we consider the Round Robin (RR) scheduler that works by dividing the available resources among the UEs in a fair manner. In particular the scheduler allocates a set of consecutive resource blocks, called resource block groups (RBGs), whose size $P$ depends on the system bandwidth [69]. Consequently the number of RB assigned to each UE is simply given by

$$n_k = \max\left\{P, \left\lfloor\frac{q}{n}\right\rfloor\right\} \quad . \tag{4.14}$$

The scheduler also interacts closely with the HARQ protocol. Typically, a non-adaptive HARQ mechanism is implemented in LTE systems, which implies that the scheduler should maintain the same RBG and MCS configuration of the original TB when scheduling the retransmissions.

---

[5] Empirical BLER curves can be obtained through field measurements or detailed link-level simulations.

**Figure 6: RR operations with q=12, P=2 and n=8.**

As discussed in Section 3.1.1, the scheduler can control up to 8 HARQ processes for generating new packets and managing the retransmissions. However, the actual number of HARQ processes that are activated by the scheduler is bounded by the number of times the same UE is scheduled during the HARQ period $\tau_{ARQ}$. In turn, this depends on the total number of available resources during a time window of duration $\tau_{ARQ}$, the RGB size and the number of UEs. To illustrate this dependency in Figure 6 we exemplify the scheduling decisions that are cyclically performed by the RR scheduler during an HARQ period with q=12, P=2 and n=8. As shown in the figure, each UE is scheduled six times. In general, the average number of times each UE is scheduled in $(1+\tau_{ARQ})$ TTIs is simply given by

$$n_{RR} = \frac{q(1+\tau_{ARQ})}{nP} \quad . \tag{4.15}$$

However, not all the transmission opportunities allocated by the eNB to an UE result into a successful transmission. In particular let us denote with $P_s(m,k,r)$ the probability that the $k^{th}$ UE correctly decodes a TB after $r$ retransmissions when $m$ is the adopted MCS. It holds that

$$P_s(m,k,r) = \left[ \prod_{i=0}^{r-1} P_e(m,k,i) \right] \times \left[ 1 - P_e(m,k,r) \right] \quad . \tag{4.16}$$

The above equation is easily derived by observing that the $r^{th}$ retransmission is a success only if the previous $(r-1)$ transmissions were TBs received erroneous and the $r^{th}$ transmission is correctly decoded.

To perform the throughput analysis we observe the system behaviour only during the HARQ period because the HARQ processes that define the occupancy pattern of the channel (i.e., new transmissions and retransmissions) regenerates after each such period. Then, it follows that the MAC-layer throughput for the $k^{th}$ UE is

$$\rho_k = \frac{n_k^{succ} \cdot E[TB]}{1+\tau_{ARQ}} \quad , \tag{4.17}$$

where $n_k^{succ}$ is the average number of HARQ processes for the $k^{th}$ UE which terminate with a success in the HARQ period, while $E[TB]$ is the average number of information bits that are delivered with a successful transmission. To compute $n_k^{succ}$ we can note that in an HARQ period there are at most $n_{RR}$ active HARQ

processes. Since RR equally distributes transmission opportunities to each UE, then all UEs have the same number of active HARQ processes. Note that only a fraction $P_s(k)$ of the $n_{RR}$ HARQ process that are on average active in each HARQ period terminates with a successful transmissions. Hereafter we derive exact expressions for the unknowns in (4.17)

*Theorem 4.2:* By assuming an homogenous cell with Rayleigh-distributed fading, and a RR scheduling policy

$$E[TB] = \sum_{j=0}^{L} TBS(m(j), n_k) g_k[j] \quad , \tag{4.18a}$$

$$n_k^{succ} = n_{RR} \sum_{j=0}^{L} \left[ \sum_{r=0}^{r_{max}} \frac{[P_s(m,k,r)]^2}{[1 - P_e(m,k,r)]} \right] g_k[j] \quad . \tag{4.18b}$$

### 3.1.3 Validation

In this section we show a preliminary validation or our modelling approach using the ns-3 simulator with the LENA module for LTE. The main simulation parameters are reported in Table 3. Specifically, we consider an *Urban Macro* scenario, in which path loss and shadowing are modelled according to the COST231-Hata model [22], which is widely accepted in the 3GPP community. The fading is Rayleigh distributed. To limit the computation complexity of the simulator pre-calculated fading traces are included in the LTE model. Given the downlink system bandwidth (see Table 3) a RBG comprises two RBs [69], i.e., P=2. Regarding the network topology, we considered a single cell with a varying number of static UEs, chosen in the range [10,50]. One *tagged* UE is positioned at a fixed distance from the eNB, while the other UEs are randomly deployed within the cell. The cell radius is 2 Km. Note that, in our settings a maximum number of 96 (i.e., 8q/P) unique UEs can be scheduled within an HARQ period. Indeed, if n>96 the RR period is longer than the HARQ period. All results presented in the following graphs are averaged over multiple simulation runs with different fading traces and topology layouts. Confidence intervals are generally very tight and are not shown in the figures when below 2%. Each simulation run lasts 300 seconds.

**Table 3: Simulation parameters**

| Parameter | Value |
|---|---|
| Carrier frequency (GHz) | 2.14 |
| DB bandwidth (MHz) | 5 |
| *q* | 25 |
| eNB TX Power (dBm) | 43 |
| CQI Processing time (TTI) | 2 |
| CQI transmission delai (TTI) | 4 |
| Antenna scheme | SISO |
| PDCCH & PCFICH (control ch.) | 3 OFDM symbols |
| PDSCH (data ch.) | 11 OFDM symbols |
| *n* | [10,50] |

The accuracy of our modelling approach is validated considering the throughput of the tagged UEs. Specifically, Figure 7 shows the MAC throughput obtained by the tagged UE by varying its distance from the eNB and the number of competing UEs. The shown results have been obtained by setting the maximum

number of transmission equal to one. In other words, transport blocks that are not correctly received are discarded without being retransmitted. The plots clearly indicate that our analysis is very accurate in all the considered settings. At the time of the writing of this deliverable an extensive validation is in progress.



**Figure 7: Comparison of simulation and analytical results versus the distance of the tagged UE.**

## 3.2    Robust AMC in LTE using Reinforcement Learning

Adaptive Modulation and Coding (AMC) in LTE networks is commonly employed to improve system throughput by ensuring more reliable transmissions. Traditional AMC schemes rely on CQI feedbacks that are periodically reported by UEs to their eNBs. AMC schemes typically exploit static mappings between these link quality metrics and the BLER performance of each MCS to select the best MCS (in terms of link throughput). In other words, for each MCS a range of LQM values is associated via a look-up table, over which that MCS maximises link throughput. Either link-level simulations or mathematical models can be used to generate such static BLER curves under a specific channel model. Unfortunately, past research has shown that it is difficult to derive accurate link performance predictors under realistic channel assumptions [26][10][13][44]. Furthermore, a simulation- based approach to derive the mapping between LQM values and BLER performance is not scalable since it is not feasible to exhaustively analyse all possible channel types or several possible sets of parameters [36]. The second main problem with table-based AMC solutions is that a delay of several transmission time intervals (TTIs) may exist between the time when a CQI report is generated and the time when that CQI feedback is used for channel adaptation. This is due to process- ing times but also to the need of increasing reporting frequency to reduce signalling overheads. This mismatch between the current channel state and its CQI representation, known as CQI ageing, can negatively affect the efficiency of AMC decisions [35][5].

To deal with the above issue we have propose a new flexible AMC framework, called RL-AMC, that autonomously and at run-time decides upon the best MCS (in terms of maximum link-layer throughput) based on the knowledge of the outcomes of previous AMC decisions. To this end we exploit reinforcement learning techniques to allow each eNB to update its MCS selection rules taking into account past observations of achieved link-layer throughputs. In this section we outline the general design principles of the proposed solution and we show the main results of the performance evaluation performed in ns3 (complete details are available in Appendix A, which is a reprint of [14]). Overall, our RL-based scheme not only improve the LTE system throughput compared to other schemes that use static mappings between

SINR and MCS, but it is also capable of discovering the best MCS even if the CQI feedback provides a poor prediction of the channel performance.

### 3.2.1 Protocol design

In order to apply the Q-learning approach to the MCS selection problem it is necessary to define: i) the state space of the problem, ii) the feedbacks that the decision agent receives from the LTE network, and iii) the admissible actions for the agent with the action selection strategy. In our RL-based AMC framework, the problem state consists of CQI feedbacks and their evolution trends. The reward is the instantaneous link throughput obtained by a user after each transmission. Finally, an action is the selection of a correction factor to be applied to each CQI feedback to identify the best MCS under the current channel conditions.

Intuitively, a straightforward approach to define the state of the MCS selection problem would be to use the SINR values of received segments of data as state variables, as in [17]. However, the SINR is a continuous variable and it should be discretised to be compatible with a discrete MDP formulation. The main drawback is that a fine discretisation leads to a large-dimensional state space, which increases convergence and exploration times. To avoid this problem, we directly use CQI-based metrics for the state representation. Specifically, we adopt a two-dimensional space to characterise the LTE communication channel. The first state variable represents the CQI value (called $CQI^m$) that the UE should select using the internal look-up table that associates BLER and MCS and received SINR. The second state variable represents the $\Delta CQI^m$ value, which is defined as the difference between the last two consecutive $CQI^m$ estimates. In other words, $\Delta CQI^m$ provides a rough indication of the trend in channel quality evolution. For instance, $\Delta CQI^m < 0$ implies that the channel quality is temporarily degrading.

Since the objective of the MCS selection procedure should be to maximise the link throughput it is a natural choice to define the reward function as the instantaneous link-layer throughput achieved when taking action $a$ in a given state. Thus, a key aspect in the design of the Q-learning algorithm is represented by the set $A$ of admissible actions. In our learning model we assume that an action consists of applying a correction factor to the CQI value that is initially estimated by means of the internal look-up table. As discussed above, the mapping relationship between SINR values and MCS may be inaccurate and the correction factor allows the agent to identify the best modulation and coding scheme (in the sense of maximising the link throughput) for the given channel conditions. For instance, it may happen that the SINR-to-MCS mapping is too conservative for the current channel conditions and an MCS with a higher data rate can be used without violating the target BLER requirement. In this case the correction factor should be positive. Furthermore, a correction factor is also needed to compensate eventual errors due to CQI feedback delay. More formally, we assume that an action taken by the AMC decision agent at time t is one possible choice of an integer number in the set $(-k,\ldots,-2,-1,0,1,2,\ldots,k)$, that we denote as $a_t$ in the following. This index is added to the original $CQI^m$ value to compute the CQI to be sent to the eNB. We argue that $\Delta CQI^m < 0$ we should prefer conservative MCS selections (and thus use values of at lower than 0) because the channel trend is negative, while if $\Delta CQI^m \geq 0$ we can try to use MCSs offering higher data rates (and thus positive values for at). Thus, the set of admissible actions is different whether the channel-quality trend is negative or non-negative. Before proceeding it is useful to point out that the choice of the $k$ value determines how aggressively we want to explore the problem state space. In general, the selection of the $k$ value could take into account the CQI difference statistics, i.e., to what extent a current CQI may be different from the reported CQI after a feedback delay [44]. Finally, a very important learning procedure is the action selection rule, i.e., the policy used to decide which specific action to select in the set of admissible actions. There is a clear trade-off between exploitation (i.e., to select the action with the highest Q-value for the current channel state) and exploration (i.e., to select an action randomly). In our solution we adopt a *softmax action-selection rule* [62] that assigns a probability to each action by applying a Boltzmann-like function to the current Q-value for that action (see Appendix A for more details).

## 3.2.2 Performance evaluation

The simulation setup is the same as the one already described in Section 3.1.3. Regarding the simulation scenarios, we consider two cases. In the first one, ten UEs are randomly deployed in the cell and they are static. Then an additional tagged user is moving with pedestrian speed from the centre of the cell to its boundaries. However, independently of the UE position the *CQI feedback is constant*. Then, Figure 8 shows a comparison of the throughput achieved by the tagged user with and without reinforcement learning. This is obviously a limiting case which is analysed to assess the robustness of our RL-AMC scheme even when CQI provides a very poor prediction of channel performance. As expected with fixed MCS the user throughput is constant when the MCS is over provisioned, while it rapidly goes to zero after a critical distance. On the contrary, our RL-AMC is able to discover the correction factor that should be applied to the initial CQI to force the selection of a more efficient MCS. In addition, the throughput performance of RL-AMC is almost independent of the initial CQI value. Note that in this case RL-AMC must explore the full range of CQI values and we set the k parameter for action selection equal to 15.



**Figure 8: Average throughput as a function of the distance of the tagged user from the eNB in a pedestrian scenario.**

In the second scenario, each UE implements the SINR to CQI mapping described in [69][68], called SE-AMC because it uses spectral efficiency to estimate the transport block error rates. Then, Figure 9 shows a comparison of the throughput achieved by the tagged user with both SE-AMC and RL-AMC schemes at different distances of the tagged UE from the eNB. We can observe that the MCS selection in SE-AMC is too conservative and this results in a throughput loss. On the contrary, RL-AMC method is able to discover the MCS configuration that can ensure a more efficient use of the available channel resources. This is more evident at intermediate distances from the eNB when short-term fading may lead to use more frequently low-rate MCSs. As shown in the figure, the throughput improvement varies between 20% and 55% in the range of distances between 200 meters and 800 meters.

**Figure 9: Average throughput as a function of the distance of the tagged user from the eNB in a pedestrian scenario**

Finally, we also consider a more dynamic environment in which there is an increasing number of UEs in the cell, and all the UEs are moving according to a random waypoint mobility (RWM) model with speed 30 km/h and pause time equal to 5 seconds. Figure 10 shows a comparison of the aggregate cell throughput with both SE-AMC and RL-AMC schemes as a function of the network congestion (i.e., number of UEs). The results clearly indicate that the throughput improvement provided by RL- AMC is almost independent of the number of UEs and it is about 10%. We can also observe that the cell capacity initially increases when going from 10 to 20 UEs. This is due to two main reasons. First, RR is able to allocate RBs in a more efficient way when the number of UEs is higher. Second, the higher the number of UEs and the higher the probability that one of the UEs is close to the eNB and it can use high data-rate MCSs.



**Figure 10: Average cell throughput as a function of the number of UEs in an urban vehicular scenario.**

# 4 Capacity analysis: Assessing capacity of an integrated offloaded network

In this section we present the status of our work on the capacity gains that can be obtained in an integrated offloaded network. Also in this case, we build upon and extend the initial results presented in D3.1 [1], in particular those related to the Push&Track and Droid offloading systems, presented in Section 3 of D3.1.

As it was the case of solutions presented in D3.1, in WP3 we are investigating in parallel several mechanisms and configurations for offloading. The rationale, as also described in the DoW of the project, is to design and quickly test such solutions in dedicated settings, in order to initially compare them, understand their feasibility and indicative performance. This part of the work in this WP is, therefore, preliminary to the more complex work of integration and comprehensive evaluation that is being carried out in WP5, and that will be fully developed during the third year of the project. Through these investigations we are able to already understand the most useful offloading mechanisms and solutions, that are then integrated in selected reference use cases, demonstrated and validated in WP5.

In Section 4.1 we start from the basic Droid system presented in D3.1 (and published in [50]), and extend its basic operations by considering the opportunity of using also a WiFi Access Point infrastructure in a city environment. This is one of the first studies – to the best of our knowledge – that takes into consideration all the three main content delivery enabling technologies that are nowadays considered for offloading, i.e. cellular networks, WiFi networks and opportunistic networks. We compare "vanilla" Droid (which uses only cellular+opportunistic) with content delivery solutions based only on WiFi APs, and on WiFi APs and opportunistic dissemination. Note that we consider a real WiFi AP development, i.e. the one currently available and open to the general users (managed by the municipality) in the city of Bologna. This analysis shows that (i) augmenting WiFi-based delivery with opportunistic networks is very important, as opportunistic delivery allows us to complement the limited pervasiveness of the WiFi AP deployment. However, comparing WiFi+opportunistic with Droid, we show that in case of stringent delivery deadlines, Droid possibility to exploit a more pervasive cellular network (as opposed to the WiFi network) to disseminate the content becomes a winning factor. Moreover, we also test the case of a fully combined solution based on Droid, and also exploiting WiFi APs, and therefore characterise the additional gain, in terms of offloading, of using it. Finally, we investigate the energy-capacity trade-off in this case. Complementary to the approach presented in Section 2.2, in this case we assume that the constraint in terms of maximum energy consumption is represented by a maximum number of copies of the message that each node can disseminate in the opportunistic network (i.e., a form of limited epidemic diffusion). Again, we show that, for a giv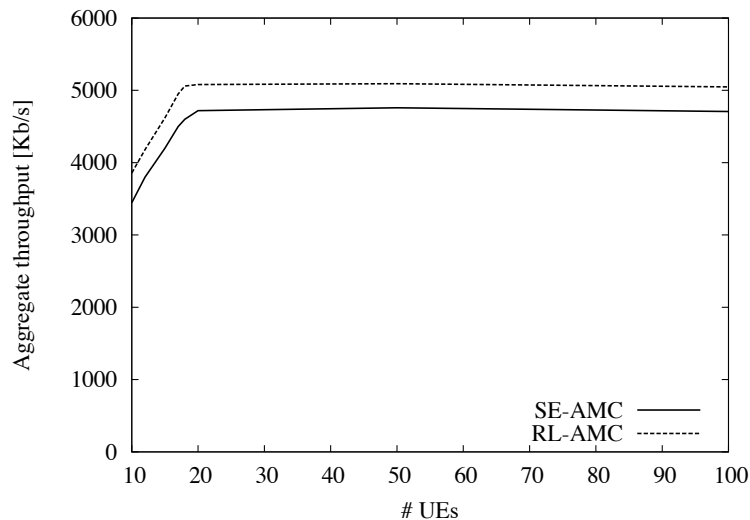en such constraint, using an offloading solution integrating cellular and opportunistic network is winning, with respect to one that uses WiFi APs and opportunistic dissemination.

In Section 4.1 we focus on a typical publish/subscribe scenario, i.e. one where content becomes available at some point in time, and is automatically (and implicitly) requested by all interested nodes at that time. In Section 4.2, instead, we complement these results by considering the case where content is not requested by all mobile nodes in a synchronised way. Specifically, we again consider a vehicular scenario, and assume that content of interest is geo-localised, and becomes interested for users once they enter a specific geographical area. Therefore, requests for content are generated asynchronously from each other, and schemes such as Droid should be modified. We therefore define simple algorithms to support this case, compatible with the overall MOTO architecture described in D2.2.1 [3]. Then, we assess the feasibility of offloading, by showing how much capacity operators can gain if opportunistic dissemination is used also in this case. We show that, clearly, this depends on the features of the mobility patterns, on the amount of time during which users keep the content locally (after having received it), and by how many users are interested in it. Overall, the additional capacity that can be gained ranges between 20% and 90%. These results are presented in more detail in [15], also included as Appendix A.

We refer the reader to Section 8 of this document for a discussion on how we are completing these activities, according to the overall logical structure of the WP activities described in Section 1.

## 4.1   Offloading with both opportunistic and Wi-Fi networks

In this section, we evaluate through simulation the performance of several offloading strategies. We compare opportunistic and AP-based offloading strategies under tight delays. We consider a location-based service in a vehicular context.  Content with traffic information or some infotainment announcement must be distributed to a multitude of users within a given maximum reception delay (in order to guarantee a minimal QoS on a per-content basis). We assume that nodes are equipped with several wireless interfaces, so that they are able to communicate through multiple interfaces simultaneously. Possible combinations involve 3G and 4G to communicate with the cellular infrastructure, Bluetooth or Wi-Fi ad hoc to communicate with neighboring devices, and Wi-Fi in infrastructure mode to communicate with fixed APs.



**Figure 11: Offloading model: The dissemination process is kick started through cellular and/or AP transfers. Content is diffused among vehicles through subsequent opportunistic contacts. Upon reception, users acknowledge the offloading agent using the feedback cellular channel. The coordinator may decide at any time to re-inject copies through the cellular channel to boost the propagation. 100% delivery ratio is reached through fallback re-injections.**

### 4.1.1 Mobility trace and AP position

We employ a large-scale vehicular mobility trace representing the city of Bologna (Italy). The Bologna dataset consists of 10,333 nodes, covering a total of 20.6 km^2 and 191 km of roads. The simulated traffic in the dataset mimics the everyday road activity in the two metropolitan areas. From the mobility trace, we derive a contact trace that features contacts between nodes when the distance between them is below a given threshold (we consider a range of 100 meters, in line with IEEE 802.11p specifications). The resulting trace has a duration of about one hour; on average, 3500 nodes are present at the same time (because some nodes leave while others join during the observation period). The advantage of using this large-scale trace is that, differently from other available datasets, we have a clear high turnover rate, due to vehicles entering and exiting the interest area, and no apparent social links between nodes. The distribution of contact durations is exponential. Most contacts are very short, confirming the highly dynamic nature of the trace. In addition, only few contacts last for more than a few minutes.

We extracted the location of an existing Wi-Fi Hotspot public deployment from http://www.comune.bologna.it/wireless. Figure 12 shows the position of 93 APs along with a map of the town center. We merged the location of APs with the vehicular mobility trace to extract a completely new dataset that includes vehicles mobility and AP positions. From this trace, we derived the connectivity traces between vehicles and fixed APs. The outcome is a completely new time-variant graph with unidirectional links connecting vehicles with the APs and the other vehicles.

**Figure 12: Bologna map with fixed AP positions.**

To characterize this new dataset, we study the pairwise interactions among mobile nodes (vehicles) and fixed APs only. Figure 13 presents the distributions of contact and inter-contact times between vehicles and APs. Contact and inter-contact times follow a lognormal distribution.



**Figure 13: CCDF of contact and inter-contact times with fixed Aps.**

Two anomalies make the study of the dataset interesting. First, a relevant number (around 20%) of inter-contact times are zero s, meaning that when a vehicle exits from the coverage area of an AP, it is already under the range of another. We may explain this aspect by noting that in the city center several APs are very close together. Nevertheless, the sole contact distribution does not tell us the whole story, as very few APs are located in southern and western town districts, and many vehicles passing there enter and exit the system without falling into the coverage zone of any APs.

Moreover, we note that many AP meetings occur in bursts. We infer a strong correlation between the geographic position and the expected duration of contact and inter-contact times with APs. In addition, we remark that around 80% of the contacts with APs last for more than 10 seconds. While this may be an acceptable duration for data transfers, short-lived contacts lasting less than that value could suffer from the duration of authentication and address granting procedures with an AP.

**Figure 14: Time evolution of the total number of contact between vehicles and APs for different message lifetime. Note that, in average, 3500 users are active at the same time.**

Finally, Figure 14 depicts the evolution of the number of vehicles in contact with at least an AP during a given time window (equivalent to the delay-tolerance of content reception in this case). The figure indicates the amount of vehicles that enter in the transmission range of least one AP during the considered distribution period. Augmenting the delay tolerance, the chances that a vehicle enters in the range of at least one AP increase. Still, the number of vehicles benefiting from this transfer opportunity is limited, if measured against their total number in the system, lying always between 5 and 25% of present users.

## 4.1.2 Simulation setup and scenario

In our implementation, we consider a simple contact-based ad hoc MAC model, where a node may transmit only to a single neighbor at a time. Transmission times are deterministic since we do not take into account complex phenomena that occur in the wireless channel such as fading. Communications consist of two different classes of messages (content and control). All transfers, including *ack* messages, may fail due to nodes moving out of each other's transmission range or exiting the simulation area. In addition, it is possible for the same message to be concurrently received through the two interfaces. In that case, we consider the one processed first. The ad hoc routing protocol employed by nodes to disseminate the content is the epidemic forwarding.

Parameters in simulation are set to mimic the functioning of communication technologies currently available to consumers. In each simulation run, the downlink bit-rate for the infrastructure network is set to 100 KB/s, while uplink is fixed at 10 KB/s. These values are in line with the average bit-rate experienced by users of a typical 3.5G network. The bit-rate for the ad hoc link is set to 1 MB/s, also in line with the advertised bit-rate of the IEEE 802.11p standard. The size of each content update is set at 100 KB. The size of the acknowledgement messages is 256 bytes, as it carries very little information (content and node identifiers).

## 4.1.3 Opportunistic or AP-based offloading?

How opportunistic offloading behaves when compared to most traditional AP-based offloading strategies?

To answer this question, we run network level simulations to benchmark the opportunistic offloading algorithm DROiD [50], which is at the base of MOTO, against other more conventional strategies based on direct offload from Wi-Fi hot spots. AP-based offloading takes advantage of the presence of fixed infrastructure that can serve to offload the cellular network. Nevertheless relying only on fixed deployment typically lacks of the flexibility of the pervasive cellular networks, since transmission range and spatial density are limited for physical reasons.

We exploit the Bologna dataset, considering also the AP deployment, with simulation parameters previously described in Section 4.1.2.

**Figure 15: Bologna trace - offloading efficiency comparison between *DROiD*, *DROiD + AP*, *AP-based,* and *AP + opportunistic* distribution strategies. 95% confidence interval plotted.**

In Figure 15 we may appreciate the efficiency of three alternative approaches to data offloading making use of fixed APs, simulated on the Bologna dataset and compared with the DROiD strategy as benchmark. We note that the **AP only** strategy gives extremely poor results. Even with larger tolerance to delays, this strategy is never capable of saving more than *20%* of traffic. Thus, this strategy turns out to be unable to relieve substantially a large fraction of the cellular load. As already hinted during the dataset analysis phase, the main problem in this case is that APs are not ubiquitously available in each town district. Vehicles traveling in areas without Wi-Fi coverage cannot download data from nearby APs, and will likely reach the panic zone without the content. An increase in delay-tolerance of the content only partially mitigates the issue. The analysis, carried out employing the real-world deployment in the city of Bologna, suggests that in order to offload a substantial part of traffic through fixed hot spots their deployment should be carefully planned, without black holes in spatial coverage.

As a second option, we consider that our nodes be able to communicate with both APs and other vehicles through direct ad hoc links, without considering any cellular re-injection (the **AP + Opportunistic** strategy). This proves to be very beneficial for the overall offloading performance, increasing efficiency up to 50% with respect to AP *only*. The possibility to exchange data directly among users, together with mobility, allows spreading the infection in many areas not covered by fixed APs through store-carry-forward routing. Gains, as expected, rapidly improve as the reception delay increases. Fixed APs, in this case, act as fixed (and free) infection source. Still, the benefits of re-injections through the cellular link emerge for shorter reception delays (up to 120s). For shorter deadlines, DROiD results always preferable than AP-based strategies. Once again, re-injections prove essential in the case of lagging infection evolution, and this is particularly true when opportunistic contacts among vehicles are scarce. For short distribution intervals, infected nodes hardly can carry the content far from AP location before its expiration. For longer delay tolerance instead, the continuous use of the APs to infect neighbor nodes gives the *AP + opportunistic* strategy an edge.

In order to take advantage of the fixed hot-spot infrastructure along with the feedback-based re-injection through derivative strategy, we evaluate **DROiD + APs**. In this case, the APs pre-fetch the content at $t_0$. Cellular re-injections intervene only when the diffusion lags, so to overcome the difficulties encountered by users located away from APs range. This strategy emerges to be always the best, guaranteeing more than 65% of offloaded content. APs guarantee a steady infection rate to vehicles passing in their transmission range, letting the cellular infrastructure to target users in more isolated areas.

Further analyses of the traffic flowing on each interface reveal interesting and unexpected information. It turns out that in joint opportunistic/AP-based offloading strategies (**AP + opportunistic** and **DROiD + APs**), the fixed APs tend to kick-start the dissemination, which is then carried over with subsequent direct communications between mobile nodes. The aggregate amount of data flowing through APs in these cases

is roughly 10 times less than when hot spots are the only offloading options. Remarkably, this small fraction of data transferred through APs results very important to bootstrap the dissemination, offering an advantage compared to the non-AP based solutions.

### 4.1.4 Energy savings and fairness

A critical challenge to make mobile data offloading potentially attractive to end-users is to attenuate the impact of opportunistic communications on the battery of devices, concurring thus to increase their lifetime. For this reason, we analyze the impact that simple energy-saving methods have on offloading performance.

In our analysis, we compare **DROiD** [50] with the **AP + opportunistic** strategy, fixing the maximum number of possible opportunistic transmission that a node can do for each message. To put energy saving strategies in practice, we offer to users only a fixed amount of *tokens* for each content to be distributed. The local token count is decreased each time the content is forwarded by a node. When the token count is equal to zero, the node stops forwarding, and waits for the next content to appear.



**Figure 16: Offloading efficiency as a function of the number of transmission tokens for DROiD and AP + Opportunistic. Confidence intervals omitted for clarity.**

Simulation results, presented in Figure 16, show that a limited number of tokens affects offloading performance, wasting possible contacts and lowering the aggregate capacity.

Energy saving strategies have a more pronounced effect on the AP + *opportunistic* schema, which sees its performance highly lowered. The performance gap stretches as the delay-tolerance increases because nodes are more likely to run out of tokens. From the figure, we may appreciate that restricting the number of tokens to 20 does not bring a substantial performance hit for *DROiD*, while its influence is more pronounced in *AP + opportunistic.*

The energy saving scheme should trade off offloading efficiency for battery life, while ensuring at the same time to split the overall energy cost equally between all the nodes involved in the dissemination. However, in opportunistic networks, contacts are typically imbalanced between nodes, and it is common to find nodes that during the message lifetime sustain an important number of data forwarding.

|      | 3    | 5    | 10   | 20   | ∞    |
|------|------|------|------|------|------|
| 30s  | 0.22 | 0.21 | 0.19 | 0.15 | 0.15 |
| 60s  | 0.27 | 0.25 | 0.23 | 0.2  | 0.17 |
| 90s  | 0.29 | 0.27 | 0.27 | 0.23 | 0.18 |
| 120s | 0.3  | 0.33 | 0.3  | 0.25 | 0.21 |
| 150s | 0.35 | 0.33 | 0.32 | 0.27 | 0.2  |
| 180s | 0.39 | 0.37 | 0.34 | 0.29 | 0.21 |

(a) AP + Opportunistic

|      | 3    | 5    | 10   | 20   | ∞    |
|------|------|------|------|------|------|
| 30s  | 0.43 | 0.4  | 0.34 | 0.29 | 0.28 |
| 60s  | 0.44 | 0.41 | 0.35 | 0.3  | 0.28 |
| 90s  | 0.46 | 0.42 | 0.36 | 0.31 | 0.29 |
| 120s | 0.46 | 0.42 | 0.36 | 0.31 | 0.28 |
| 150s | 0.47 | 0.43 | 0.37 | 0.31 | 0.28 |
| 180s | 0.47 | 0.43 | 0.38 | 0.31 | 0.27 |

(b) DROiD

**Figure 17: Jain's fairness index for different combinations of number of tokens and delay tolerance.**

To evaluate **fairness**, we use the Jain's index to compare the fairness in the number of opportunistic transmission for the two schemes under evaluation, with different token values and content reception delay. Figure 17 shows that DROiD always presents better fairness indexes than AP-based strategy. The analysis of nodes forwarding reveals that in the *AP + opportunistic* strategy the number of users participating in content forwarding is sensibly lower than in *DROiD*. This depends on a mix of two factors: the efficiency is lower in general, so a greater number of nodes do not physically store the content to forward. The effect relies also on the fact that AP-based strategies tend to concentrate data forwarders among those nodes that receive the content first. This is the reason why in the 180 s scenario with infinite tokens, DROiD shows better fairness but lower efficiency.

The use of fixed Wi-Fi APs as the only data offloading strategy proves ineffective in the case of medium density AP deployment. Coupling it with opportunistic distribution could improves instead data offloading, namely to kick start the distribution process and to deliver free copies of the content to users located inside their coverage range. Nevertheless, if we consider tight delivery times, the use of the pervasive cellular infrastructure is still required to target isolated node. Offloading strategies relying on random re-injections result intrinsically more fair. This can improve the battery duration of mobile users with respect to AP-based strategies, which always target nodes in their spatial proximity.

## 4.2 Offloading with non-synchronised content requests

Most of the literature on opportunistic-based offloading investigates the scenario where a specific piece of content is generated, and the set of users to whom it has to be delivered is known already at that time and does not change subsequently, i.e. requests are *synchronised*. While significant, this scenario only partially captures relevant use cases. In particular, it does not cover cases where content demand is *dynamic*, i.e. users' requests for the same piece of content can arrive at different time instants. In the latter scenario offloading can still be applied: upon a request, content can reach the requesting user either through the opportunistic network, exploiting an ongoing dissemination process, or through the cellular network, in case the opportunistic dissemination does not reach the user in time. Offloading may even be more needed in case of dynamic requests, as synchronised requests could in principle be served also through multicast transmissions (although [51] – also presented in Section 5 of this document – shows that offloading is beneficial also when multicast is applied).

We have started investigating dynamic content requests, with a particular focus on vehicular scenarios. We deliberately use a very simple offloading scheme, whereby resources provided by mobile nodes are minimally used. Nodes interested in a content store it for a limited amount of time after receiving it. New requests from other users are satisfied either when the requesting user encounters another user storing a copy of the content, or through the cellular network upon expiration of the delivery deadline.

As opposed to most of the literature looking at offloading through opportunistic networks, in our scheme we do not use any epidemic dissemination mechanism. On the one hand, this allows us to test a minimally invasive offloading scheme from the mobile users' perspective. As additional resources spent by mobile devices are sometimes considered a possible roadblock for offloading, our results show the offloading

efficiency when this additional burden is extremely low. On the other hand, this simple scheme allows us to stress the efficiency of offloading in a particularly unfavorable configuration, thus providing a worst-case analysis, all other conditions being equal.

We focus on two complementary scenarios. In the first one, users move in a given physical area, and *all* request a piece of content, though at different points in time. This scenario is representative of users moving inside a limited area, and accessing very popular content, though not particularly time critical (i.e., content that does not generate a surge of requests immediately when it is generated). In the second scenario, users enter and exit (after a short amount of time) a given geographical area, and request content after a random amount of time after they entered the area. This complementary scenario is thus representative of users traversing a geographical area, as opposed to roaming there. Finally, in this scenario we also consider the case where content is requested only with a certain probability, i.e., when content has different levels of popularity.

We analyse the offloading efficiency in these scenarios, defined as the fraction of nodes receiving content through the opportunistic network. We characterise efficiency as a function of key parameters such as the number of users, the deadline of content requests, the time after which users drop the content after having received it, the popularity of the content. Even with an unfavourable opportunistic dissemination scheme, we find that offloading can be very efficient, as it is possible to offload up to more than 90% of the traffic. In other configurations, we find that the considered offloading scheme is less efficient, resulting in an offloading of only about 20%. In such cases, however, there is ample room for improvement, by further leveraging opportunistic networking resources, e.g., through more aggressive content replication schemes.

### 4.2.1 Offloading algorithms for non-synchronised requests

The algorithms we have defined are compliant with the general MOTO architecture presented in D2.2.1 [3]. We assume the existence of a Central Dissemination Manager (CDM), that can communicate with all nodes through the cellular network and keeps track of the dissemination process. The offloading mechanism is defined by the actions taken by requesting nodes and by the CDM, as described by Algorithms 1 and 2, respectively, shown in Figure 18.

---

**Algorithm 1** Actions taken by requesting nodes

▷ Run by a tagged node $k$

1: **Upon** request for content $C$
2: content_received = **false**
3: **Send** content_request to CDM
4: **if** $C$ not received immediately from CDM **then**

▷ try with opportunistic contacts

5:     **while** content_timeout is not over **do**
6:         request $C$ to encountered nodes
7:         **if** content received **then**
8:             content_received = **true**
9:             **Send** ACK to CDM
10:            **break**
11:        **end if**
12:    **end while**
13:    **if** content_received == **false then**
14:        **Receive** $C$ from CDM
15:        content_received = **true**
16:    **end if**
17: **end if**
18: **while** sharing_timeout is not over **do**

▷ available for opportunistic sharing

19:    **Send** $C$ to encountered nodes upon request
20: **end while**
21: **Cancel** content $C$

---

**Algorithm 2** Actions taken by CDM

▷ Run by the CDM for content $C$

**Init** #nodes_with_$C$ = 0

1: **Upon** request from node $k$
2: $k$_served = **false**
3: **if** #nodes_with_$C$ == 0 **then**
4:     **Send** $C$ to $k$
5:     #nodes_with_$C$++
6:     **Set** sharing_timeout for node $k$
7: **else**
8:     **while** content_timeout is not over **do**
9:         **if** ACK received by $k$ **then**
10:            #nodes_with_$C$++
11:            $k$_served = **true**
12:            **Set** sharing_timeout for node $k$
13:            **break**
14:        **end if**
15:    **end while**
16:    **if** $k$_served = **false then**
17:        Send $C$ to $k$
18:        #nodes_with_$C$++
19:        **Set** sharing_timeout for node $k$
20:    **end if**
21: **end if**

22: **Upon** sharing_timeout for node $k$ over
23: #nodes_with_$C$ = #nodes_with_$C$-1

---

**Figure 18. Algorithms for offloading in non-synchronised requests.**

Let us focus first on the actions taken by requesting nodes (Algorithm 1). When a request is generated at a node, the node sends it to the CDM via the cellular network (line 3). The node is guaranteed to receive the content within a given *content timeout*. During the timeout, the node tries to get the content from encountered nodes (lines 5-12). If the timeout expires, it receives it directly from the CDM (lines 13-16). Upon receiving the content, the node sends an ACK to the CDM (line 9 and, implicitly, line 14). In addition, it keeps the content for a *sharing timeout*, during which it can share the content with other encountered nodes (lines 18-20). After the expiration of the *sharing timeout* the content is deleted from the local cache. Note that requests and ACKs are supposed to be much shorter than the content size, and thus do not significantly load the cellular network.

Let us now focus on the actions taken by the CDM (Algorithm 2). Thanks to requests and ACKs, the CDM is always aware of the status of content availability in the network. Upon receiving a request, it checks whether some other node is already storing a copy of the content or not. In the latter case (lines 4-6) there is no chance that the user can get the content opportunistically through another node, and the CDM sends the content directly through the cellular network. In the former case (lines 7-21), it waits to receive an ACK during the *content timeout* (lines 8-15), indicating that the node has received the content. If this does not happen, it sends the content directly to the node (lines 16-20). Finally, upon expiration of the *sharing timeout* for a given node the CDM updates the view on the number of nodes with the content (lines 22-23).

### 4.2.2 Evaluation when all users request the content

In the first scenario we have considered, all users entering the physical area covered by the cell request the content, though they clearly do it at different points in time.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

(a) $N = 20, content\ timeout = 60s,$

**Figure 19. Example of results in the first scenario.**

Figure 19 shows an example of the results we have obtained in this scenario. In the figure, λ represents the rate at which nodes generate the request for the content (i.e., two requests from two different users are spaced by an exponential interval with average 1/ λ). We observe two regimes. When the *sharing timeout* is low, higher request rates result in higher offloading. This is intuitive, because higher request rates results in requests being more concentrated in time. When nodes share the content only for very short amounts of time (see for example the case of 5s), concentrating the requests in time increases the probability of encountering other nodes sharing the content. Less intuitive is the behaviour for large sharing timeouts, where higher request rates results in *lower* offloading efficiency. Intuitively, when requests are more concentrated in time, *content timeouts* for nodes that do not get the content via the opportunistic network are also more concentrated. When a timeout expires and content is delivered via the cellular network, this kicks off a fast increase in the dissemination of content via the opportunistic network in the region of the node whose *content timeout* has expired. When expirations are less concentrated in time (i.e., when request rates are lower), the opportunistic diffusion process has more time to spread content, and therefore the offloading efficiency increases. Additional aspects are analysed in [15] (see Appendix A).

### 4.2.3 Evaluation when all users request the content

Figure 20 shows the offloading efficiency in the second scenario we have considered. In this case, when entering the area covered by the dissemination system, vehicles become interested in the content with a given probability p. If they are interested, they generate a request after a time interval uniformly distributed between the time when they enter and the time when they reach the centre of the area.



**Figure 20. Example of results in the second scenario.**

Figure 20 shows the offloading efficiency for two considered densities of nodes and the different content popularities (p). Results basically confirm previous observations. This is nevertheless important, as this scenario is more representative of a "steady state" behaviour of the offloading system, as nodes constantly enter and exit the area at a given rate, and continuously generate requests (with a given probability). Denser networks (N = 40) achieve higher offloading efficiency. The effect of the popularity parameter is similar to that of the request rate in the first: the higher the popularity, the higher the number of nodes sharing content, the higher the offloading efficiency. It is interesting to note, however, that, due to the mobility of the nodes, they stay within the area only for about 30s in total, and, on average, stay in the area for about 22s after having generated a request. This is the "useful time window" during which they can receive content via opportunistic dissemination. Even though this time window is rather short, offloading is very efficient, even at quite low popularities (p = 0.2). Additional aspects are analysed in [15] (see Appendix A).

# 5 Intra-technology scheduling: Joint use of multicast and D2D in cellular networks

This is the first section where we present results on scheduling. As discussed in Section 1, we are exploring multiple directions from this standpoint. In this section, we start addressing one of the fundamental questions for the practical applicability of offloading, i.e. whether using cellular multicast would not be used in most of the cases. This is particularly relevant for multimedia (non real-time) popular content, that may be required by multiple users simultaneously in the same physical area (i.e., covered by the same cell). In fact, among general multimedia services, some involve delivering the same piece of data to a community of interested users. Examples that fit this use case are software updates, on-demand videos, and road traffic information. When a multitude of co-located users is interested in the same content, two possible approaches could help operators to relieve their cellular infrastructures: **mobile data offloading** and **multicast**.

**Multicast** makes use of a single unidirectional link, shared among several users inside the radio cell, allowing, in principle, a more efficient use of network resources with respect to the case where each user is reached through dedicated bearers.[6] To ensure the coexistence between multicast and unicast services, operators must reserve a fixed amount of resources for multicast transmissions. Despite its attractive features, multicast presents intrinsic and still unresolved issues that limit its exploitation due to the difficult adaptation to radio channel conditions.

**Mobile data offloading** is an alternative low cost solution to reduce the burden on the infrastructure network. Direct device-to-device (D2D) communications help lowering the load on the infrastructure. The increase in the density of mobile users gives rise to an abundance of contact opportunities and represents a strong argument to support opportunistic offloading strategies. In order to encourage subscribers to offer their battery and storage resources to this end, mobile providers may offer monetary incentives and pricing discounts. As a counterpart, users should accept a delayed content reception.

In the present section*, we explore the combination of opportunistic traffic offloading with multicasting*. As we will see later, this allows significant reduction in the load on the access part of the cellular network. Multicast is not intended for retransmissions, and performance suffers and resources are wasted in the case of a single bad channel user inside the cell, due to trade-offs in coverage and efficiency. By including D2D communications into the picture, we obtain additional performance gains in terms of radio resources. Well-positioned users participate in mitigating the inefficiencies of multicast, by sharing their short-range resources to hand over content to users in bad cellular channel conditions. Depending on the number of participants requesting data, we find a break-even point that achieves a good trade-off in terms of covered users and reception delay.

To assess the performance of this joint multicast/D2D approach it is necessary to evaluate the amount of radio resources consumed at the base station. This leads us to introduce a finer model of radio resource consumption than previous works in the literature. Existing proposals do not consider heterogeneous channel conditions and assume that delivering a given amount of data to different users has always the same cost. Such an assumption does not hold in reality, as radio resources vary according to the channel condition experienced by each user. In other words, transmitting the same piece of content to users with different channel conditions do lead to uneven costs at the base station. To the best of our knowledge, we are the first to evaluate this aspect in the context of data offloading. Note that our results are one of the basis for designing cellular scheduling policies mixing together multicast and D2D transmissions. D2D should be intended in a broad sense, and includes both the standardized LTE-D2D technique, as well as solutions exploiting other technologies for D2D communications, such as WiFi and Bluetooth. Nevertheless,

---

[6] Note that a more precise terminology would be "multicast/broadcast", because only a subset of nodes is concerned by the content (multicast), and the shared nature of the wireless medium (broadcast) is exploited to transmit data. For the sake of readability, in the following we will only employ the term "multicast".

we consider this part of the work more related to intra- than inter-technology scheduling, as we see a more direct applicability to scheduling policies implemented by a cellular operator inside its network.

In the following of the section we provide an extended summary of the work undertaken on this topic. Additional details are available in [51], also included as Appendix A.

## 5.1  Multicast in 4G networks

LTE proposes an optimized broadcast/multicast service through **eMBMS** (*enhanced Multimedia Broadcast Multimedia Service*), a point-to-multipoint specification to transmit control/data information from the cellular base station (eNB) to a group of user entities (UEs).

Cellular UEs can use different modulation and coding schemes (**MCS**) to deal with variable channel characteristics. Each UE experiences different radio conditions, depending on path loss, interference from other cells, and wireless fading. UEs that are closer to the base station are able to decode data at a higher rate, while others located near the edge of the cell have to reduce their data rate and use a degraded MCS. This heterogeneity (time-varying and user-dependent) reduces the effectiveness of multicast because the eNB uses a single MCS to multicast downlink data. The selected MCS for multicast should be robust enough to ensure the successful reception and decoding of the data-frame for each recipient inside the cell. Thus, the worst channel among all the receivers dictates performance. An increase in the number of UEs boosts the probability that at least one UE experiences bad channel conditions, degrading the overall throughput.

To quantify this effect, we simulate a 500 x 500 square meters single LTE cell with an increasing number of randomly located receivers using the ns-3 simulator. Figure 21 presents the average minimum channel quality, in terms of CQI (Channel Quality Indicator), reported at the eNB by UEs (static). The reported CQI is a number between *zero* (worst) and *15* (best). The CQI indicates the most efficient MCS giving a Block Error Rate (BLER) of *10%* or less.



**Figure 21: Minimum CQI for different multicast group sizes. 100 runs, confidence intervals are tight (not shown).**

The average minimum CQI value decreases as the number of users in the multicast group increases. The result is that augmenting the number of multicast receivers clearly affects the attainable cell throughput. This greatly motivates us to investigate methods to cope with the inefficiencies of multicast.

## 5.2   Joint D2D / multicast offloading

We address the distribution of popular content to a set of *N* mobile UEs inside a single LTE cell. We want to transmit data to each UE with a guaranteed maximum deadline D at the minimum cost for the cellular infrastructure. We exploit D2D connectivity and store-and-carry forwarding at UEs.

A UE with good channel quality can obtain higher bit-rates with the same amount of resource blocks (RBs), while bad channel users consume more RBs in order to transmit the same amount of data. We capture the allocation expenditure, dynamically ranking the cost of transmitting the content to UEs according to their instantaneous CQI values.

The principles behind our approach are:

(1) at initial time, the eNB sends data to the $I_0$ UEs with the best radio conditions through a single multicast emission;

(2)  the UEs that have received the data $I_0$ start disseminating it in a D2D (epidemic) fashion;

(3) Before the deadline, we define a time interval, a panic zone where all the nodes that have not yet retrieved the content receive it through unicast LTE emissions.

The proposed scheme allows all UEs to receive data by the deadline (as long as the panic zone is sufficiently large). It adapts to different deadlines -- the larger ones allowing for more D2D dissemination. Its performance relies essentially on one key parameter $I_0$ that characterizes the number of UEs reached by the initial multicast transmission. This immediately improves the usage of resources at the eNB, because it excludes the $N - I_0$ worst-channel UEs.



**Figure 22: UEs can decode data with a maximum modulation schema depending on their position in the cell. The eNB may decide to multicast at higher rate (E.g., MCS index 12). UEs unable to decode data are reached through out-of-band D2D links.**

Figure 22 offers a representative example of the proposed strategy with 6 UEs in the cell. In the D2D dissemination phase, *outaged* UEs benefit from nearby nodes, fetching data directly from them through D2D transmissions. This cooperative strategy is by far more efficient in terms of cellular resource consumption than multicast alone, given that the transmission rate increases and the D2D links typically exploit a much larger bandwidth than cellular communications.

## 5.3   Performance Evaluation

We compare the performance of the proposed joint distribution system with the one achieved by the classic cellular multicast alone. All the results presented in this section are averages over 25 independent simulation runs.

We consider a static number of UEs within the cell for each simulation run, to prove the validity of the concept. Future work will tackle the case where UEs can enter and exit the distribution area. Node mobility is implemented according to the random waypoint model with speed fixed at 27 m/s and pause-time set at 0.5 s. We simulate UDP constant bit-rate downlink flows, each one with packet size $s_k$ = 2048 bytes and a total load of 8 Mb.

We implemented our joint D2D/multicast strategy employing the MOTO simulator [4]. Since the MOTO simulator does not natively support cellular multicast, we implemented an additional module that interacts with the packet scheduler emulating single-cell multicast. The multicast module receives the CQI reports of UEs and decides the transmission rate following the steps explained in Section 5.2. Further parameters for the simulations can be found in Appendix A.

**Reference Strategies: No D2D** is the basic strategy, where UEs have no direct connectivity options, and multicasting through the cellular infrastructure is the only means of distributing content. We compare this base case to our joint D2D/multicast strategy. We assess the performance for three different values of **N** -- the number of users inside the cell -- respectively **10**, **25**, and **50**, so to evaluate performance under different loads. We also consider various values for the parameter **$I_0$** -- the number of direct multicast recipients. In order to be consistent with the notation, we evaluate this value as a percentage of **N**.

**Reception Methods:** Figure 23 provides the fraction of packets partitioned by their reception method. For now, we focus only on their relative weight. As expected, the fraction of packets delivered through multicast follows **$I_0$**. The fraction of panic and D2D messages strongly depends on the parameters *D* and *N*. Tight service delays leave less time to opportunistic distribution to reach *outaged* UEs, resulting in a more intense use of panic retransmissions.

We can find a small amount of packet retransmitted during panic zone even in the **No D2D** strategy. These are packets incorrectly decoded by UEs during the initial multicast emission. In the other strategies, D2D allows not to make use of retransmissions where possible, because UEs can retrieve missing packets from other UEs. For instance, the strategies *No D2D* and *100%* have the same fraction of multicast reception, but differ on the amount of panic and D2D messages. We note also that for sufficiently long deadlines, panic zone is never triggered, and D2D transmissions meet the goal of guaranteeing total data diffusion.



**Figure 23: Data packet ranked by reception method. Multicast and Panic flows through the cellular infrastructure, D2D is on the Wi-Fi channel.**

**Cellular Resources:** Mobile operators are primarily concerned about radio resource usage. This gives hints on the actual amount of RBs devoted to distribute data in the considered scenarios. Unlike the previous case, Figure 24 focuses on the amount of consumed radio resources at the eNB.

We note that the parameter *N* strongly affects the number of employed resources. This is even more evident if we consider very short deadlines. While the amount of resources devoted to multicast only slightly increases with the number of UEs, the impact of unicast re-injections heavily depends on the number of UEs in the cell. In some cases, a small fraction of unicast transmissions could translate into great resource usage. For large *N*, the choice of good values of $I_0$ becomes fundamental in order to avoid congesting the cell with too many panic retransmissions.

*The interesting result is that for any possible value of N and D we may always find a joint strategy that offers better efficiency than No D2D.*

**Figure 24: Average resource blocks employed at eNB to reach 100% dissemination. Note that even few panic zone retransmissions (in unicast) result very costly in resources.**

In this section, we have presented a hybrid distribution system for popular content with guaranteed delays. Multicast is a valuable option to distribute popular data into a cellular network. However, performance is limited by the channel quality of the worst UE in the cell. In this context, we propose a framework that exploits D2D capabilities at UEs to counter the inefficiencies of cellular multicast.

The performance of a joint D2D/multicast strategy is evaluated by varying the number of UEs in the cell and the maximum reception deadline. Simulation results prove that the use of D2D communications allows increasing the multicast transmission rate, saving resources and improving the overall cell throughput.

# 6 Intra-technology scheduling: Towards energy efficiency in the LTE network

In the past decades, the design of cellular networks was primarily aimed at delivering the maximum performance in terms of coverage and/or capacity performance neglecting the energy efficiency of the solutions adopted in the core or the access parts of the network. Energy efficiency in cellular network has become a central point in recent research efforts as it became evident that the operation of this sector is responsible of a non-negligible part of the energy bill of an operator. Furthermore, the access part of the cellular networks has drawn significant attention being the biggest contributor to this energy footprint. Consequently various approaches have been employed to improving its energy efficiency, including putting to sleep the electronic components in the eNB (or the entire eNB) whenever possible combined with dynamic coverage adjustment, data offloading through D2D communications and so on. However, the problem of improving the energy efficiency of eNBs is complicated by the recent adoption in LTE networks of a multi-tier architecture, in which a mix of macrocells and smaller cells (namely microcells, picocells and femtocells) coexist and cooperate. Indeed, the deployment of heterogeneous LTE networks is seen as a cost-effective solution to ever-growing traffic demands, because small cells have lower implementation costs, use less expensive equipment and consume less energy than traditional macrocells. On the other hand, the introduction of small cells may potentially lead to increased ICI (Inter Cell Interference) due to intense frequency reuse in neighbouring or overlapping cells. For this reasons, LTE standards envisage the use of various techniques, such as ICIC (Inter-Cell Interference Coordination) to mitigate the negative effect of the interference between overlaying macro- and microcells.

In this section, we use the ns3 simulator extended with LENA module to start exploring the trade-off between capacity increase and energy efficiency in a heterogeneous LTE cell in which a macrocell coexists with multiple outdoor picocells. In particular, we assume that in each cell there is an *energy scheduler* that decides which eNBs to activate and when. Then, we investigate different strategies for switching off eNBs based on network status without degrading the overall cell capacity. Note that in the considered network scenario we also assume that the picocells employ Cell Range Extension (CRE) techniques to offload data from nearby macrocells. Such technique consists in adding a positive range expansion bias to the pilot downlink signal strength received from picocell so that more users connect to them. A novelty of our study is that we consider different requirements for different data types. In particular we assume that existence of **delay tolerant traffic**, which is amenable to be offloaded using terminal-to-terminal communications.

In this initial study we consider energy saving at two levels. On the one hand, we compare configurations with a single macrocell with others where the macrocell uses a lower transmit power and exploits picocells to serve additional UEs. On the other hand, we consider the effect of policies that switch off some picocells (those that are serving fewer users). Our results show that, in general, reducing the transmit power of the macro cell and complementing this with picocells is viable, as the overall throughput provided to the users is not negatively affected. The additional gain of switching off picocells is still to be better investigated, as for now the opportunity to switch them of does not arise too frequently. We are currently carrying over this analysis to better understand when this can, instead, provide an advantage.

## 6.1 Power consumption model

We use the power consumption model of a BS with the e3F model introduced by Auer et al. [32]. This model provides a relation between $P_{out}$ (the output power radiated by antenna) and $P_{in}$ (the total power needed by the eNB to operate) for different types of eNBs in an LTE system with 10 MHz of bandwidth and 2x2 MIMO antenna configuration. The energy model is well approximated by this linear model:

$$P_{in} = \begin{cases} N_{TRX}(P_0 + \Delta_P P_{out}) & 0 < P_{out} < P_{max} \\ N_{TRX}P_{sleep} & P_{out} = 0 \end{cases}$$

where $N_{TRX}$ represents the number of transceiver chains, $P_0$ represents the power consumption at the minimum non-zero output power (empty eNB), $\Delta_P$ is the slope of the load dependent power consumption, $P_{max}$ represent the maximum transmission power achievable by the BS and $P_{sleep}$ represents the power consumption of the eNB in sleep mode. Table 4 reports the typical values for the parameters of the power consumption model for different types of eNBs.

**Table 4: Power model for different eNB types**

| eNB type | $N_{TRX}$ | $P_{max}$[W] | $P_0$[W] | $\Delta_p$ | $P_{sleep}$[W] |
|----------|-----------|--------------|----------|------------|----------------|
| Macro | 6 | 20 | 130.0 | 4.7 | 75.0 |
| RRH | 6 | 20 | 84.0 | 2.8 | 56.0 |
| Micro | 2 | 6.3 | 56.0 | 2.6 | 39.0 |
| Pico | 2 | 0.13 | 6.8 | 4.0 | 4.3 |
| Femto | 2 | 0.05 | 4.8 | 8.0 | 2.9 |

## 6.2   Simulation set-up

The simulations were performed using ns3 simulator extended with the LENA module for LTE. The simulation parameters are reported in Table 5, while Figure 25 illustrates the network scenario that we have used in the simulations. Specifically, we consider a heterogeneous network in which there is one macrocell and a varying number of picocells are deployed within the macrocell. Furthermore, we consider two different types of deployment for the picocells. The first one is a *total random deployment* in which no constraint is set on the location of the picocell within the macrocell. The second one is a *planned deployment*. In particular, to avoid excessive interference from the macrocell to the picocels we avoid deploying picocells at a distance shorter than 200 meters from the centre of the macrocell. Furthermore, we also avoid deploying two picocells at a distance shorter than 150 meters. Note that in the simulations we have used the *Hybrid Building Propagation Loss* model that is used in ns3 to model propagation losses due to the presence of different types of buildings (i.e., residential, office and commercial). More details on the implementation details of this propagation loss model are provided in Section 6.2.2.



**Figure 25: Network scenario**

To evaluate the benefits of introducing an energy scheduler within a heterogeneous network we have chosen the following metrics.

- *Overall network throughput:* the total throughput measured over the entire network including macro- and small-cells
- *Network energy consumption:* the total energy consumed by the eNBs.

We remind that our ultimate goal is to design a scheduling policy that achieves a good comprise between power consumption reduction and network capacity degradation. Thus, before presenting the simulations results we describe in greater details: *a)* how UEs are connected to macrocells and picocells, and *b)* the algorithm that we have used to decide when a picocell should be switched-off and to which picocell to move its associated UE.

**Table 5: Simulation parameters for energy efficiency in LTE core network**

| eNodeB type | Macro | Pico |
|---|---|---|
| Bandwidth | 5Mhz | 5Mhz |
| Tx Power | 46, 43, 36 dbm | 28, 24 dbm |
| Antenna | parabolic | omnidirectional |
| bias$_{CRE}$ value | 0 | 10 db |
| Number of antenna | 3 | 1 |
| Cell layout | Hexagonal, ISD =500m | 0,6,8,10 per macrocell sector |
| Min dist to picocells | 200 m | 150 |
| Ue max speed | 54 km/h | |
| Mobility Model | Random Waypoint | |
| Pathloss | Urban/Buildings | |
| BS Distruibution | | Random |
| Data Rate | 8.2Mb/sec | |
| UE Distribution | Random (density per m$^2$ 0.0002, 0.00015, 0.0001, 0.00005) | |
| Packets dimension | 1024 KBytes | |
| Simualtion time | 20 s | |

## 6.2.1 Switch-off procedures

At the beginning of each experiment each UE is connected to the closest eNB. Then, standard handover procedures implemented in the LENA module are activated and UEs automatically hand over towards the eNBs to which they receive the highest RSSI. After the initial configuration phase the switch-off algorithm decides which are the picocells to temporarily deactivate by looking at the network load distribution. Specifically, picocells with less than 4 users are switched off and their users are handed over to the closest picocell that has less than 25 users already connected. Note that after this tentative association of UEs to picocells based on network load, the automatic handover procedures implemented in ns3 are used to trigger additional handovers towards picocells with better SINR. For the sake of clarity, the various phases of the switch-off procedures are also illustrated in Figure 26. Note that in this scenario we have introduced a bias CRE value equal to 10 dBm, to enable picocells to capture more traffic and to favor macrocell offload.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

**Figure 26: Switch-off procedures.**

Next, we provide an outline of the implementation of the automatic handover procedures in the LENA module. Specifically, the RRC (Radio Resource Control) model implemented in the simulator provides handover functionality through the shared X2 interface. There are two ways to initiate a handover procedure:

- The handover could be triggered explicitly by the simulation program by scheduling an execution of the method *LteEnbRrc::SendHandoverRequest()*
- The handover could be triggered automatically by the eNB RRC entity. The eNB executes the algorithm Figure 27 for a UE providing measurements in its serving cell and the neighbour cells the UE measures:

**Figure 27:Algorithm to automatically trigger the Handover procedure**

Furthermore, the following parameters can be adjusted to control the handover decision process:

- *servingHandoverThreshold*: if the RSRQ (Reference Signal Received Quality) value measured by the UE in its serving cell is less or equal to the s*ervingHandoverThreshold* parameter (i.e. the conditions of the UE in the serving cell are getting bad or not good enough), then the eNB considers this UE to hand it over to a new neighbor eNB. The handover will eventually be triggered depending on the measurements of the neighbor cells.

- *neighbourHandoverOffset*: if the UE is considered for handover, and the difference between the best neighbor RSRQ and the RSRQ difference between the neighbor and the serving cell is greater or equal to the *neighbourHandoverOffset* parameter, then the handover procedure is triggered for this UE.

The X2 interface interconnects two eNBs though a point-to-point link between the two eNB. The X2 interface implemented in the LENA module provides detailed implementation of the following elementary procedures of the Mobility Management functionality:

- Handover Request procedure

- Handover Request Acknowledgement procedure

- SN Status Transfer procedure

- UE Context Release procedure

The above procedures are involved in the X2-based handover. We note that the simulator model currently supports only the *seamless handover* while *lossless handover* is not supported. Generally speaking, seamless handover is used for channels transporting traffic that is very sensitive to delay and jitter and would rather accept retransmissions than delay, like VoIP, while the lossless handover is used for channels that transport traffic that does not care too much about delay but is sensitive to retransmissions, like FTP and HTTP. Figure 28 shows the interaction of the entities of the X2 model in the simulator.

**Figure 28: Sequence diagram of the X2-based handover**

### 6.2.2 Buildings Module

The Buildings module used in our scenario and implemented in ns3, provides:

1. a new class (*Building*) that models the presence of a building in a simulation scenario;

2. a new class(*MobilityBulidingInfo*) that allows to specify the location, size and characteristics of buildings present in the simulated area, and allows the placement of nodes inside those buildings;

3. a container class with the definition of the most useful pathloss models and the correspondent variables called *BuildingPropoagationLossModel*.

4. a new propagation model (*HybridBuildingsPropagationLossModel*) working with the mobility model just introduced, that allows to model the phenomenon of indoor/outdoor propagation in the presence of buildings.

5. a simplified model working only with the Okumura Hata propagation model (*OhBuildingsPropagationLossModel*), which consider the phenomenon of indoor/outdoor propagation in the presence of buildings.

The models have been designed with LTE in mind, though their implementation is in fact independent from any LTE-specific code, and can be used with other ns-3 wireless technologies as well (e.g., wifi, wimax).

*HybridBuildingsPropagationLossMode*l pathloss model included is obtained through a combination of several well-known pathloss models in order to mimic different environmental scenarios such as urban, suburban and open areas. Moreover, the model considers both outdoor and indoor indoor and outdoor communication has to be included since HeNB might be installed either within building and either outside. In case of indoor communication, the model has to consider also the type of building in outdoor <-> indoor communication according to some general criteria such as the wall penetration losses of the common materials; moreover it includes some general configuration for the internal walls in indoor communications.

*OhBuildingsPropagationLossModel* pathloss model has been created for simplifying the previous one removing the thresholds for switching from one model to other. For doing this it has been used only one propagation model from the one available (i.e., the Okumura Hata). The presence of building is still

considered in the model; therefore all the considerations of above regarding the building type are still valid. The same consideration can be done for what concern the environmental scenario and frequency since both of them are parameters of the model considered.

## 6.3 Simulation Results

In the following figures we report a subset of the most interesting simulation results. Figure 29 shows the total throughput with 8 picocells in the following configurations:

1. 1 macrocell transmitting with a power equal to 43 dBm (red line)
2. 1 macrocell transmitting with a power equal to 36 dBm and 8 picocells with a *random* distribution and transmitting with a power equal to 24 dBm (green line)
3. 1 macrocell transmitting with a power equal to 36 dBm and 8 picocells with a *planned* distribution and transmitting with a power equal to 24 dBm (blue line).

In case 2 and 3 the throughput is lower than the one obtained in case 1 if the number of UEs is small. Indeed, in this case the bandwidth of the macrocell is sufficient to serve all the users and adding the smallcells negatively affects the level of interference in the cell. Furthermore, if the number of picocells is small the probability that a UE has a better connection with a picocell than with the macrocell is typically low. An important finding of Figure 29 is also that a planned deployment of picocell inside the macrocell is beneficial to reduce the ICI, and intuitively this effect is more evident for high number of UEs. For these reasons, in the following we only show results for planned picocell deployments and for larger number of picocells, which provide the most interesting and significant results.



**Figure 29: Throughput measured with 8 picocells.**

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

**Figure 30: Throughput measured with 18 picocells with and without switch-off of lightly loaded picocells.**

In Figure 30 we show the total cell throughput in a network scenario with 18 picocells with and without switch-off of lightly loaded picocells. We can observe that even if the transmission power of the macrocell has been reduced from 43dBm to 36 dBm the overall throughput provided to the users increases due to the additional capacity provided by the picocells. Moreover, in the considered scenario, switching off lightly loaded picocells (i.e., picocells with less than 4 UEs) has a negligible impact on the throughput. Finally, Figure 31 shows the energy consumption in the same network scenarios of Figure 30. Two important conclusions can be derived from the plots. The first one is that in a heterogeneous network with a macrocell and multiple picocells the energy consumption can be significantly less than in a network without picocells because in the former case the macrocell can use a significantly lower transmission power than in the latter case. The second observation is that the additional energy gain due to switching off picocells is way lower than that due to reducing the transmission power of the macrocell. This is due to the fact that the picocells are low power nodes and their energy consumption can be an order of magnitude less than the one of a macrocell. It is also important to point out that a more aggressive scheduling strategy might entail to also switch off the macrocell. This could provide a much higher energy reduction with respect to the case of switching off only the picocells, but at the cost of reduced coverage. On the other hand, when the macro cell cannot be switched off, and its transmission power cannot be further reduced, switching off picocells can still provide a noticeable benefit in terms of lower energy consumption.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

**Figure 31: Power consumption with 18 picocells when switching off lightly loaded picocells**

In conclusion, our preliminary evaluation of switch-off strategies for heterogeneous LTE networks with a macrocell and multiple picocells have shown that it is possible to reduce the transmit power of macrocell without negatively affecting the throughput performance experienced by users, as the effect of a reduction in the macrocell transmit power is compensated by additional picocells. On the other hand, to switch off lightly loaded picocell also does no negatively affects the overall capacity of the LTE network. Furthermore, we have also shown that the throughput obtained by deploying picocells is heavily dependent on the deployment strategy of picocells (random vs. planned deployments) and the density of UEs. Clearly, the efficiency of the switch-off strategy depends on several parameters (e.g, CRE bias, macrocell transmission powers, thresholds to decide to switch off a picocell) and we plan to extend our study to identify an optimal configuration of such parameters. Finally, we also intend to extend our initial study to investigate the energy saving potential of adding D2D communications in an heterogeneous network scenario.

# 7 Inter-technology scheduling: Multi-user offloading in heterogeneous wireless network infrastructures

In this section we present the status of the work on intra-technology scheduling. This is, for now, a complementary yet very important research direction for the WP, as it allows us to better understand the option of using multiple wireless infrastructure to offload traffic from possibly congested cellular networks.

Specifically, we consider the case where users may either be served by an operator through a cellular or a WiFi infrastructure, and investigate how to best schedule users across these technologies. The objective of this work is to investigate handover decision making algorithms in heterogenous networks and point out the metrics and factors influencing data offloading and related open research issues to the research community.

To this extent, two multi-user multiple attribute decision making (MADM) algorithms have been developed. These algorithms are built on the TOPSIS framework [42] adapting it to the case of offloading across wireless infrastructure. This is basically an optimization framework where we can express decision criteria based on QoS metrics. We thus identify appropriate QoS metrics for our case, and show how to customize the optimization framework for our purposes, with specific reference to real-time traffic. Specifically, we consider a straightforward application of TOPSIS, where we optimize individual users' performance. On the other hand, in the second algorithm we take into consideration network-wide performance, trying to optimize the minimum bit-rates by assigning the users in the intersection area to WLAN where cellular utilization is high. The objective of this part of the work is to investigate handover decision making algorithms in heterogeneous networks and point out the metrics and factors influencing data offloading and related open research issues to the research community. To this extent, a capacity aware multi-user multiple attribute decision making (MADM) algorithm based on QoE metrics has been developed and evaluated. The proposed capacity aware multi-user load balancing algorithm optimizes total benefit of the system that is balanced according to total channel utilization among different heterogeneous networks. We performed initial simulation evaluation in the NS environment. The main purpose of the simulations is not to provide a comprehensive assessment of the algorithm, but to show that the basic idea behind its definition is appropriate in simple – and thus easy to understand – configurations. The proposed algorithm is shown to enhance total channel utilization of heterogeneous networks compared to standard single-user decision making algorithms. Specifically, we show that the algorithm optimizing global benefit (as opposed to individual user metrics) avoids to overload the cellular network. The alternative algorithms, instead, drastically overload the network, and this therefore results in uncontrolled ping-pong effects (as users need anyway to switch to a different technology), which is in the end very detrimental also from the standpoint of the individual user.

## 7.1 TOPSIS-based solutions in the context of ongoing research

Resource management in heterogeneous Wireless Networks or WiFi integrated cellular networks could be coordinated through user-centric models, network-centric models or collaborative schemes. User centric models offer ease if implementation and scalability, as opposed to the other two approaches, at the expense of reduced –system wide- efficiency. The network-centric models, on the other hand, provide more efficient solutions that improve "social welfare" (addressing both 3GPP and non-3GPP subsytems), at the cost of increased control overhead and risk of "single point of failure". Collaborative solutions, on the other hand, introduce a little bit complexity; however, in return, offer drastic performance difference with respect to network-centric solutions in terms of QoE. Basically, in a collaborative solution, UE data such as RSS or CQI along with access network metrics obtained from an operator server are combined in decision making phase and based on the implementation choice either on network-side or user-side the decision is executed.

The handoff decision algorithm aims at selecting a network for a particular service that can satisfy objectives based on some criteria (such as low cost, good Received Signal Strength (RSS), optimum bandwidth, low network latency, high reliability and long life battery) and taking into account the preferred access network of user. Some techniques used for network-centric solutions such as stochastic programming, game theory and utility function could be performed in this respect [48].

In network-centric approaches, the goal is often to acquire maximum total allocation in 3GPP and non-3GPP networks while minimizing cost of underutilization and demand rejection. Stochastic linear programming obtains maximum allocation in each network by using probabilities related to allocation, underutilization, and rejection in Heterogeneous Wireless Networks (HWNs). Game theoretic approaches take advantage of the bankruptcy game, and efficient bandwidth allocation and admission control algorithms are developed by utilizing available bandwidth in each network. In utility function, operator prioritizes users and classifies services to allocate bandwidth for the users [45][63][66].

In user-centric solutions, the users themselves (or their agents) make the decisions, often prioritizing the needs and objectives of the individual users. Analytical hierarchy processes help ranking the networks based on induced QoS indicators, by checking user's requirements and network conditions. Proposed approaches make use of the consumer surplus model and similar economic theory based techniques. Users are often modelled to have profit functions amounting to the difference between bandwidth gain and handoff cost for each network is computed. The most appropriate network is found through utility maximization [40][46][57].

As for the collaborative models, fuzzy logic controller ranks the candidate networks based on the user's selection criteria, network data rate and SNR. In objective function, user's RSS, network's queue delay and policy preferences such as cost are fed as input parameters, and the function provides the allocation of services to APs and terminals. Lastly, in TOPSIS, the best path for flow distribution on muti-homed end-hosts is computed. Also, network's QoS (delay, jitter, and BER), user's traffic class and most importantly QoE are also considered [34][65][7].

Other options could be to harness impatient or patient algorithms which are based on user-centric solutions. The Impatient algorithm uses a very simple policy: use 3G whenever WiFi is unavailable; else use WiFi. The Patient waits and sends data on WiFi until the delay tolerance threshold, and only switches to 3G if all of the data are not sent on WiFi before the delay tolerance threshold [6]. This are the standard approaches investigated in WiFi-based offloading.

## 7.2   Performance Metrıcs Influencing Data Offloadıng and System model

As for decision making functionality, UE or Mobile Network Operator (MNO) selects the access network by considering probabilistic demands. Network related, terminal related, user related and application related metrics need to be considered pertaining to handover decisions. However, the paramount elements amongst them are the user-related ones as QoE is at the very hearth of contemporary mobile business performance expectations. Related parameters include throughput, energy consumption of the terminal, security etc. It is interesting to note that an adult's preferences along these dimensions would potentially differ from that of a young person. For instance, security-wise an adult might not prefer to watch videos through WEP or WPA on WiFi networks but EAP-SIM on 3GPP network. Maybe this choice could be trivial for a young person and actually he would prefer a free communication band, but considering recently emerging security challenges, operators need to pay importance to the subjects of security and privacy pertinent to each and every user they serve [30].

Handoff decision criteria can be categorized as below:

- **Network-related** considering coverage, bandwidth, latency, link quality (RSS (Received Signal Strength), BER (Bit Error Rate), cost, security level.

- **Terminal-related** considering, e.g., velocity or battery powe

- **User-related** considering user profiles and preferences

- **Service-related** considering service capabilities, QoS, QoE, security level [30].

The Quality of Service (QoS) and Quality of Experience (QoE), mobility and network architecture are important factors during decision making or network selection phase. The following QoS and QoE metrics are important to be checked while offloading the data traffic due to the nature of real-time applications:

(a) **End to end delay (s):** This includes processing, queuing in both ingress and egress, and propagation delay. The end-to-end delay of a video signal is the time taken for the packets to enter the transmitter at one end, be encoded into a digital signal, travel through the network, and be regenerated by the receiver at the other end.

(b) **Data received (Kbps):** This is calculated based on the successfully received packets.

(c) **Packet Loss (%):** This is calculated based on the dropped packets due to either network problems or some queuing problems.

(d) **Throughput (Kbps):** this is the total traffic where packets are successfully received by the destination excluding packets for other destinations.

(e) **MOS Value (Mean Opinion Score):** This corresponds to a numerical value, ranging between 1(worst) and 5(best) expressing the quality perceived by user. It is also used as a QoE metric.

(f**) Jitter (s):** In IP networks, jitter is the variation in the time-of-arrival of consecutive packets. Jitter results from a momentary condition where more packets are trying to get on a particular link than the link can carry away [41].

Considering these performance metrics as reference, we assume the following notation and model to represent the multiple user multiple attribute decision making problem:

- The total users set in the system is denoted as $U = \{u_1, u_2, u_3, ..., u_k\}$ where k (k>=2) denotes number of users.

- The multiple users' set involved in the decision making process are denoted as $V = \{v_1, v_2, v_3, ..., v_{k'}\}$ where k' (k'<=k) denotes number of users under multiple coverage.

- The multiple attribute set is denoted as $S = \{s_1, s_2, s_3, ..., s_m\}$ where m (m>=2) denotes number of possible attributes.

- The multiple decision point set is denoted as $E = \{e_1, e_2, e_3, ..., e_P\}$ where there are p (p≥2) possible decision points.

- The weight set is denoted as $w = \{w_1, w_2, w_3, ..., w_m\}$, where each weight $w_i$ is the weight assigned to attribute $s_i$ i ∈ {1,2,...,m}.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

**Figure 32. A sample user distribution map under multiple wireless technology coverage**

We use this model to exploit the TOPSIS framework in order to decide how to best allocate users to the possible wireless technologies under consideration. As can be seen in Figure 32, the set of users that are in the coverage area of both WLAN and LTE are shaded with gray area. These users have the high potential of handover and have to make a smart decision to select the best access point. Therefore, our method runs on the scenarios based on the users that are concentrated on this region.

## 7.3 Multiuser offloading algorithms for heterogeneous networks

In this Section we first recall the main features of the standard TOPSIS framework, and then describe how we adapt it to our specific cases.

### 7.3.1 TOPSIS

TOPSIS (Technique for Order Preference by Similarity to Ideal Solution) [42][24], due to its easy implementation, is a suitable candidate to select the optimal target network for a given a set of given observed attributes for a user.

In the first step of TOPSIS algorithm a decision matrix **A** is created:

$$\mathbf{A} = \begin{bmatrix} a_{ij} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & ... & a_{1m} \\ a_{21} & a_{22} & ... & a_{2m} \\ . & . & ... & . \\ . & . & ... & . \\ . & . & ... & . \\ a_{p1} & a_{p2} & ... & a_{pm} \end{bmatrix} (i = 1,....,p; j = 1,....,m)$$

In matrix $A^{i'}$ matrix, m refers to size of the multiple attribute set such as link quality, MOS of the target network for the given application, user preference (cost security), etc and p refers to size of the multiple decision points target networks which can be LTE, WLAN or D2D (device-to-device). Note that that all the attributes are transformed to have positive impact if necessary.

In second step, a normalized decision matrix is formed by using the following equation:

$$r_{ij} = \frac{a_{ij}}{\sqrt{\sum_{k=1}^{p} a_{kj}^2}}$$

Then the normalize matrix R is obtained as:

$$\mathbf{R} = [r_{ij}] = \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1m} \\ r_{21} & r_{22} & \dots & r_{2m} \\ . & & & . \\ . & & & . \\ . & & & . \\ r_{p1} & r_{p2} & \dots & r_{pm} \end{bmatrix}$$

In third step, a weighted normalized decision matrix is created by multiplying each column of the matrix by corresponding weight wi where $\sum_{i=1}^{m} w_i = 1$ by using the following equation:

$$\mathbf{v_i} = w_i * \mathbf{r_i}, \qquad \mathbf{r_i} = [r_{1i}, ...., r_{pi}]^T, \qquad i = \{1, 2, ..., m\}$$

In fourth step, the positive ($A^*$) and negative ($A^-$) solutions are formed by using the following formulas:

$$A^* = \left\{ (\max_i v_{ij} \,|\, j \in \{1, 2, ..m\}) \right\}$$

$$A^- = \left\{ (\min_i v_{ij} \,|\, j \in \{1, 2, ...m\}) \right\}$$

At the end of fourth step, we end up with the following sets: $A^* = \{v_1^*, v_2^*, ..., v_n^*\}$ and $A^- = \{v_1^-, v_2^-, ..., v_n^-\}$

By calculating the Euclidean distance $S_i^*$ of each multiple decision point from the positive point $A^*$ and $S_i^-$ of each multiple decision point from the negative point $A^-$.

$$S_i^* = \sqrt{\sum_{j=1}^{p} (v_{ij} - v_j^*)^2}, \quad i = \{1, ..., p\}$$

$$S_i^- = \sqrt{\sum_{j=1}^{p} (v_{ij} - v_j^-)^2}, \quad i = \{1, ..., p\}$$

In the final step, the relative similarity of the alternatives from the positive and negative point is calculated as:

$$C_i = \frac{S_i^-}{S_i^- + S_i^*}, \quad i = \{1, ..., p\}$$

where $0 \le C_i \le 1$ the final solution is selected by:

$$e^* = e_{i^*} \text{ where } i^* = \arg\max_i C_i, \quad i = \{1,...,p\}$$

### 7.3.2 Multiple attribute sets in TOPSIS algorithm

In this study, we used TOPSIS as our core algorithm [42], due to its easy implementation, as a way of selecting the best target network for a given user's video application. The decision to use this algorithm was made based on the other multiple attribute decision making (MADM) algoritms' performance comparison results. In [61], four different MADM algorithms (MEW, SAW, GRA, TOPSIS) were evaluated and it was concluded that they all performed very similar.

Decision parameters of the TOPSIS are as follows:

(i) MOS: Mean Opinion Score is considered as a subjective measure. Currently, it is more often used to refer to one or another objective approximation of subjective MOS. ITU P.800 and P.830 define the MOS scale as showed in Table 6.

(ii) PSNR (dB): The peak signal-to-noise ratio is used as an objective measurement of the restored quality. PSNR is most commonly used to measure the quality of reconstruction of lossy compression codecs. PSNR is defined as follows:

$$PSNR = 20\log\frac{V_{peak}}{MSE}$$

where Vpeak = $2^k - 1$ and k is equal to number of bits per pixel. MSE is standard mean squared error. In case of multimedia real-time traffic, we calculate the PSNR frame by frame and map it to the corresponding MOS value as follows:

**Table 6- PSNR to MOS mapping**

| PSNR [dB] | MOS |
|-----------|-----|
| > 37 | 5 (Excellent) |
| 31 - 37 | 4 (Good) |
| 25 - 31 | 3 (Fair) |
| 20 - 25 | 2 (Poor) |
| < 20 | 1 (Bad) |

(iii) CQI: Channel quality indicator is reported by UE and is calculated using BLER and SNR values. It is a vital parameter to estimate the UMTS air interface quality. The UE type that is assumed in the simulator is 3GPP UE category 1 to 6. In our simulation, the highest CQI value was accepted as 22. However, it varies between 1 and 22.

(iv) QoS: Quality of service level of the access point (AP) is utilized in the algorithm to determine the link-quality of WiFi network. Voice = Platinum = 6, Video = Gold = 5, Best Effort = Silver = 3, Background = Bronze = 1

(v) Security Policy used in WiFi network: WPA or WPA2 cannot be used for a seamless solution. EAP-SIM is required to do so.

(vi) Channel Utilization: Channel utilization is a term which is used to measure the channel usage taking into account the throughput as well as the overhead. For example, in latency-sensitive applications over wireless, such as voice, channel utilization is used to measure the usage percentage of the RF channel. It is a network parameter, and is monitored for a stable traffic level and to prevent under or over utilization.

Note that we also define individual Quality-of-Experience (QoE) value of users as the weighted sum of these attributes.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

In the next two sub-sections, we define two algorithms based on TOPSIS, configured according to the decision parameters above. The first algorithm is a Multi-user TOPSIS with capacity-aware characteristic where channel utilization parameter is of the utmost importance for the 3GPP network to balance the channel allocations. With this type of multi user algorithm, the total system benefit is considered as important. The second algorithm is a Standard TOPSIS (ST) algorithm. With this method each user's individual benefits are considered individually as they arrive.

### 7.3.3 Capacity aware multi-user iterative TOPSIS (CAT) algorithm

In order to obtain certain benefits for access channel selection and resource allocation problem between multiple users, we propose Capacity aware iterative multi-user TOPSIS algorithm:

**Input:** Set of technology E, total channel utilization threshold for each technology $e \in E$, $CU_{th}^e$ , and the TOPSIS matrix of user $v_{i'} \in V$ denoted by

$$\mathbf{A}^{i'} = \begin{bmatrix} a_{ij} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & ... & a_{1m} \\ a_{21} & a_{22} & ... & a_{2m} \\ . & . & ... & . \\ . & . & ... & . \\ . & . & ... & . \\ a_{p1} & a_{p2} & ... & a_{pm} \end{bmatrix} (i = 1,....,p; j = 1,....,m)$$

**Output:** Capacity-aware channel utilization vector $\mathbf{CU^e}$ = $[CU_1^e, CU_2^e, ...., CU_{k'}^e]$, $e \in E$.

**Step1:** Set $\mathbf{CU^e}$=[0] and i'=0 (i' ≤ k' is the user number)

**Step2:** Put i'= i'+1, as user $v_{i'}$ arrives.

**Step3:** Run TOPSIS algorithm using $A^{i'}$ and select the optimal decision point e* =$e_n \in$ E and construct coincidence coefficient

$$\delta_{i'}^e = \begin{cases} 1 \ if \ e = e* \\ 0, otherwise \end{cases} , \forall \ e \in E$$

**Step4:** Update the temporary channel utilization vector $\widehat{CU}^{e*} = \mathbf{CU}^{e*}$ and put $\widehat{CU}_{i'}^{e*} =$ $a_{nc}$ where c denotes the column number in $A^{i'}$ for attribute corresponding to CU $\in$ S..

**Step5:**

- If ($\sum_{j=1}^{i'} \delta_j^{e*} CU_j^{e*} \leq CU_{th}^{e*}$ )
  - o $\mathbf{CU}^{e*} = \widehat{CU}^{e*}$
- else
  - o E = E \ e*
    - If E = {} then e* = WLAN
    - else go to **Step 3**

Using this multi user algorithm, the total system benefit is considered as the first criteria to optimize and we are trying to maximize the minimum bit-rates by assigning the users in the intersection area to WLAN where 3GPP utilization is high. In our analysis, we will be using channel utilization as the selected parameter to be of utmost importance for the 3GPP network to balance the channel allocations.

### 7.3.4 Standard TOPSIS (ST) method

With standard TOPSIS method, user's individual's benefits are considered. The method details are explained in the following steps:

**Input:** Set of technology E, and the TOPSIS matrix of user $v_{i'} \in V$ denoted by

$$\mathbf{A}^{i'} = \begin{bmatrix} a_{ij} \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & ... & a_{1m} \\ a_{21} & a_{22} & ... & a_{2m} \\ . & . & ... & . \\ . & . & ... & . \\ . & . & ... & . \\ a_{p1} & a_{p2} & ... & a_{pm} \end{bmatrix} (i = 1,....,p; j = 1,....,m)$$

**Output:** Standard TOPSIS channel utilization vector $\mathbf{CU}^e = [CU_1^e, CU_2^e, ...., CU_{k'}^e]$, $e \in E$.

**Step1:** Set $\mathbf{CU}^e=[0]$ and i'=0 (i' ≤ k' is the user number)

**Step2:** Put i'= i'+1, as user $v_{i'}$ arrives.

**Step3:** Run TOPSIS algorithm using $A^{i'}$ and select the optimal decision point e* =$e_n$ ∈ E and construct coincidence coefficient

$$\delta_{i'}^e = \begin{cases} 1 \ if \ e = e * \\ 0, otherwise \end{cases} , \forall \ e \in E$$

**Step4:** Update the channel utilization vector by $CU_{i'}^{e*} = a_{nc}$ where c denotes the column number in $A^{i'}$ for attribute corresponding to CU ∈ S.

## 7.4  Performance Results

### 7.4.1 Simulation Scenario

For simulations, we used the NS simulation environment. EURANE (Enhanced UMTS Radio Access Network Extension), NIST (National Institute of Standards and Technology) and EVALVID packages are used in order to evaluate the video performances in a heterogeneous network during a handover execution where the above algorithms' results are utilized.

We assume a video streaming use case scenario. In order to offload a video streaming seamlessly, only most relevant parameters were selected such as Channel utilization, MOS, QoS, delay, and energy as shown in Table 7.

**Table 7 - Network Selection Criterias for Video Streaming**

| Parameters | Rank | Weight |
|---|---|---|
| MOS | 2 | 0.25 |
| QoS | 3 | 0.1 |
| Energy | 5 | 0.05 |
| Channel Utilization | 1 | 0.5 |
| Delay | 4 | 0.1 |

For our scenario, channel utilization and MOS of service are of utmost importance and therefore the weight coefficients are distributed accordingly as shown in Table 7 which also ranks access networks based on their weight coefficients.

For the purpose of comparison, the same attribute values and the same weights are assigned to attributes for the different algorithms presented above. Note also that assignment of weights could be initiated by either user or operator or collaboratively.

In order to compare algorithms, we consider an environment where k'= 4 users are under multiple coverage and their respective attribute weight values are the same. The decision points (i.e., the wireless technologies where users can handover to) is limited to WLAN and 3GPP networks.  i.e. E = {LTE,WLAN} and

$$\mathbf{A}^1 = \begin{bmatrix} 6 & 3 & 5 & 6 & 7 \\ 4 & 3 & 4 & 6 & 3 \end{bmatrix}, \ \mathbf{A}^2 = \begin{bmatrix} 1 & 5 & 5 & 4 & 6 \\ 4 & 3 & 3 & 7 & 6 \end{bmatrix}, \ \mathbf{A}^3 = \begin{bmatrix} 1 & 5 & 5 & 7 & 6 \\ 7 & 1 & 2 & 2 & 1 \end{bmatrix}, \ \mathbf{A}^4 = \begin{bmatrix} 1 & 5 & 5 & 2 & 6 \\ 7 & 6 & 6 & 4 & 4 \end{bmatrix}$$

The channel utilization threshold for LTE is,  $CU_{th}^{LTE}$ = 8 units,  the channel utilization threshold for WLAN is $CU_{th}^{WLAN}$ = 12 units,

### 7.4.2 Results

The total user distributions on a HetNET comprising 3GPP and WLAN access networks for CAT, ST algorithms as well as ALL 3GPP and ALL WLAN scenarios is shown in **Table 8**. When CAT algorithm is applied, the number of users among different acces technologies are %25 and %75 for 3GPP and WLAN respectively. For the ST algorithm, the distributions become %75 and %25 for 3GPP and WLAN respectively.

**Table 8- TOTAL USER DISTRIBUTION AND CHANNEL UTILIZATIONS (%) ON A HETNETS OF ALL ALGORITHMS:**

|  | USERS DISTRIBUTION | | TOTAL CHANNEL UTILIZATION (%) | |
|---|---|---|---|---|
|  | **3GPP** | **WLAN** | **3GPP** | **WLAN** |
| **CAT** | %25 | %75 | %50 | %66 |
| **ST** | %75 | %25 | %225 | %50 |
| **ALL 3GPP** | %100 | %0 | %275 | 0 |
| **ALL WLAN** | %0 | %100 | 0 | %75 |

Similarly, the total channel utilization distributions for CAT, ST algorithms as well as ALL 3GPP and ALL WLAN scenarios are also shown in **Table 8**. The total channel utilization percentages is calculated by dividing sum of the demands of users for channel utilization over channel utilization thresholds of each technology.

When the CAT algorithm is applied, where total benefit of the system is optimized according to multiple attributes descibed above, balancing the total channel utilizations among 3GPP and WLAN technologies provides lower channel utilizations yielding high capacity. The CAT algorithm has the final total channel utilization percentages of %50 and %66 over 3GPP and WLAN technologies respectively. On the other hand, when ST algorithm is applied, the TOPSIS algorithm will prioritize individual user benefits or indivual QoE (Quality-of-Experience). It is clearly seen that the ST algorithm and/or a random assignments of users could lead to high channel utilization which consequently would decrease MOS substantially for the corresponding access networks.  For ST algorithm, channel utilization percentage of %225 represents over channel utilization for 3GPP.

The important thing to notice for ST algorithm is that even though the expected individual QoE will be high with this type of algorithm, due to overallocation in one access network after the handover decisions are executed, the users will suffer from either ping-pong effect or real-time network changes which will induce additional burden into the system both in terms of network and terminal. However with CAT algorithm, after prioritizing channel utilization and MOS attributes, channel utilizations are optimized between 3GPP and WLAN access networks, which in return increases the QoE of users compared to simple ST algorithm.

*D3.1 Design and evaluation of enabling techniques for mobile data traffic offloading*
*(release a)*
*WP3 – Offloading foundations and enablers*

From the operator point of view, CAT algorithm works best in terms of channel utilization or load balancing; however, with this type of scheme some attributes (other than channel utilization) observed by users can be diminished compared to ST algorithm.

Lastly, in terms of total channel utilization, we compare CAT and ST with ALL 3GPP and ALL WLAN scenarios where no algorithm is implemented and all users are either on 3GPP or WLAN networks. One can observe that channel utilizations for ALL 3GPP exceeds channel utilization thresholds which in return will overallocate the system, adding additional burden to the operators of these wireless access technologies.

## 7.5  Using TOPSIS with D2D technology

Ongoing work on this topic is mainly related to how to include also D2D technologies into the picture. TOPSIS is a centralized scheme, and therefore the most suitable type of D2D technology that we are considering is inspired by LTE-D2D, i.e. where the operator controls the fine details of direct communications between mobile devices.

We hereafter provide a brief sketch of how we are extending the algorithms presented before in this sense.

Assume that we list all the combinations for p different technologies as the decision point in decision matrix as

$$D = \begin{array}{c} \\ sce_1 \\ sce_2 \\ sce_3 \\ \dots \\ sce_l \end{array} \begin{array}{ccccc} T_1 & T_2 & T_3 & T_4 \dots & T_m \\ \left[\begin{array}{ccccc} . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ \dots & \dots & \dots & \dots & . & . \\ . & . & . & . & . & . \end{array}\right] \end{array}$$

Where $sce_i$ is the i-th scenario, $l = \widehat{CU}^{\theta*} = p^{k'}$ and k' denotes the users under multiple coverage.

For example, if we have three technologies, LTE, WLAN and D2D and two users, U1, U2 are under multiple coverage ,then there are 9 possibilities on the rows of D matrix.

$sce_1$= [LTE, LTE]  $\qquad$ $sce_4$= [WLAN, LTE ]  $\qquad$ $sce_7$ = [D2D, LTE]

$sce_2$= [LTE, WLAN]  $\qquad$ $sce_5$ = [WLAN, D2D]  $\qquad$ $sce_8$ = [D2D, WLAN]

$sce_3$= [LTE, D2D]  $\qquad$ $sce_6$ = [WLAN, WLAN]  $\qquad$ $sce_9$ = [D2D, D2D]

Then on the columns of D matrix, there are attributes such as:

$T_1$) Total QoS= $\sum QoS_i$

$T_2$) Total Delay = $\sum \tau_i$

$T_3$) Total Energy = $\sum E_i$

$T_4$) Total Throughput = $\sum Th_i$

$T_5$) Total MOS values = $\sum MOS_i$

…

$T_m$) …

After constructing the decision matrix, we can run TOPSIS to decide for the appropriate possibility of scenarios ($sce_1$, …, $sce_9$) for users that are under multiple coverage. Clearly, careful evaluation of this extension is needed, as the D2D alternatives are typically much more dynamic. Even more importantly, the number of possible combinations increases exponentially with the number of technologies, and this needs to be taken into consideration when the number of users increases.

# 8 Open issues

The results presented in this document already provide solid results about enabling techniques for offload networks. However, open issues still remain. Hereafter, for each of these topics, we briefly discuss the main open points, as well as ongoing work to address them.

*Capacity analysis of opportunistic networks*. Ongoing work here consists basically in refinements and generalisation of what has been presented in this document. Specifically, we are working on studying if convergence can be assessed through aggregate statistics about contact patterns, instead of pairwise statistics. This would be important, among others, because aggregate statistics provide much less information about the behaviour of individual nodes, and are therefore more robust from a privacy standpoint. We are then developing models for end-to-end delay of multi-path opportunistic routing protocols, both when duty cycling is used and when it is not, to better characterise the trade-off between energy consumption and end-to-end delay, and thus additional capacity.

*Capacity analysis of LTE networks*. Ongoing work here is mainly devoted to complete the detailed model of LTE throughput perceived by the user. While the model presented in this document is already a good starting point, several additional aspects should be taken into account, including different scheduling algorithms and the presence of MIMO technologies.

*Capacity analysis of integrated offload networks*. Thanks to the analysis of individual components carried out until now, we are now in the position of better analysing the capacity performance of integrated offload networks. We are progressing at multiple levels of abstractions, primarily as far as the LTE network analysis is concerned. We are progressing both by using quite detailed modelling of all the LTE features, but we are also developing more agile capacity models where we abstract several details of the LTE internals, in order to reduce the complexity of analysing an integrated offload network.

*Intra-technology scheduling.* We are extending the analysis on the use of multicast and D2D technologies together. We are progressing along two main directions. On the one hand, we are defining algorithms to exploit some of the main outcomes of the initial stage of analysis presented in this document, i.e. how to automatically select the optimal set of users for the initial multicast transmissions. On the other hand, we are working to derive capacity results also in this configuration of the MOTO solutions.

*Energy-saving scheduling of LTE elements*. The initial simulation results have shown that in general there is ample room for defining smart algorithms that can switch off part of the LTE network to conserve energy. Defining and evaluating these algorithms is the main objective of the rest of this activity. Note that this type of solutions is very well aligned with a global trend of how to manage operator networks, whereby components are entirely and dynamically switched off in the core of the network, in order to reduce the carbon footprint of the operated network.

*Inter-technology scheduling*. The TOPSIS framework, adapted as presented in this document, is a very flexible tool to analyse inter-technology scheduling in a more innovative way. In particular, we are extending it to also consider the possibility of scheduling terminal-to-terminal communications. In principle, terminal-to-terminal could be seen as an additional technology (in addition to WiFi and cellular) that is available to some of the nodes, and therefore can be put as part of the TOPSIS optimisation framework. This may result in a drastic increase of the complexity of the framework, as the opportunistic network is much more dynamic and features many more parameters that need to be taken into account. We are working in order to better understand this complexity, and to identify solutions to cope with it.

Last, but not least, work that is ongoing and will be started in the third year of the project in WP4 and WP5 will provide very useful feedback in order to better understand if the level of characterisation of the enabling techniques was sufficient or not, and update it if needed. In addition, definition of precise terminal-to-terminal protocols (WP4) and their evaluation in integrated simulation and test environments

(WP5) will allow us to stress test the inevitable abstractions used in this WP to characterise the performance of the key MOTO components, and refine them as needed.

# References

[1] The MOTO consortium, Deliverable D3.1 "Initial results on offloading foundations and enablers", available at http://www.fp7-moto.eu/wp-content/uploads/2013/10/moto_D3.1_v1.0_PU1.pdf

[2] The MOTO consortium, Deliverable D3.2 "Spatiotemporal characterization of contact patterns in dynamic networks" (delivered, under review at the time of writing)

[3] The MOTO consortium, Deliverable D2.2.1 "General Architecture of the Mobile Offoading System (Release a)", available at http://www.fp7-moto.eu/wp-content/uploads/2014/01/Deliverable_D21_1.01.pdf

[4] The MOTO Consortium, "D5.1.1 – Description and development of MOTO simulation tool environment – Release a", available at http://www.fp7-moto.eu/wp-content/uploads/2014/01/moto_D5.1.1v1.01.pdf

[5] R. Akl, S. Valentin, G. Wunder, and S. Stanczak, "Compensating for CQI Aging By Channel Prediction: The LTE Downlink," in *Proc. of IEEE GLOBECOM'12*, 2012, pp. 4821–4827.

[6] A. Balasubramanian, R. Mahajan, and A. Venkataramani. Augmenting mobile 3G using WiFi. In Proceedings of the 8th international conference on Mobile systems, applications, and services, pages 209-222. ACM, 2010

[7] A. Ben Nacef, N. Montavont, A generic end-host mechanism for path selection and flow distribution, in: IEEE 19th International Symposium on Personal, Indoor and Mobile Radio Communications PIMRC, 2008, pp. 1–5

[8] Elisabetta Biondi, Chiara Boldrini, Andrea Passarella, and Marco Conti, "Optimal duty cycling in mobile opportunistic networks with end-to-end delay guarantees", *European Wireless*, Barcelona, Spain, 14-16 May 2014

[9] Elisabetta Biondi, Chiara Boldrini, Marco Conti, Andrea Passarella, "Duty Cycling in Opportunistic Networks: the Effect on Intercontact Times", *The 17th ACM International Conference on Modeling, Analysis and Simulation of Wireless and Mobile Systems (ACM MSWiM 2014)*, Montreal, Canada, September 21-26 2014.

[10] Y. Blankenship, P. Sartori, B. Classon, V. Desai, and K. Baum, "Link error prediction methods for multicarrier systems," in *Proc. of IEE VTC- Fall'04*, vol. 6, 2004, pp. 4175–4179.

[11] Chiara Boldrini, Marco Conti, Andrea Passarella, "The stability region of the delay in Pareto opportunistic networks", *IEEE Transactions on Mobile Computing*, to appear, available online at http://dx.doi.org/10.1109/TMC.2014.2316506

[12] C. Boldrini, M. Conti, and A. Passarella, "Performance modelling of opportunistic forwarding under heterogenous mobility," *Computer Communications*, pp. 1–17, 2014.

[13] K. Brueninghaus, D. Astely, T. Salzer, S. Visuri, A. Alexiou, S. Karger, and G.-A. Seraji, "Link performance models for system level simulations of broadband radio access systems," in Proc. of IEEE PIMRC'05, 2005.

[14] R. Bruno, A. Masaracchia, A. Passarella, "Robust Adaptive Modulation and Coding (AMC) Selection in LTE Systems using Reinforcement Learning", Proc. of IEEE VTC2014-Fall, 14–17 September 2014, Vancouver, Canada.

[15] Raffaele Bruno, Antonino Masaracchia, and Andrea Passarella, "Offloading through Opportunistic Networks with Dynamic Content Requests", The IEEE Workshop on CellulAR Traffic Offloading to Opportunistic Networks (IEEE CARTOON 2014), Philadelphia, Pennsylvania, USA, October 27, 2014

[16] Raffaele Bruno, Antonino Masaracchia, and Andrea Passarella, "Analysis of MAC-layer Throughput in LTE Systems with Link Rate Adaptation and HARQ Protocols", IIT-CNR Technical Report, 2014.

[17] V. Buenestado, J. Ruiz-Aviles, M. Toril, S. Luna-Ramirez, and A. Mendo, "Analysis of Throughput Performance Statistics for Bench- marking LTE Networks," IEEE Communications Letters, vol. 18, no. 9, pp. 1607–1610, September 2014.

[18] H. Cai and D. Eun, "Crossing over the bounded domain: From exponential to power-law intermeeting time in mobile ad hoc networks," *IEEE/ACM Trans. on Netw.*, vol. 17, no. 5, pp. 1578– 1591, 2009.

[19] F. Capozzi, G. Piro, L. Grieco, G. Boggia, and P. Camarda, "Downlink Packet Scheduling in LTE Cellular Networks: Key Design Issues and a Survey," IEEE Communications Surveys & Tutorials, vol. 15, no. 2, pp. 678–700, 2013.

[20] A. Chaintreau, P. Hui, J. Crowcroft, C. Diot, R. Gass, and J. Scott, "Impact of human mobility on opportunistic forwarding algo- rithms," *IEEE Trans. Mobile Comput.*, pp. 606–620, 2007.

[21] V. Conan, J. Leguay, and T. Friedman, "Characterizing pair- wise inter-contact patterns in delay tolerant networks," in *Au- tonomics'07*, 2007.

[22] COST Action 231, "Digital mobile radio future generation systems," Final Report - EUR 18957, 1999. "ns-3 Model Library – Release ns-3.17", May 14, 2013.

[23] S. Donthi and N. Mehta, "An Accurate Model for EESM and its Application to Analysis of CQI Feedback Schemes and Scheduling in LTE," IEEE Transactions on Wireless Communications, vol. 10, no. 10, pp. 3436–3448, October 2011.

[24] Dutta, A., et al.: Seamless Handover across Heterogeneous Networks - An IEEE802.21 Centric Approach. In: IEEE WPMC (2006).

[25] A. Elnashar and M. El-Saidny, "Looking at LTE in Practice: A Perfor- mance Analysis of the LTE System Based on Field Test Results," IEEE Vehicular Technology Magazine, vol. 8, no. 3, pp. 81–92, September 2013.

[26] J. Francis and N. Mehta, "EESM-Based Link Adaptation in Point-to- Point and Multi-Cell OFDM Systems: Modeling and Analysis," *IEEE Transactions on Wireless Communications*, vol. 13, no. 1, pp. 407–417, January 2014.

[27] S. Gamboa A. Pelov, P. Maillé X. Lagrange, N. Montavont, "*Energy efficient cellular networks in the presence of delay tolerant users*" Global Communications Conference (GLOBECOM), 2013 IEEE, Dec. 2013

[28] Weisi Guo, Siyi Wang , T O'Farrell,S. Fletcher, "*Energy Consumption of 4G Cellular Networks: A London Case Study*", Vehicular Technology Conference (VTC Spring), 2013 IEEE 77[th], June 2013

[29] Z. Haas and T. Small, "A new networking model for biological applications of ad hoc sensor networks," *IEEE/ACM Trans. on Netw.*, vol. 14, no. 1, pp. 27–40, 2006.

[30] Hadiji, F.; Zarai, F.; Kamoun, A., "Architecture of mobile node in heterogeneous networks," Communications and Information Technology (ICCIT), 2012 International Conference on , vol., no., pp.260,264, 26-28 June 2012

[31] Z. He and F. Zhao, "Performance of HARQ With AMC Schemes In LTE Downlink," in Proc. of IEEE CMC'10, 2010, pp. 250–254.

[32] M.A. Imran, E. Katranaras, "Energy Efficiency analysis of the reference system, areas of improvements and targets breakdown," Tech.Rep., Dec.2010

[33] T. Karagiannis, J.-Y. Le Boudec, and M. Vojnovic and, "Power law and exponential decay of intercontact times between mobile devices," *IEEE Trans. Mobile Comput.*, vol. 9, no. 10, pp. 1377 –1390, 2010.

[34] G. Koundourakis, D. Axiotis, M. Theologou, Network-based access selection in composite radio environments, in: IEEE Wireless Communications and Networking Conference, 2007. WCNC'2007, 11–15 March 2007, pp. 3877–3883

[35] A. Kuhne and A. Klein, "Throughput analysis of multi-user ofdma- systems using imperfect cqi feedback and diversity techniques," *IEEE Journal on Selected Areas in Communications*, vol. 26, no. 8, pp. 1440– 1450, October 2008.

[36] J. Ikuno, S. Pendl, M. Simko, and M. Rupp, "Accurate SINR estimation model for system level simulation of LTE networks," in *Proc. of IEEE ICC'12*, 2012, pp. 1471–1475.

[37] C. Lee, "Heterogeneity in contact dynamics: helpful or harmful to forwarding algorithms in DTNs?" in *WiOPT'09*. IEEE, 2009, pp. 1–10.

[38] J. Leinonen, J. Hamalainen, and M. Juntti, "Capacity Analysis of Down- link MIMO-OFDMA Resource Allocation with Limited Feedback," IEEE Transactions on Communications, vol. 61, no. 1, pp. 120–130, January 2013.

[39] Z. Lin, P. Xiao, and B. Vucetic, "SINR distribution for LTE downlink multiuser MIMO systems," in Proc. of IEEE ICASSP'09, April 2009, pp. 2833–2836.

[40] X. Liu, V.O.K. Li, P. Zhang, Joint radio resource management through vertical handoffs in 4G networks, in: IEEE Global Telecommunications Conference, 2006. GLOBECOM'06, November 2006, pp. 1–5

[41] M. Logothetis, K. Tsagkaris, P. Demestichas, Application and mobility aware integration of opportunistic networks with wireless infrastructures, Computers & Electrical Engineering, Available online 24 August 2012, ISSN 0045-7906, 10.1016/j.compeleceng.2012.07.014

[42] Z. Markovic, (2010). Modification of TOPSIS method for solving of multi-criteria tasks. YugoslavJournal of Operations Research. 20, p.117-143

[43] M. Mezzavilla, M. Miozzo, M. Rossi, N. Baldo, and M. Zorzi, "A Lightweight and Accurate Link Abstraction Model for the Simulation of LTE Networks in Ns-3," in Proc. of ACM MSWiM '12, 2012, pp. 55–60.

[44] M. Ni, X. Xu, and R. Mathar, "A channel feedback model with robust SINR prediction for LTE systems," in *Proc. of EuCAP'13*, 2013, pp. 1866–1870.

[45] D. Niyato, E. Hossain, A Cooperative game framework for bandwidth allocation in 4G heterogeneous wireless networks, in: IEEE International Conference on Communications, 2006. ICC'06, vol. 9, June 2006, pp. 4357–4362

[46] O. Ormond, P. Perry, J. Murphy, Network selection decision in wireless heterogeneous networks, in: IEEE 16th International Symposium on Personal, Indoor and Mobile Radio Communications, 2005. PIMRC 2005, vol. 4, 11–14 September 2005, pp. 2680–2684

[47] A. Passarella and M. Conti, "Analysis of individual pair and aggregate inter-contact times in heterogeneous opportunistic net- works," *IEEE Trans. Mobile Comput.*, vol. 12, no. 12, pp. 2483 – 2495, 2013.

[48] K. Piamrat, A. Ksentini, J.Bonnin, C. Viho, Radio resource management in emerging heterogeneous wireless networks, Computer Communications, Volume 34, Issue 9, 15 June 2011, Pages 1066-1076, ISSN 0140-3664, 10.1016/j.comcom.2010.02.015.

[49] A. Picu, T. Spyropoulos, and T. Hossmann, "An analysis of the information spreading delay in heterogeneous mobility dtns," in *IEEE WoWMoM*, 2012, pp. 1–10.

[50] Filippo Rebecchi, Marcelo Dias de Amorim, Vania Conan: DROid: Adapting to individual mobility pays off in mobile data offloading**.**Networking 2014: 1-9

[51] Filippo Rebecchi, Marcelo Dias de Amorim, and Vania Conan. 2014. Flooding data in a cell: is cellular multicast better than device-to-device communications?. In *Proceedings of the 9th ACM MobiCom workshop on Challenged networks* (CHANTS '14).

[52] M. Rinne and O. Tirkkonen, "LTE, the radio technology path towards 4G," Computer Communications, vol. 33, pp. 1894–1906, 2010.

[53] H. Seo and B. Lee, "Proportional-fair power allocation with CDF- based scheduling for fair and efficient multiuser OFDM systems," IEEE Transactions on Wireless Communications, vol. 5, no. 5, pp. 978–983, May 2006.

[54] P. Serra, A. Rodrigues, "*Picocell positioning in an LTE network*" Anancom, Nov 2013

[55] M. Sharif and B. Hassibi, "On the capacity of MIMO broadcast chan- nels with partial side information," IEEE Transactions on Information Theory, vol. 51, no. 2, pp. 506–522, February 2005.

[56] H. Song, R. Kwan, and J. Zhang, "General results on SNR statistics involving EESM-based frequency selective feedbacks," IEEE Transactions on Wireless Communications, vol. 9, no. 5, pp. 1790–1798, May 2010.

[57] Q. Song, A. Jamalipour, A network selection mechanism for next generation networks, in: IEEE International Conference on Communications, 2005. ICC'2005, vol. 2, 16–20 May 2005, pp. 1418–1422

[58] T. Spyropoulos, K. Psounis, and C. Raghavendra, "Efficient routing in intermittently connected mobile networks: The single copy case," *IEEE/ACM Trans. on Netw.*, vol. 16, no. 1, pp. 63–76, 2008.

[59] T. Spyropoulos, K. Psounis, and C. Raghavendra, "Efficient routing in intermittently connected mobile networks: The multiple- copy case," *IEEE/ACM Trans. on Netw.*, vol. 16, no. 1, pp. 77–90, 2008.

[60] T. Spyropoulos, T. Turletti, and K. Obraczka, "Routing in Delay- Tolerant Networks Comprising Heterogeneous Node Populations," *IEEE Trans. Mobile Comput.*, pp. 1132–1147, 2009.

[61] E. Stevens-Navarro and V. W. S. Wong, "Comparison between vertical handoff decision algorithms for heterogeneous wireless networks," in Vehicular Technology Conference, 2006. VTC 2006-Spring. IEEE 63rd, vol. 2, Melbourne, Vic., May 2006, pp. 947

[62] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, March 1998.

[63] A. Taha, H. Hassanein, H. Mouftah, On robust allocation policies in wireless heterogeneous networks, in: First International Conference on Quality of Service in Heterogeneous Wired/Wireless Networks, 2004. QSHINE 2004, 18–20 October 2004, pp. 198–205.

[64] N. Varanese, J. Vicario, and U. Spagnolini, "On the Asymptotic Throughput of OFDMA Systems with Best-M CQI Feedback," IEEE Wireless Communications Letters, vol. 1, no. 3, pp. 145–148, June 2012.

[65] A. Wilson, A. Lenaghan, R. Malyan, Optimising wireless access network selection to maintain QoS in heterogeneous wireless environments, Wireless Personal Multimedia Communications 2005. WPMC'05, 18–22 September 2005

[66] X. Yang, J. Bigham, L. Cuthbert, Resource management for service providers in heterogeneous wireless networks, in: IEEE Wireless Communications and Networking Conference 2005, vol. 3, 13–17 March 2005, pp. 1305–1310.

[67] "Key Performance Indicators (KPI) for Evolved Universal Terrestrial Radio Access Network (E UTRAN): Definitions," TS 32.450, Version 9.1.0 Release 9, June 2010.

[68] 3GPP, "Conveying MCS and TB size via PDCCH," R1-081483, Septem- ber 2010.

[69] 3GPP, "Evolved Universal Terrestrial Radio Access (E-UTRA): Physical layer procedures (Release 9)," TS 36.213 V9.3.0, September 2010.

[70] IEEE 802.16 Broadband Wireless Access Working Group, "Evaluation Methodology for P802.16m-Advanced Air Interface," IEEE 802.16m- 08/004r2, 2008.

[71] WiMAX Forum, "WiMAX System Evaluation Methodology," V2.1, July 2008.

[72] "ns-3 Model Library – Release ns-3.17", May 14, 2013.

[73] "ns-3 Model Library – Release ns-3.18", August 29, 2013

# Appendix A

This appendix contains the reprint of the following papers

- Chiara Boldrini, Marco Conti, Andrea Passarella, "The stability region of the delay in Pareto opportunistic networks", *IEEE Transactions on Mobile Computing*, to appear, available online at http://dx.doi.org/10.1109/TMC.2014.2316506

- Elisabetta Biondi, Chiara Boldrini, Andrea Passarella, and Marco Conti, "Optimal duty cycling in mobile opportunistic networks with end-to-end delay guarantees", *European Wireless*, Barcelona, Spain, 14-16 May 2014

- R. Bruno, A. Masaracchia, A. Passarella, "Robust Adaptive Modulation and Coding (AMC) Selection in LTE Systems using Reinforcement Learning", Proc. of IEEE VTC2014-Fall, 14–17 September 2014, Vancouver, Canada.

- Raffaele Bruno, Antonino Masaracchia, and Andrea Passarella, "Offloading through Opportunistic Networks with Dynamic Content Requests", The IEEE Workshop on CellulAR Traffic Offloading to Opportunistic Networks (IEEE CARTOON 2014), Philadelphia, Pennsylvania, USA, October 27, 2014

- Filippo Rebecchi, Marcelo Dias de Amorim, and Vania Conan. 2014. Flooding data in a cell: is cellular multicast better than device-to-device communications?. In *Proceedings of the 9th ACM MobiCom workshop on Challenged networks* (CHANTS '14).

# The stability region of the delay in Pareto opportunistic networks

Chiara Boldrini, Marco Conti, *Member, IEEE*, and Andrea Passarella, *Member, IEEE*

**Abstract**—The intermeeting time, i.e., the time between two consecutive contacts between a pair of nodes, plays a fundamental role in the delay of messages in opportunistic networks. A desirable property of message delay is that its expectation is finite, so that the performance of the system can be predicted. Unfortunately, when intermeeting times feature a Pareto distribution, this property does not always hold. In this paper, assuming heterogeneous mobility and Pareto intermeeting times, we provide a detailed analysis of the conditions for the expectation of message delay to be finite (i.e., to converge) when social-oblivious or social-aware forwarding schemes are used. More specifically, we consider different classes of social-oblivious and social-aware schemes, based on the number of hops allowed and the number of copies generated. Our main finding is that, in terms of convergence, allowing more than two hops may provide advantages only in the social-aware case. At the same time, we show that using a multi-copy scheme can in general improve the convergence of the expected delay. We also compare social-oblivious and social-aware strategies from the convergence standpoint and we prove that, depending on the mobility scenario considered, social-aware schemes may achieve convergence while social-oblivious cannot, and vice versa. Finally, we apply the derived convergence conditions to three popular contact datasets available in the literature (Cambridge, Infocom, and RollerNet), assessing the convergence of each class of forwarding protocols in these three cases.

**Index Terms**—Opportunistic networks, DTN, Routing protocols, Performance models, Delay convergence

✦

## 1 INTRODUCTION

THE great popularity of the delay tolerant networking paradigm is due to its ability to cope with challenged network conditions, such as high node mobility, variable connectivity, and disconnected subnetworks, that would impair communications in traditional Mobile Ad Hoc Networks. Opportunistic networks are an instance of the delay tolerant paradigm applied to networks made up of users' portable devices (such as smartphones and tablets). In this scenario, user mobility becomes one of the main drivers to enable message delivery. In fact, according to the store-carry-and-forward paradigm, user devices store messages and carry them around while they move in the network, exchanging them upon encounter with other nodes, and eventually delivering them to their destination.

An opportunistic forwarding protocol defines the strategy according to which messages are exchanged during encounters. Two main approaches can be identified: *social-oblivious* protocols, which do not exploit any information about the users' context and social behaviour, but just hand over the message to the first node encountered, and *social-aware*[1] protocols, which make use of information on how users behave or which

social relations they share in order to make predictions on users' future behavior that might be useful for forwarding messages. Depending on the number of copies generated for the same message, forwarding protocols can be further classified into single-copy or multi-copy schemes. Forwarding protocols may also differ in the number of intermediate hops that they exploit. Simpler strategies may be single-hop or two-hop strategies (e.g. Direct Transmission and Two Hop [2]), while others can allow multi-hop paths to bring the message to the destination.

Modelling the performance of social-oblivious and social-aware forwarding protocols for opportunistic networks is still an open research issue. Knowing the distribution of intermeeting times and the rules applied by the forwarding algorithm used in the network, one could - in principle - model the distribution of the delay experienced by messages and compute its expectation. In practice, modeling analytically the delay of the various forwarding protocols for general distributions of inter meeting times is very hard, and models exist only for some specific cases, typically assuming exponential intermeeting times [1] [3] [4] [5] [6] [7]. A related modeling challenge is to assess the convergence of routing protocols, i.e. whether a specific protocol yields finite or infinite expected delay. Assessing convergence allows us to understand whether a particular protocol can be safely used or not given a pattern of intermeeting times and how to configure it so that it converges, if possible. Although less informative than a complete delay model, convergence models can be derived for a large class of routing protocols releasing the exponential intermeeting time assumption, as shown in this paper.

• *C.Boldrini, M.Conti, and A.Passarella are with the Institute for Informatics and Telematics of the Italian National Research Council, Via G. Moruzzi 1, 56124 Pisa, Italy.*
*E-mail: chiara.boldrini, marco.conti, andrea.passarella@iit.cnr.it*

1. These policies are also referred to as utility-based [1], in contrast to randomized strategies, corresponding to our social-oblivious schemes.

The convergence of the expected delay is not guaranteed in all cases in which the expectation of the inter-meeting times may diverge. In fact, being the delay the result of the composition of the time intervals between node encounters, depending on the convergence of inter-meeting times, the expectation of the delay itself might diverge. This can happen, for example, when intermeeting times feature a Pareto (also known as power law) distribution, as first highlighted in [8]. The problem with Pareto distributions is that their expectation is finite only for certain values of their exponent $\alpha$. More specifically, the expectation is finite if $\alpha > 1$, while for $\alpha \leq 1$ it diverges to infinity. The first to postulate the existence of Pareto intermeeting times in real mobility scenarios (i.e., analyzing real traces of human mobility) were Chaintreau et al. in their seminal work in [8]. The relevance of Pareto intermeeting times in opportunistic networks is both theoretical and empirical. Cai and Eun [9] have mathematically derived that heavy-tailed intermeeting times can emerge depending on the relationship between the size of the boundary of the considered scenario and the relevant timescale of the network, showing that, at least in principle, Pareto intermeeting times are something that one may be faced with when studying opportunistic networks. Empirical evidence for the presence of Pareto intermeeting times was first suggested by [8], but it has been later criticised, arguing that the tail of the distribution is in fact exponential (e.g., [10]). Typically, these results are derived focusing on the aggregate inter-contact time distribution, while convergence depends on pairwise distributions. As proved in [11], the aggregate and pairwise distributions can be in general very different, and therefore analysis of pairwise inter-contact times are necessary, which are however mostly missing in the literature. To address this issue, we have performed a pairwise hypothesis testing on three popular publicly available contact datasets (Cambridge, Infocom'05, and RollerNet – see Section 8) and we have found that the Pareto hypothesis for intermeeting times cannot be rejected for $80\%$, $97\%$, and $85.5\%$ of pairs, respectively. We believe that these results provide a strong case for Pareto intermeeting times in opportunistic networks and substantially motivate analyses like the one presented in this paper.

Under the Pareto intermeeting times assumption, in this work we derive the stability region (i.e., the Pareto exponent values of pairwise intermeeting times for which finite expected delay is achieved) of a broad class of social-oblivious and social-aware forwarding protocols (single- and multi-copy, single- and multi-hop). The starting point of our paper is the work by Chaintreau et al. [8], where such conditions have been studied for the two-hop scheme (see Section 2 for more details) under the assumption of homogeneous mobility (i.e., i.i.d. intermeeting times across all pairs). However, measurement studies [12] [8] have shown that real networks are intrinsically heterogeneous. Thus, in this paper, we investigate whether heterogeneity in contact

patterns helps the convergence of the expected delay of a general class of social-oblivious (Section 5) and social-aware (Section 6) forwarding protocols, and whether convergence conditions can be improved using multi-copy strategies and/or multi-hop paths. In general, we find that there is no protocol or family of protocols that always outperform the others (Section 7). More specifically, the key findings presented in the paper are the following:

- For *social-oblivious strategies*, if convergence can be achieved, *two hops are enough* for achieving it.
- Using $n$ *hops* can help *social-aware schemes*, and make them converge in some cases when all other social-aware or social-oblivious schemes diverge.
- In both the social-oblivious and the social-aware case, we find that *multi-copy strategies* can achieve a finite expected delay even when single-copy strategies cannot.
- Comparing *social-oblivious and social-aware multi-copy solutions*, we are able to prove mathematically that there is no clear winner between the two, since either one can achieve convergence when the other one fails, depending on the underlying mobility scenario.

In addition to these results, we discuss the related work in Section 2 and the network model we refer to in Section 3. Our reference forwarding policies are described in Section 4. Finally, in Section 8 we showcase the main results of this work by applying the derived convergence conditions to three popular contact datasets available in the literature.

## 2 RELATED WORK

This work is orthogonal to the literature on models of delay in opportunistic networks, since we provide the conditions for the *existence* of a finite delay. Once convergence has been verified, the expected delay value can be computed using complete delay models like the one in [13]. End-to-end delay models are typically much more difficult to obtain than convergence models, and therefore convergence models are a very useful first step in the analytical characterisation of forwarding protocols in opportunistic networks. Most existing models assume that intermeeting times are approximately exponentially distributed [5] [6] [14], and in these cases convergence is never an issue. However, when this assumption does not hold, convergence becomes a critical evaluation aspect, and should be studied preliminarily to any additional analysis of the exact value of the expected delay. To the best of our knowledge, there is no other contribution, besides that of Chaintreau et al. [8], that considers the problem of the convergence of the expected delay when intermeeting times feature a Pareto distribution. Our work differs from that of Chaintreau et al. both in the mobility settings and in the forwarding schemes considered. More specifically, we focus on the more general case of heterogeneous intermeeting times (as opposed to

the homogeneous mobility considered in [8]), we extend the set of social-oblivious policies considered and we add the social-aware case.

Forwarding protocols for opportunistic networks can be classified as social-oblivious or social-aware protocols, depending on whether they use information on the way nodes behave in order to make forwarding decisions. In this paper we abstract the detailed mechanisms of both classes of protocols, in order to study their convergence properties, as discussed in Section 4. The simplest social-oblivious protocol is Direct Transmission [2], in which the source node is only allowed to deliver the message directly to the destination, if ever encountered. At the opposite side of the spectrum, with Epidemic routing [15] a new copy of the message is generated and handed over (both by the source and intermediate relays) any time a new node is encountered. In order to mitigate the side effects of Epidemic-style forwarding schemes in resource constrained environments, controlled flooding solutions have been proposed (e.g., Spray&Wait [3], gossiping [7]). Another popular social-oblivious forwarding protocol is the Two Hop scheme [2], in which a message is forwarded by the source node to the first node encountered, which is then allowed only to pass the message directly to the destination. The Two Hop strategy has been shown to guarantee the maximum capacity in a homogeneous network [2].

Social-aware strategies can have different levels of awareness. Simplest approaches exploit information such as time since the last encounter (Spray&Focus [3]) or frequency of encounters (PROPHET [16]). This information is used to predict future meetings between pairs of nodes and thus to select relays that can guarantee a quick delivery according to the heuristic in use. In more complex strategies, the centrality of nodes in the social graph connecting the users of the network is used as an indicator of the ability to deliver messages (see, e.g., BUBBLE [17], SimBet [18]). Alternatively, as in the case of HiBOp [19] and SocialCast [20], the fitness of a node as a forwarder is computed from information on the context the users live in, e.g., information on the people they meet, the friends they have, the places they visit.

This paper extends our previous work in [21], which was only focused on social-oblivious forwarding strategies. In this work, besides extending the convergence conditions for the $m$-copy 2-hop case that we derived in [21], we include the analysis of social-aware forwarding strategies and a detailed comparison between social-aware and social-oblivious strategies from the convergence standpoint.

## 3 Preliminaries

Our model considers a network with $N$ mobile nodes. For the sake of simplicity, we hereafter assume that messages can be exchanged only at the beginning of a contact between a pair of nodes and that the transmission of the relayed messages can be always completed within the duration of a contact. In addition, we assume that each message is a bundle [22], an atomic unit that cannot be fragmented. We also assume infinite buffer space on nodes. All the above assumptions allow us to isolate, and thus focus on, the effects of node mobility from other effects, and are common assumptions in the literature on opportunistic networks modelling (they are used in most of the literature reviewed in Section 2). In addition, for the sake of comparison with [8], we also assume that the probability that two nodes meet is greater than zero for all node pairs. This ensures that, in principle, all nodes can meet with each other. Therefore, cases of deadlock (a message reaches a node which is impossible to leave due to the total absence of contacts with either other possible relays or the destination) are not possible. The only cause of divergent expected delay are the distributions of intermeeting times.

As we neglect transmission time, the actual duration of the contact is not critical. Thus, the main role in the experienced delay is played by intermeeting times, which are defined as the time between two consecutive meetings between the same pair of nodes[2]. We denote with $M_{ij}$ the intermeeting times between nodes $i$ and $j$. For the sake of tractability, we assume that the network is stationary and that intermeeting times for a specific node pair $i, j$ are i.i.d.. Under these assumptions, the encounter process between two nodes $i$ and $j$ can be seen as a renewal process with renewal intervals distributed as $M_{ij}$.

The message generation process and the mobility process are assumed to be independent. Thus, the time at which a new message is generated can be treated as a random time in the evolution of the mobility process, and thus the message sees the network as an observer arriving at a random point in time would. For this reason, in our analysis we will often use the concept of residual intermeeting time, defined as the time two nodes that are not in contact at a random time $t_0$ have to wait before meeting again. We denote the residual intermeeting time for the $i, j$ node pair as $R_{ij}$.

Under our assumption of Pareto intermeeting times, the intermeeting time $M_{ij}$ between a generic pair of nodes $i$ and $j$ is described by the CCDF $F_{M_{ij}}(t) = \left(\frac{t_{min_{ij}}}{t+t_{min_{ij}}}\right)^{\alpha_{ij}}$, in which we use the definition of the Pareto distribution that allows for values arbitrarily close to zero, usually denoted as American Pareto [23] or Pareto distribution of the second kind [24][3]. Parameters $\alpha_{ij}(> 0)$ and $t_{min_{ij}}$ are the shape and scale of the Pareto distribution and, similarly to the reference literature [8][10], in the following we restrict to the case of power law random variables having the same

---

2. Without loss of generality, here we assume a deterministic unit disk graph model for radio propagation, which is a common assumption in the literature on opportunistic networks. The proposed framework still applies for every other model of radio propagation.

3. The stability region derived in this paper holds also for the other version of the Pareto distribution, usually denoted as European Pareto, as discussed in [25]. Content in [25] not included in this paper is provided as supplemental material.

scale, i.e., $t_{min_{ij}} = t_{min}, \forall i, j$. We will use the following properties of Pareto intermeeting times throughout the paper (please refer to Appendix A in [25] for more details):

P1  $E[M_{ij}]$ converges (i.e., is finite) if and only if $\alpha_{ij} > 1$.

P2  The residual intermeeting times $R_{ij}$ associated to $M_{ij}$ feature an American Pareto distribution with shape $\alpha_{ij} - 1$ and scale $t_{min}$ [23], hence their expectation converges if and only if $\alpha_{ij} > 2$.

P3  $\min_j\{R_{ij}\} \sim \text{Pareto}\left(\sum_j(\alpha_{ij}-1), t_{min}\right)$ and its expectation converges if and only if $\sum_j(\alpha_{ij}-1) > 1$.

P4  Assume $R_{ij}^{t_i}$ denotes the residual conditioned to the fact that $i$ and $j$ met at time $t_i$. Then the expectation of $R_{ij}^{t_i}$ converges if and only if the expectation of $R_{ij}$ converges.

Furthermore, in the mathematical analysis in Sections 5 and 6, we will also heavily rely on the result in Lemma 1 below. The intuition behind it is the following. In the general case, the time before a node $i$ currently holding a copy of the message hands it over to another node $j$ depends on whether nodes $i$ and $j$ met in time interval $(t_0, t_i)$, where $t_0$ is the message generation time. In fact, meetings correspond to renewals in the encounter renewal process between $i$ and $j$, hence, from the meeting time on, we should consider the intermeeting time and not its residual. However, Lemma 1 below tells us that, when intermeeting times feature a Pareto distribution, we can simply study the case in which nodes $i$ and $j$ did not meet in $(t_0, t_i)$ (i.e., model the time to the next encounter as a residual time), thus simplifying our analysis (for more details, see the proof of Lemma 1 in [25]).

*Lemma 1 (Worst-case waiting time):* Assume that node $i$ has received a copy of the message at time $t_i$. In the worst case (happening with a non negligible probability), the time before node $i$ hands over the message to another node $j$ can be modeled as $R_{ij}^{t_i - t_0}$ (i.e., $R_{ij}$ conditioned to be greater than $t_i - t_0$) or, equivalently from a convergence standpoint, as $R_{ij}$.

## 4 FORWARDING STRATEGIES

In this section we summarise the main variants of opportunistic forwarding schemes that will be later evaluated against each other as far as the convergence of their expected delay is concerned. We identify two main strategies that forwarding protocols can adopt in order to improve their forwarding performance, namely the number of copies generated and the number of hops allowed. As we show later on in the section, it is easy to place any of the most popular routing protocols proposed in the literature in this classification.

As far as multi-copy strategies are concerned, here we only allow the source node to create and hand over multiple copies. Other possible configurations (e.g.,

intermediate relays allowed to generate new copies, like in the Spray&Wait case [3]) are left as future work. With respect to the number of hops $n$, we assume it to be either limited arbitrarily (e.g., using the TTL field) or naturally constrained by the forwarding strategy (e.g., social-aware schemes can exploit a number of intermediate relays that is at most equal to the number of nodes that are better forwarders - according to some social-aware metric - than the source node). In all cases, the last relay can only deliver the message to the destination directly. Please note that in this paper we only consider the case in which both the source node and intermediate relays refuse the custody of copies that they have already relayed (i.e., we assume that nodes are *memoryful*)[4]. For this to be feasible, we assume that the identity of previous relays is enclosed into the copy's header. In the case of multiple copies, we assume that the source node does not use the same relays multiple times, and that relays do not accept the custody of the same copy of the message more than once. They can be used, however, as relays for different copies of the same message (as avoiding this would need to keep track of all forwarded messages at each relay, which would make protocols not scalable).

Due to the variety of social-aware schemes available in the literature (see Section 2) and the limited space, here we only consider an abstract social-aware protocol that measures how good a relay is for a given destination in terms of its *fitness*. The fitness $fit_i^d$ is assumed to be a function of how often node $i$ meets the destination $d$, thus $fit_i^d$ can be taken as proportional to the rate of encounter $\frac{1}{E[M_{id}]}$ between node $i$ and the destination. Under this abstract and general social-aware strategy, upon encounter, a node $i$ can hand over the message to another node $j$ only if its fitness is lower than the fitness of the peer, i.e., if $fit_j^d > fit_i^d$ holds (in the following we drop superscript $d$). The fitness function considered here uses only information on contacts between nodes, which have a direct dependence on the intermeeting time distribution. This lets us clearly show what is the impact of the contact dynamics on the performance of opportunistic forwarding protocols. How such simple fitness function can be extended to more complex forwarding strategies has been discussed in [13].

The combinations of the forwarding characteristics described above can be found in well known routing strategies. For example, the 1-hop 1-copy forwarding corresponds to Direct Transmission [2], while the 2-hop 1-copy forwarding is equivalent to the Two Hop forwarding introduced in [2]. The 2-hop $m$-copy forwarding is equivalent to the multi-copy version of the Two Hop protocol studied in [8]. Note that for most of the social-aware protocols, the number of copies and the maximum hops are also defined as parameters of the algorithm.

4. In [21] we have derived the convergence conditions for the memoryless version of the class of social-oblivious forwarding protocols considered here, showing that the absence of memory always penalizes the convergence.

# 5 EXPECTED DELAY CONVERGENCE FOR SOCIAL-OBLIVIOUS SCHEMES

In this section we study under which conditions the expected delay of the social-oblivious schemes described in Section 4 converges for a *tagged* source-destination pair. We denote it as $E[D_{sd}]$, where $s$ and $d$ are the source and destination node, respectively. Simultaneous convergence for all source-destination pairs simply requires combining the conditions derived in the paper.

Recall that according to social-oblivious forwarding a message is handed over to the first feasible relay encountered. In the following, we denote with $\mathcal{P}_i$ the set of all nodes that can be encountered by node $i$. For the sake of comparison with [8], we assume that the probability of an encounter between any pair of nodes is strictly greater than zero (hence, we have that $|\mathcal{P}_i| = N - 1$ for all nodes $i$) and that $\alpha_{ij} > 1$ for all $i, j$ node pairs.

## 5.1 Single-copy schemes

Theorem 1 below focuses on the 1-copy 1-hop social-oblivious scheme, which corresponds to the popular Direct Transmission scheme. In the following, we omit the proof since this result follows directly from property P2 and we move to the analysis of the 1-copy 2-hop scheme immediately.

*Theorem 1 (1-copy 1-hop scheme):* $E[D_{sd}]$ converges if and only if $\alpha_{sd} > 2$.

*Theorem 2 (1-copy 2-hop scheme):* $E[D_{sd}]$ converges if and only if both the following conditions hold true:

C1 $\quad \sum_{j \in \mathcal{P}_s} \alpha_{sj} > 1 + |\mathcal{P}_s|$
C2 $\quad \alpha_{jd} > 2, \forall j \in \mathcal{P}_s - \{d\}$.

*Proof:* The protocol converges if and only if (iif) both the delay at the first hop converges and the delay at the second hop converges. We analyse the former, first. The delay of the first hop converges iff the time required by the source to hand over the message, which is the time to encounter the first node in set $\mathcal{P}_s$, is finite. The source node $s$ can either deliver the message directly to the destination or hand it over to an intermediate relay. The time before the source node hands over the message is distributed as $\min_{j \in \mathcal{P}_s}\{R_{sj}\}$, which is the time before the first node (possibly including the destination) is encountered. From property P3, we know that $\min_{j \in \mathcal{P}_s}\{R_{sj}\}$ has a finite expectation iff $\sum_{j \in \mathcal{P}_s} \alpha_{sj} > 1 + |\mathcal{P}_s|$, thus obtaining condition C1. We now consider the convergence of the second hop. As social oblivious protocols cannot control which relay is used, the delay of the second hop is finite iff delays from all possible relays to the destinations are finite. If the node to which the message has been handed over is not the destination but another generic node $j$, the expected delay from $j$ to $d$ is finite iff the expectation of the time before $j$ meets $d$ is finite. Exploiting Lemma 1, we can model the time before node $j$ hands over the message to $d$ as $R_{jd}$, whose expectation is finite iff $\alpha_{jd} > 2$. Please note that from

here on, due to lack of space, we will not prove again the necessity of the convergence conditions we derive. In all cases it will be straightforward to prove it using the same argument outlined above. The complete proofs are however available in [25]. □

According to Theorem 1, the Direct Transmission protocol yields a convergent expected delay only if the source node meets the destination with a residual inter-meeting time whose expectation converges. This clearly follows from the fact that the source node cannot exploit any other relays for the forwarding of the message. In the case of the two-hop scheme, the expectation converges even if the source node is not able to ensure convergence with a direct delivery. This can happen if the source node is able to hand over the message to any of the possible relays within a convergent expected time (Condition C1) and if the meeting process between this relay and the destination has a residual whose expectation converges (Condition C2). Please note that condition C1 alleviates the convergence condition on the source node at the expense of the additional condition C2 on intermediate relays.

With Theorem 3 we extend the analysis of single-copy schemes by studying their $n$-hop version.

*Theorem 3 (1-copy $n$-hop scheme):* $E[D_{sd}]$ converges if and only if conditions C1 and C2 in Theorem 2 hold true.

*Proof:* See Appendix C of [25] for a complete proof. The intuitive reason behind Theorem 3 is that, since the first hop (from source to first relay) and last hop (from last relay to destination) are equivalent to those in Theorem 2, they also share the same convergence conditions (C1 and C2). For intermediate hops, it is possible to prove that convergence conditions are looser than C1 and C2, which are then sufficient and necessary for convergence. □

Theorem 3 tells us that, when using single-copy social-oblivious schemes, letting the message traverse more than two hops does not improve the convergence of the expected delay. Thus, when convergence is the only goal, network resources can be saved using the two-hop social-oblivious scheme without impairing the convergence of the expected delay.

## 5.2 Multi-copy schemes

As discussed in Section 2, when multiple copies of the same message can travel in parallel the opportunities to reach the destination are multiplied. In this section we investigate whether this also positively affects the convergence of the expected delay. Please note that hereafter we discuss the complete proof only for the first lemma, which provides the rationale for obtaining the other results of this section (for which just an intuitive proof is sketched, while more details can be found in Appendix C of [25]).

### 5.2.1 Two-hop forwarding

Recall that, according to the multi-copy version of the two-hop forwarding scheme, the source node hands over a copy of the message to the first $m$ encountered nodes, which will then be only allowed to deliver the message directly to the destination, if ever encountered. In Lemmas 2 and 3 we study separately the first hop and the second hop, then putting together their results in Theorem 4. The goal is to derive how many convergent copies the source node can send out at the first hop and how many are needed for having a convergent second hop. In fact, as we demonstrate below, the higher the number of copies on the intermediate relays, the easier the convergence at the second hop. Thus, the number of copies that the source node is able to hand over within a finite expected time is critical to the convergence of the whole path. It is possible to prove that first-hop convergence becomes more difficult as the number of available relays decreases. Hence, after a certain point, the number of relays left does not allow to deliver one more copy within a finite expected time, setting an upper bound on the maximum number of first-hop convergent copies that the source node can send. This number depends also on the order in which relays are used (i.e., on the Pareto exponents of the available relays), which in turn depends on the sequence of encounters at the source node. Clearly, this order cannot be controlled and it is only the result of the evolution of the meeting process. Since the source node can meet at most $N-1$ nodes, the possible sequences of distinct encounters are $(N-1)!$. Let us denote as $\pi_i$ the $i$-th of these permutations. For each possible permutation $\pi_i$, in Lemma 2 we are able to compute the maximum number ($max_i^{so}$) of convergent copies that can be sent at the first hop by the social-oblivious source node. Then, considering all possible permutations $\pi_i$, we can identify (Corollary 1) a range of values (specifically, $[max_{lo}^{so}, max_{up}^{so}]$) within which $max_i^{so}$ can vary, and under which permutations $\pi_i$ the extreme values of the interval are achieved.

*Lemma 2 ($max_i^{so}$):* When intermediate relays are selected by the source node according to sequence $\pi_i$, the source node is able to deliver at most $max_i^{so}$ copies to as many relays with finite first hop expected delay, with $max_i^{so}$ being equal to the following:

$$max_i^{so} = \arg\min_m \{f_{max}^{so}(m, \pi_i) > 0\}, \qquad (1)$$

where $f_{max}^{so}(m, \pi_i) = m + \sum_{z=m}^{|\mathcal{P}_s|} \alpha_z^{(i)} - (2 + |\mathcal{P}_s|)$ and $\alpha_z^{(i)}$ denotes the $\alpha_{sj}$ exponent of the $z$-th node belonging to $\pi_i$.

*Proof:* At the first hop $m$ copies are relayed to the first $m$ distinct encountered nodes. Thus, the delivery process at the first hop is a selection without repetitions: every time a relay is selected, it is removed from the set of future relays for the same message.

Let us define $\mathcal{P}_s^k$ as the set of relays still available to $s$ when the source node is delivering the $k$-th copy, $t_0$ the time at which the message is generated at the source, and $t_k$ the time at which the $k$-th copy is handed over. Given that we assume that the probability that any two nodes meet is greater than zero, we have that $|\mathcal{P}_s| = N-1$ and $|\mathcal{P}_s^k| = N-1-(k-1) = N-k$. Exploiting Lemma 1, the time before the $k$-th copy is relayed is given by $\min_{j \in \mathcal{P}_s^k}\{R_{sj}\}$, which converges (Property P3) as long as $\sum_{j \in \mathcal{P}_s^k}(\alpha_{sj} - 1) > 1$, or equivalently, $\sum_{j \in \mathcal{P}_s^k} \alpha_{sj} > 1 + |\mathcal{P}_s^k|$, with $|\mathcal{P}_s^k| = N - k$. In order to achieve convergence for the $m$ copies, this condition should be satisfied for all $k$ from 1 to $m$.

We start by finding whether convergence is achieved for a fixed $m$. Lemma C1 in [25] tells us that the smaller the cardinality of the set of random variables of which we take the minimum, the tougher the convergence. This implies that the strictest condition for the convergence of the expected delay of the first hop is imposed by the $m$-th copy, i.e., by the one that sees the smallest set of nodes left for relaying. Thus, if we are able to define a convergence condition for the $m$-th copy, then it follows that the finiteness of the expected time to relaying for all previous copies is automatically guaranteed. Let us thus focus on the relaying of the $m$-th copy. When the $(m-1)$-th copy has been delivered, there are $N - 1 - (m-1) = N - m$ potential relays left for the $m$-th copy. The identities of these $N - m$ potential relays depend on the previous evolution of the forwarding process (i.e., which nodes have already been used). More specifically, there can be $(N-1)!$ different[5] permutations of the $N-1$ nodes in $\mathcal{P}_s$, while there can be $\binom{N-1}{N-m}$ possible combinations for the relays in $\mathcal{P}_s^m$. Let us denote with $\pi_i$ the $i$-th of the $(N-1)!$ permutations and with $\nu_i$ its corresponding combination. That is, taken sequence $\pi_i = \{a, e, c, b, f, h, g\}$ of encounters (where $a, b, c, e, f, g, h$ are the nodes that the source node can meet) and assuming $m = 3$ we denote with $\nu_i$ the set $\{c, b, f, h, g\}$, i.e, the set of nodes available as relays once the first and second copies have been handed over. Let us now define a mapping $g^{(i)}$ that goes from set $\{\alpha_{sj}\}_{j \in \mathcal{P}_s}$ to set $\{\alpha_z^{(i)}\}_{z \in \{1, \dots, |\mathcal{P}_s|\}}$, where $\alpha_z^{(i)}$ corresponds to the exponent $\alpha_{sj}$ of the $z$-th element in $\pi_i$. Using the above notation, the time before the $m$-th copy is handed over is described by $\min_{j \in \nu_i} R_{sj}$. Using again property P3 and the mapping defined above, we have that the convergence condition for the expected delay of the $m$-th copy is given by the following:

$$\sum_{z=m}^{|\mathcal{P}_s|} \alpha_z^{(i)} + m - (N+1) > 0. \qquad (2)$$

As discussed before, since the $m$-th copy experiences the worst conditions for convergence, guaranteeing convergence for the $m$-th copy implies automatic convergence of all previous copies. Hence, Equation 2 characterizes the stability region for first hop convergence.

---

5. Please note that any of these permutations happen with non negligible probability, since we assume that all nodes can meet with each other. A rigorous proof can be obtained exploiting the same argument used in Part 4 of the proof of Theorem 3.

The above equation defines the convergence condition for the $m$-th copy when relays are encountered according to encounter sequence $\pi_i$. For a given node permutation $\pi_i$, we can also compute the greatest $m$ value for which convergence is achieved, and in the following we discuss how. Recall that, according to Lemma C1 in [25], convergence becomes more difficult as $m$ increases. This is highlighted also by Equation 2. In fact, the left-hand side of the equation (hereafter denoted as $f_{max}^{so}(m, \pi_i)$) decreases as $m$ increases (the formal demonstration is at the end of the proof). This implies that either $f_{max}^{so}(m, \pi_i)$ is always above/below zero or $f_{max}^{so}(m, \pi_i)$ crosses the x-axis at a certain point. If $f_{max}^{so}(m, \pi_i)$ is always below zero, the source node is not able to send any copy with finite first hop expected delay. Otherwise, the maximum number of convergent copies (for a given node encounter sequence $\pi_i$) that the source node can send is equal to the greatest $m$ for which $f_{max}^{so}(m, \pi_i)$ is still above zero. Hence, Equation 1 follows.

To conclude the proof, let us now demonstrate that $f_{max}^{so}(m, \pi_i)$ decreases with $m$. To this aim, consider moving from $m$ to $m + 1$. Function $f_{max}^{so}(m + 1, \pi_i)$ can be rewritten as $\sum_{z=m}^{|\mathcal{P}_s|} \alpha_z^{(i)} - \alpha_m^{(i)} + m + 1 - (N + 1)$. Thus, the difference between $f_{max}^{so}(m + 1, \pi_i)$ and $f_{max}^{so}(m, \pi_i)$ is $1 - \alpha_m^{(i)}$. $1 - \alpha_m^{(i)}$ is always smaller than zero, since we have assumed $\alpha_{ij} > 1$ for all $i, j$ node pairs. This implies that the left-hand side of Equation 2 decreases as $m$ increases. □

*Corollary 1:* Quantity $max_i^{so}$ derived in Lemma 2 takes values in the interval $[max_{lo}^{so}, max_{up}^{so}]$. The upper and lower bound on $max_i^{so}$ (corresponding to the best and worst case for convergence) are reached when $\pi_i$ corresponds to nodes encountered in increasing and decreasing order of $\alpha_{sj}$, respectively.

*Proof:* Let us provide an intuitive explanation for this result. We can divide the set $\mathcal{P}_s$ of possible relays at the source node into two disjoint sets, one containing the nodes that have already been used as relays and one containing those that have not. Clearly, as copies are handed over by the source node, nodes move from the second subset to the first subset. Convergence is determined by the exponents associated with nodes in the second subset (nodes still to be encountered). The higher the exponents in this subset, the easier the convergence, and vice versa. When convergence is easier, the source node can send more copies. Conversely, when convergence is more difficult less convergent copies can be sent. Thus, in the best case the exponents associated with nodes in the second subset are the highest among the nodes in $\mathcal{P}_s$, while in the worst case such exponents are the lowest. From this, Corollary 1 follows. □

Let us now focus on the second hop. The sequence $\pi_i$ according to which the source node meets the other nodes affects not only the first hop delay but also the second hop delay. In fact, the relays picked by the source node according to $\pi_i$ are those that are in charge of bringing the message to its final destination. It is possible

to prove (Lemma C1 in [25]) that the higher the number of relays the easier the convergence. However, given a sequence of encounters $\pi_i$, there exists a minimum number of relays that is enough for guaranteeing convergence. We denote this number as $min_i^{so}$, and we derive it in Lemma 3.

*Lemma 3 ($min_i^{so}$):* Assuming that intermediate relays are selected in the order specified by sequence $\pi_i$, the expected delay from intermediate relays to the destination $d$ will converge if and only if there are *at least $min_i^{so}$ intermediate relays*, with $min_i^{so}$ being equal to the following:

$$min_i^{so} = \arg\min_m \{f_{min}^{so}(m, \pi_i) > 0\}, \quad (3)$$

where $f_{min}^{so}(m, \pi_i)$ is defined as $\sum_{z=1}^{m} \alpha_z^{(i)} - (1 + m)$ and $\alpha_z^{(i)}$ denotes the exponent $\alpha_{jd}$ associated with the $z$-th node in encounter sequence $\pi_i - \{d\}$.

*Proof:* The second hop can be modelled as a parallel delivery from $m$ relays to the destination. Let us consider the $i$-th relay, assuming that it receives its copy of the message at time $t_i$. The time before the $i$-th relay hands over its copy to the destination can be modeled as a residual intermeeting time (Lemma 1). Considering all $m$ relays, the time before the first of the $m$ copies reaches the destination can be modeled as the minimum of the residual intermeeting times between the relays and the destination. Once we focus on a specific sequence $\pi_i - \{d\}$ of encounters at the source node, it is clear that the first $m$ relays correspond to the first $m$ nodes in the sequence. We denote with $\alpha_z^{(i)}$ the $\alpha_{jd}$ exponent associated with the $z$-th node in $\pi_i$. Then, applying property P3, we obtain the convergence condition $\sum_{z=1}^{m} \alpha_z^{(i)} - m > 1$. Since, as discussed before, convergence becomes easier as $m$ increases, the minimum number $min_i^{so}$ of copies required at the second hop for convergence under sequence of encounters $\pi_i - \{d\}$ corresponds to the first (integer) $m$ value in the above equation for which the condition is satisfied. Hence Equation 3 follows. For a detailed proof, see [25]. □

*Corollary 2:* Quantity $min_i^{so}$ derived in Lemma 3 takes values in $[min_{lo}^{so}, min_{up}^{so}]$. The upper and lower bounds on $min_i^{so}$ (corresponding to the worst and best case for convergence) are reached when $\pi_i$ corresponds to the sequence of nodes ordered in increasing and decreasing order of their exponents $\alpha_{jd}$, respectively.

Now, we use the results in Lemmas 2 and 3 for deriving the stability region of the delay under $m$-copy 2-hop social-oblivious forwarding (with $m < N - 1$).

*Theorem 4 ($m$-copy 2-hop scheme):* $E[D_{sd}]$ converges if and only if the following condition holds true:

C3 $\quad m \geq min_{up}^{so} \wedge max_i^{so} \geq min_i^{so}, \forall i \in \{1, \ldots, |\mathcal{P}_s^p|\}$,

where set $\mathcal{P}_s^p$ is the set of all permutations for elements in $\mathcal{P}_s$.

*Proof:* In order to derive C3, first we notice that, e.g., the first hop and second hop worst case ($max_{lo}^{so}$ and $min_{up}^{so}$, respectively) in general do not happen simultaneously. In fact, meeting processes between nodes are

This article has been accepted for publication in a future issue of this journal, but has not been fully edited. Content may change prior to final publication. Citation information: DOI 10.1109/TMC.2014.2316506, IEEE Transactions on Mobile Computing

8

independent, and the fact that node $j$ meets the destination frequently (high $\alpha_{jd}$) does not generally imply that it also meets the source node frequently (high $\alpha_{sj}$), and vice versa. Thus, the order on set $\{\alpha_{sj}\}$ determined by sequence $\pi_i$ does not correspond to the same ordering on set $\{\alpha_{jd}\}$. Since worst cases are not correlated (either positively or negatively), we have to impose convergence on all the possible combinations for relay selections. This implies deriving a sequence of $\alpha_{sj}$ based on $\pi_i$ and its corresponding sequence of $\alpha_{jd}$ and verifying convergence for each of these permutations. In practice, we compute a pair $(max_i^{so}, min_i^{so})$ for each possible sequence $\pi_i$ of relays. Convergence is possible as long as $max_i^{so} \geq min_i^{so}$ for all permutations, since this means that the first hop is always able to provide to the second hop the number of copies needed for convergence. When the above condition is satisfied, convergence is ensured as long as we send a number $m$ of copies equal to or greater than the number of copies needed in the worst case at the second hop, hence[6] we set $m \geq min_{up}^{so}$. □

*Corollary 3:* A sufficient condition for the convergence of the expected delay under the memoryful $m$-copy two-hop forwarding scheme in Theorem 4 is given by the following:

C3$_{[s]}$    $m \geq min_{up}^{so} \wedge min_{up}^{so} \leq max_{lo}^{so}$.

*Proof:* The sufficient condition C3$_{[s]}$ follows directly from Lemmas 2 and 3. What these lemmas told us is that, in the worst case, the first hop delivery can at most provide $max_{lo}^{so}$ copies (with finite first hop expected delay) while, again in the worst case, the second hop delivery needs at least $min_{up}^{so}$ copies. When C3$_{[s]}$ holds true, it is guaranteed that, in all cases, the minimum number of copies needed at the second hop is provided by the first hop, thus proving the sufficiency of the condition. □

As discussed before, Chaintreau et al. [8] studied the $m$-copy two-hop scheme under homogeneous mobility patterns (corresponding to $\alpha_{ij} = \alpha, \forall i, j$). For the sake of completeness, in [25] (Appendix C.2.1) we verify that Theorem 4 confirms and extends the results in [8].

### 5.2.2  Multi-hop forwarding

Again we consider a social-oblivious protocol in which the source node generates $m$ copies of the message and hands them over to the first $m$ nodes encountered. Once the source node has handed over the $m$ copies, these copies travel along multi-hop social-oblivious paths until the destination is found. Theorem 5 describes the convergence conditions that apply in this case.

*Theorem 5 (m-copy n-hop scheme):* $E[D_{sd}]$ converges if and only if condition C1 and C2 in Theorem 2 hold true.

*Proof:* As we did before, we only sketch the proof and we refer the reader to Appendix C in [25] for

the rigorous mathematical derivation. Here, the source node is memoryful and thus it guarantees that the $m$ copies are relayed to $m$ distinct nodes. However, it is possible to prove that, after the first hop, there is a non negligible probability that all $m$ copies are relayed to the same node. This is clearly a worst case as far as the convergence of the expected delay is concerned, because the parallel delivery offered by the multi-copy approach is not exploited. Since basically the multi-copy forwarding process turns into a 1-copy $n$-hop scheme, it means that copies in addition to the first one are useless in terms of convergence. Thus, we simply need to ensure that at least one copy achieves convergence, which is guaranteed by the same conditions applying to the 1-copy $n$-hop scheme, i.e., C1 and C2. □

### 5.3  Discussion

Table 1 summarises the results derived so far for social-oblivious forwarding protocols. In the following we will informally speak about *convergent* relays to denote nodes for which the associated convergence condition is satisfied. The first interesting finding is that $n$-hop social-oblivious protocols (last two columns of Table 1) are no more effective in delivering the message with finite expected delay than the simple 1-copy 2-hop forwarding. In fact, they both share the same convergence conditions (C1 and C2), but the former consumes much more network resources than the latter. So, if we are only interested in the convergence of the expected delay, paths with more than two hops should be avoided, as two hops ensure that the available forwarding diversity between nodes is explored, while minimizing resource consumption.

With social-oblivious protocols, when the source node meets the destination with a residual intermeeting time having $\alpha_{sd} > 2$, there is no reason to exploit other relays, as this will only introduce the chance of picking a bad relay. In fact, when the number of hops is allowed to grow, we have to impose on intermediate relays additional constraints that are not needed by Direct Transmission (see, e.g., condition C2 in Theorem 3 which requires that the residual intermeeting time between any relay and the destination achieves a finite expectation).

Different is the situation in which $\alpha_{sd} \leq 2$. In this case, the source node cannot reach destination $d$ directly with a finite expected delay but it may be able to hand over the message to other nodes within a finite expected time. Hence, exploring more relays can prove convenient. If these intermediate relays are *all* able to individually deliver the message to the destination within a finite expected time, then the 1-copy 2-hop strategy guarantees convergence while minimizing resource consumption.

Instead, when there exists at least one divergent intermediate relay, the most effective strategy is the $m$-copy 2-hop forwarding, under which the source is able to send up to $max_{up}^{so}$ copies of the message. If $max_{up}^{so} = 1$, we find again conditions C1 and C2 that hold for the 1-copy 2-hop strategy. However, if the source node can

---

6. Please note that $m$ can be configured to be smaller or greater than $max_i^{so}$ for a given $\pi_i$. In the first case, the source node will simply send $m$ convergent copies rather than $max_i^{so}$. In the second case, the source node will be able to send $max_i^{so}$ with finite first hop expected delay and all other copies will be divergent.

| | 1 hop | | 2 hops | | $n$-hop | |
|---|---|---|---|---|---|---|
| | 1 copy | $m$ copies | 1 copy | $m$ copies | 1 copy | $m$ copies |
| social-oblivious | $\alpha_{sd} > 2$ | - | [C1,C2] | [C3] | [C1,C2] | [C1,C2] |
| social-aware | $\alpha_{sd} > 2$ | - | [C4,C5] | [C8] | [C6, C7] | [C6, C7] |

TABLE 1
Summary of convergence conditions for social-oblivious and social-aware routing strategies

reach operating point $max_{up}^{so} > 1$, conditions on the delivery from the relays to the destination become less restrictive since the more the copies sent out by the destination (with finite first hop expected delay) the easier the convergence at the second hop (Lemma 3).

# 6 EXPECTED DELAY CONVERGENCE FOR SOCIAL-AWARE SCHEMES

In this section our goal is to derive the convergence conditions for the social-aware approaches introduced in Section 4, which will then be used to investigate whether the social-aware approach outperforms the best social-oblivious ones. In the following, we denote with $\mathcal{R}_i$ the set of possible relays for node $i$, i.e., the set of nodes whose fitness is greater than that of node $i$. Recall that, with social-aware forwarding, nodes can hand over a message only to nodes with higher fitness.

## 6.1 Single-copy schemes

We start our discussion with the case of single copy schemes. Please recall that social-aware strategies do not make sense when only one hop is allowed, since this hop is necessarily the destination itself and Theorem 1 holds. Thus we go straight to the 1-copy 2-hop case.

*Theorem 6 (1-copy 2-hop social-aware scheme):* $E[D_{sd}]$ converges if and only if the following conditions hold:

C4 $\quad \sum_{j \in \mathcal{R}_s} \alpha_{sj} > 1 + |\mathcal{R}_s|$
C5 $\quad \alpha_{jd} > 2, \forall j \in \mathcal{R}_s - \{d\}$.

*Proof:* The proof is a step-by-step repetition of the proof of Theorem 2, with the only difference that this time relays belong to $\mathcal{R}_s$, thus we omit the proof. □

Theorem 6 mirrors Theorem 2 with the exception that only nodes with fitness higher than that of the source node can be selected. At first sight, this seems only a minor difference, but it proves extremely significant in all those cases in which the source node is already a "good" relay (from the convergence standpoint). In these cases, in fact, with social-aware forwarding we are sure that only relays better than the source node can be picked, thus ensuring that convergence can only improve, never get worse, at the second hop.

*Theorem 7 (1-copy $n$-hop social-aware scheme):* $E[D_{sd}]$ converges if and only if the following condition holds:

C6 $\quad \sum_{j \in \mathcal{R}_i} \alpha_{ij} > 1 + |\mathcal{R}_i|$ for all $i \in \mathcal{R}_s \cup \{s\}$
C7 $\quad n \geq |\mathcal{D}| + 1$,

where set $\mathcal{D}$ comprises nodes $j \in \mathcal{R}_s$ whose exponent value $\alpha_{jd}$ is smaller than or equal to 2.

*Proof:* The proof exploits the ordering guaranteed by social-aware policies. Specifically, when social-aware policies are used, messages are forwarded along a path
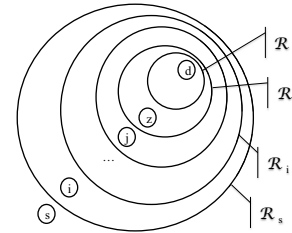


Fig. 1. Social forwarding at a glance

with increasing fitness. For the sake of simplicity, in the following we assume that there cannot be two nodes with the same fitness value. Recalling that $\mathcal{R}_i$ denotes the set of potential relays when the message is on node $i$, we have that, for a generic path $\{s, i, \cdots, j, z, d\}$ with increasing fitness, the relation $\mathcal{R}_s \supset \mathcal{R}_i \supset \cdots \mathcal{R}_j \supset \mathcal{R}_z$ holds. Exploiting Lemma 1, we know that the time before the message leaves a generic node $i$ is distributed as $\min_{j \in \mathcal{R}_i}\{R_{ij}\}$. According to property P3, the above expression has a finite expectation if and only if $\sum_{j \in \mathcal{R}_i} \alpha_{ij} > 1 + |\mathcal{R}_i|$ (condition C6).

In order to complete the proof, we have to consider the fact that when a message has reached the maximum number $n - 1$ of allowed intermediate hops, the relay currently holding the message can only deliver it to the destination directly. Thus, $\alpha_{jd} > 2$ is required after $n - 1$ relays have been reached. Let us split all possible relays in $\mathcal{R}_s$ into two subsets $\mathcal{C}$ and $\mathcal{D}$, such that $\mathcal{C} \cup \mathcal{D} = \mathcal{R}_s$. Subset $\mathcal{C}$ contains all nodes $j \in \mathcal{R}_s$ such that $\alpha_{jd} > 2$, while subset $\mathcal{D}$ contains those nodes $j \in \mathcal{R}_s$ with exponent $\alpha_{jd}$ smaller than or equal to 2. Please note that, due to social-aware forwarding rules, once a relay in $\mathcal{C}$ is picked, all subsequent relays will be also drawn from $\mathcal{C}$, since nodes in $\mathcal{C}$ are "closer" to the destination than those in $\mathcal{D}$. As far as convergence is concerned, in the worst case, all nodes in $\mathcal{D}$ are exploited before those in $\mathcal{C}$. So, if we set $n - 1$, i.e., the maximum number of intermediate hops allowed, to be greater than or equal to $|\mathcal{D}|$, we are sure that, even in the worst case, a relay in $\mathcal{C}$ is eventually selected. Since for relays in $\mathcal{C}$ convergence is guaranteed (in fact, $\alpha_{jd} > 2$, when $j \in \mathcal{C}$), the overall expected delay will converge. □

## 6.2 Multi-copy schemes

Frequently, social-aware schemes are multi-copy. In the following we analyze whether using multiple copies can help the convergence of the expected delay when social-aware schemes are in use. The proofs of this section follow the same line of reasoning of the corresponding social-oblivious versions, once substituting $\mathcal{P}_i$ with $\mathcal{R}_i$. For this reason, in the following we omit them.

### 6.2.1 Two-hop forwarding

First, we focus on the $m$-copy 2-hop scheme. To this aim, we derive Theorem 8, which is in turn based on the following lemmas. Please note that in this case sequence $\pi_i$ only contains nodes that belong to $\mathcal{R}_i$.

*Lemma 4 ($max_i^{sa}$):* When intermediate relays are selected according to sequence $\pi_i$, the source node is able to deliver at most $max_i^{sa}$ copies with finite first hop expected delay, with $max_i^{sa}$ being equal to the following:

$$max_i^{sa} = \arg\min_m\{f_{max}^{sa}(m, \pi_i) > 0\}, \qquad (4)$$

where $f_{max}^{sa}(m, \pi_i) = m + \sum_{z=m}^{|\mathcal{R}_s|} \alpha_z^{(i)} - (2 + |\mathcal{R}_s|)$ and $\alpha_z^{(i)}$ denotes the exponent $\alpha_{sj}$ of the $z$-th node in sequence $\pi_i$.

*Corollary 4:* Quantity $max_i^{sa}$ derived in Lemma 4 takes values in the interval $[max_{lo}^{sa}, max_{up}^{sa}]$. The upper and lower bound on $max_i^{sa}$ (corresponding to the best and worst case for convergence) are reached when $\pi_i$ corresponds to nodes in $\mathcal{R}_s$ encountered in increasing and decreasing order of $\alpha_{sj}$, respectively.

*Lemma 5 ($min_i^{sa}$):* Assuming that intermediate relays are selected in the order specified by sequence $\pi_i$, the expected delay from intermediate relays to the destination $d$ will converge if and only if there are *at least* $min_i^{sa}$ intermediate relays, with $min_i^{sa}$ being equal to the following:

$$min_i^{so} = \arg\min_m\{f_{min}^{sa}(m, \pi_i) > 0\}, \qquad (5)$$

where $f_{min}^{sa}(m, \pi_i) = \sum_{z=1}^m \alpha_z^{(i)} - (1 + m)$ and $\alpha_z^{(i)}$ denotes the exponent $\alpha_{jd}$ associated with the $z$-th node in encounter sequence $\pi_i - \{d\}$.

*Corollary 5:* Quantity $min_i^{sa}$ derived in Lemma 5 takes values in $[min_{lo}^{sa}, min_{up}^{sa}]$. The upper and lower bounds on $min_i^{sa}$ (corresponding to the worst and best case for convergence) are reached when $\pi_i$ corresponds to the sequence of nodes (belonging to $\mathcal{R}_s$) ordered in increasing and decreasing order of their exponents $\alpha_{jd}$, respectively.

Lemmas 4 and 5 are the social-aware equivalent of Lemmas 2 and 3. Using their results, the following theorem about the 2-hop convergence can be derived, with $m < |\mathcal{R}_s|$.

*Theorem 8 ($m$-copy 2-hop social-aware scheme):* $E[D_{sd}]$ achieves convergence if and only if the following condition holds true:

C8 $\qquad m \geq min_{up}^{sa} \wedge max_i^{sa} \geq min_i^{sa}, \forall i \in \{1, \ldots, |\mathcal{R}_s^p|\}$,

where set $\mathcal{R}_s^p$ is the set of all permutations $\pi_i$ for set $\mathcal{R}_s$.

*Corollary 6:* A sufficient condition for the convergence of the expected delay under the social-aware $m$-copy two-hop forwarding scheme in Theorem 4 is given by the following:

C8$_{[s]}$ $\quad m \geq min_{up}^{sa} \wedge min_{up}^{sa} \leq max_{lo}^{sa}$.

Comparing the social-aware $m$-copy 2-hop with its social-oblivious counter part is not straightforward. In Section 7 we prove analytically that there is no clear

winner between the two, and that either one or both can achieve convergence depending on the mobility scenario considered.

### 6.2.2 Multi-hop forwarding

Finally, in Theorem 9 we consider the most general case in which the source node generates $m$ copies for the message and each of them travel up to $n$ hops along independent paths. We find that also in the social-aware case, multiple copies used together with multiple hops do not improve convergence with respect to the simple 1-copy $n$-hop scheme. Intuitively, this is because in the worst case (which occurs with non-negligible probability) all copies, after the first hop, follow the same multi-hop path, which requires conditions C6 and C7 for convergence.

*Theorem 9 ($m$-copy $n$-hop social-aware scheme):* $E[D_{sd}]$ converges if and only if conditions C6 and C7 in Theorem 7 hold true.

## 6.3 Discussion

Table 1 summarizes the convergence conditions for social-aware schemes derived so far. As in the social-oblivious case, multi-hop schemes do not benefit from the use of multiple copies, and in fact the 1-copy $n$-hop scheme and the $m$-copy $n$-hop scheme share the same convergence conditions. Similarly, the difference between 2-hop schemes mirrors that between the corresponding social-oblivious versions. Thus, the 1-copy 2-hop scheme is effective when $\alpha_{jd} > 2$ for all $j \in \mathcal{R}_s$, since it allows us to save resources by sending a single copy. However, when condition C5 does not hold, the only chance to achieve convergence is to exploit multiple copies.

If we focus on single-copy schemes, it is interesting to note that, differently from the social-oblivious case in which using additional hops did not provide any advantage, 1-copy social-aware schemes may benefit from multiple hops. In fact, for the 1-copy 2-hop scheme we need to impose that all intermediate relays $j$ meet the destination with $\alpha_{jd} > 2$, which is a quite strong condition. On the other hand, if we use multiple hops (1-copy $n$-hop case), conditions C6 and C7 are required, which are milder than C5. More specifically, assuming that there are no limitations to the value that we can assign to $n$, condition C7 can be easily satisfied. Then, C6 relates to the convergence of the minimum of a set of Pareto random variables, which is always easier to achieve than for any single random variable from the set (corresponding to condition C5). The only constraint for the 1-copy $n$-hop case is that there must be at least one node $z$ (the one with the highest fitness) meeting the destination with $\alpha_{zd} > 2$. In fact, for $z$, $\mathcal{R}_z = \{d\}$.

Finally, we compare the $m$-copy 2-hop case with the 1-copy $n$-hop case (which is equivalent to the $m$-copy $n$-hop scheme). There is no clear winner here, as each scheme can provide convergence when the other one cannot. For example, consider the case in which the

source node is not able to send more than one copy (i.e, $max_i^{sa} = 1, \forall i \in \{1, \ldots, |\mathcal{R}_s^p|\}$). In this case, the $m$-copy 2-hop scheme becomes effectively a 1-copy 2-hop scheme, which fails to achieve convergence if some intermediate hop $j$ does not have exponent $\alpha_{jd}$ greater than 2 (condition C5). Instead, exploiting multiple hops pays off in this case, as it allows us to rely on more intermediate relays, which may not meet the destination within a finite expected time but can bring the message "closer" to nodes that do meet $d$ with $\alpha_{jd} > 2$. Vice versa, when $max_i^{sa} > 1$ for some $i$, the cooperative delivery of the multiple copies can overcome the presence of intermediate relays for which conditions C6-C7 do not hold. For example, when there is not even one relay $j$ with $\alpha_{jd} > 2$, then the $m$-copy 2-hop case is the only possible choice.

# 7 COMPARING SOCIAL-AWARE AND SOCIAL-OBLIVIOUS STRATEGIES

In the previous sections we have separately analyzed the convergence properties of social-oblivious and social-aware forwarding schemes, identifying the best strategies, from the convergence standpoint, for each of the two categories. In the following we take the champions of each class and we investigate whether there is a clear winner between social-oblivious and social-aware strategies when it comes to the convergence of their expected delay.

Let us first consider the case $\alpha_{sd} > 2$. We have seen in Section 5.3 that with this configuration the Direct Transmission scheme is the best choice from the convergence standpoint. In fact, with social-oblivious schemes using more than one hop, "bad" relays can be selected even starting from a source that is already able to reach the destination with a finite expected residual intermeeting time. This does not happen with social-aware strategies. In fact, assume that the source is the only node with $\alpha_{sd} = 2 + \varepsilon$, while all other nodes meet the destination with $\alpha_{jd} = 1 + \varepsilon$, with $\varepsilon$ being a very small quantity. In the social-aware case, $\mathcal{R}_s$ contains only the destination, as all other nodes are clearly worse than the source node as relay. This shows the adaptability of social-aware schemes: the additional knowledge that they exploit makes them able to resort to simpler approaches (in this case, $\mathcal{R}_s = \{d\}$ is equivalent to the Direct Transmission) when they realize that additional resources in terms of number of copies or number of hops would not help the forwarding process. This implies that one can safely use the $m$-copy 2-hop or the 1-copy $n$-hop social-aware protocols because in the worst case they will do no harm (they will downgrade to simpler strategies, without exploiting wrong paths), while in the best case they are able to improve the convergence of the forwarding process.

When $\alpha_{sd} \leq 2$ and $\alpha_{jd} > 2$ for all nodes $j$ in the relay set (i.e., $j \in \mathcal{R}_s - \{d\}$ for the social-aware case and $j \in \mathcal{P}_s - \{d\}$ for the social-oblivious case), the strategy of choice is the 1-copy 2-hop for both the social-oblivious

and social-aware category. However, the 1-copy 2-hop social-aware scheme is overall more advantageous than its social-oblivious counterpart. More specifically, when the source node is the worst relay for the destination (i.e., $\min_i\{\alpha_{id}\} = \alpha_{sd}$), the social-oblivious and the social-aware approaches are equivalent (given that $\mathcal{P}_s = \mathcal{R}_s$). In all other cases, instead, $\mathcal{R}_s \subset \mathcal{P}_s$, thus, for the set of nodes in $\mathcal{P}_s - \mathcal{R}_s$, social-aware forwarding does not impose any constraint, while social-oblivious forwarding needs to impose constraints, thus resulting in stricter conditions for convergence.

Let us now focus on the remaining cases, namely i) when $\alpha_{sd} \leq 2$ and not all intermediate relays have exponent greater than 2, and ii) when $\alpha_{jd} \leq 2$ for all nodes $j$. In the first case, the social-aware $m$-copy 2-hop, the social-aware 1-copy $n$-hop, and the social-oblivious $m$-copy 2-hop can achieve convergence. In the second case, the only options for convergence are the social-aware $m$-copy 2-hop and the social-oblivious $m$-copy 2-hop. We first highlight the differences between the $n$-hop approach and the 2-hop approach by discussing when the social-aware 1-copy $n$-hop outperforms the other two strategies in terms of convergence (which can only happen in case $i$), then we focus on the social-aware and social-oblivious $m$-copy 2-hop strategies, thus covering both case $i$ and $ii$.

So, assume that there exists at least one node $z$ that meets the destination with $\alpha_{zd} > 2$. The $m$-copy 2-hop strategies send multiple copies to a set of relays, which in turn can only deliver the message to the destination directly. This implies that intermediate relays must have collectively the capability of reaching the destination, for all subsets with size $m$ of possible relays. Here, only meetings with the destination are relevant, and if all relays but $z$ have very low exponent for encounters with the destination, convergence may not be achieved. Differently from the 2-hop strategies, the social-aware $n$-hop scheme do not rely exclusively on the capabilities of meeting with $d$, but it is able to generate a *path* towards the destination in which intermediate nodes may not be good relays for $d$ but good relays towards nodes with high fitness (in the extreme case, only $\alpha_{zd} > 2$ can hold). Thus, in the $n$-hop case, as long as the message can leave intermediate relays within a finite expected time, this could be enough for convergence. An example scenario is provided in Section 7.1. When all three strategies achieve convergence, the one to be preferred can be chosen based on resource consumption considerations. With the $m$-copy 2-hop strategies there can be up to $2m$ transmissions, while with the 1-copy $n$-hop scheme there are $n$. Hence, when $n < 2m$, the single-copy scheme should be preferred.

Let us finally compare the social-oblivious and the social-aware $m$-copy 2-hop schemes. Since they seem to cover similar mobility scenarios (as discussed in the previous section) and to be based on similar mechanisms (the $min_i$ and $max_i$ quantities, whose relation with $m$ determines the convergence), it may be difficult to intu-

itively evaluate which one performs better in terms of convergence. For this reason, Theorem 10 below (whose proofs can be found in [25]) we tackle this problem from an analytical perspective. In Appendix E of [25] we provide a concrete example for both cases.

*Theorem 10:* Since both the following configurations are feasible under the conditions in Lemma D1, it may happen that either the social-oblivious $m$-copy 2-hop scheme achieves convergence when the social-aware $m$-copy 2-hop scheme does not (Equation 6), or vice versa (Equation 7), depending on the underlying mobility process.

$$max_i^{so} \geq min_i^{so} \geq min_i^{sa} > max_i^{sa} \qquad (6)$$

$$min_i^{so} > max_i^{so} \geq max_i^{sa} \geq min_i^{sa} \qquad (7)$$

Intuitively, an example of the first case is when there are a lot of nodes that meet the source with high $\alpha_{sj}$ (thus resulting in high $max_i^{so}$, high enough to be greater than $min_i^{so}$); if those relays have very low $\alpha_{jd}$, they will not be used by the social-aware scheme, thus resulting in a low $max_i^{sa}$, possibly not high enough to guarantee that the second hop converges. It is easy to construct a corresponding example for the other case.

## 7.1 Example

In order to complement the theoretical discussion of the previous section, in the following we provide a concrete example for the case in which the social-aware 1-copy $n$-hop scheme is the only one achieving convergence. In [25] we also provide two concrete examples for the cases discussed in Theorem 10.

The mobility scenario we consider is described by the exponent matrix in Figure 2. Please note that such matrix has been chosen just to exemplify the behaviour of the social-aware 1-copy $n$-hop scheme. For a more realistic analysis, please refer to Section 8. Element $\alpha_{ij}$ in matrix $\boldsymbol{\alpha}$ (of size 10) gives the Pareto exponent for the $i, j$ node pair. We assume that node $i = 1$ is the source node and that node $j = 10$ is the destination. In this case the source node is the node with the lowest fitness value, thus the $m$-copy 2-hop social-oblivious and social-aware schemes overlap (in fact, $\mathcal{P}_s = \mathcal{R}_s$).

$$\boldsymbol{\alpha} = \begin{pmatrix} 0 & 1.22 & 1.22 & 1.22 & 1.22 & 1.22 & 1.22 & 1.22 & 1.22 & 1.125 \\ 1.22 & 0 & 1.33 & 1.33 & 1.33 & 1.33 & 1.33 & 1.33 & 1.33 & 1.13 \\ 1.22 & 1.33 & 0 & 1.44 & 1.44 & 1.44 & 1.44 & 1.44 & 1.44 & 1.14 \\ 1.22 & 1.33 & 1.44 & 0 & 1.55 & 1.55 & 1.55 & 1.55 & 1.55 & 1.15 \\ 1.22 & 1.33 & 1.44 & 1.55 & 0 & 1.66 & 1.66 & 1.66 & 1.66 & 1.16 \\ 1.22 & 1.33 & 1.44 & 1.55 & 1.66 & 0 & 1.77 & 1.77 & 1.77 & 1.17 \\ 1.22 & 1.33 & 1.44 & 1.55 & 1.66 & 1.77 & 0 & 1.88 & 1.88 & 1.18 \\ 1.22 & 1.33 & 1.44 & 1.55 & 1.66 & 1.77 & 1.88 & 0 & 1.99 & 1.19 \\ 1.22 & 1.33 & 1.44 & 1.55 & 1.66 & 1.77 & 1.88 & 1.99 & 0 & 2.1 \\ 1.11 & 1.12 & 1.13 & 1.14 & 1.15 & 1.16 & 1.17 & 1.18 & 2.1 & 0 \end{pmatrix}$$

Fig. 2. Exponent matrix

We start with the 1-copy $n$-hop scheme. The size of set $\mathcal{D}$ is 8, since there are eight nodes with $\alpha_{jd} \leq 2$. Thus, we need to set the maximum number of allowed hop $n$ to 9 (condition C7). Then, we compute $\sum_{j \in \mathcal{R}_i} \alpha_{ij} - (1 + \mathcal{R}_i)$ (condition C6) for all nodes $i \in \mathcal{R}_s \cup \{s\}$ (Table 2). Since the computed quantities are greater than zero for

TABLE 2
Condition C6 for each relay (including the source)

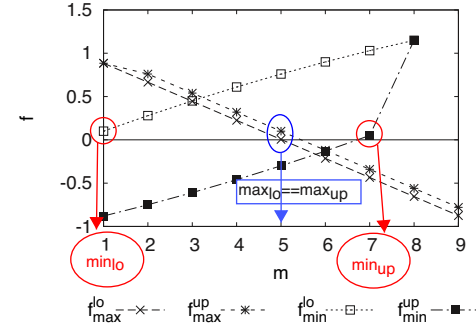| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| $C6$ | 0.89 | 1.44 | 1.78 | 1.9 | 1.8 | 1.48 | 0.94 | 0.18 | 0.1 |



Fig. 3. $min_i$ and $max_i$ for the $m$-copy 2-hop scheme.

all possible relays (including the source node), the 1-copy $n$-hop social-aware scheme achieves convergence in this scenario.

We now focus on the $m$-copy 2-hop scheme, recalling that the social-oblivious version and the social-aware version are equivalent in this case (thus we drop superscripts $so$ and $sa$). In order to verify sufficient condition C3$_{[s]}$, we need to find $max_{lo}$ and $min_{up}$, i.e., the value of $max_i$ and $min_i$ in the worst case. According to Corollary 1, $max_{lo}$ is reached when permutation $\pi_i$ corresponds to relays encountered in decreasing order of $\alpha_{sj}$, while, according to Corollary 2, $min_{up}$ is achieved when permutation $\pi_i$ corresponds to relays encountered by the source node in increasing order of $\alpha_{jd}$. We denote these two permutations as $\pi_i^*$ and $\pi_i'$ respectively. In Figure 3, we plot function $f_{max}^{lo}(m) = f_{max}(m, \pi_i^*)$ corresponding to the case in which $max_{lo}$ is reached and function $f_{min}^{up}(m) = f_{min}(m, \pi_i')$ corresponding to the case in which $min_{up}$ is achieved. Recall that $min_{up}$ corresponds to the first $m$ value for which $f_{min}^{up}$ is greater than zero, thus $min_{up} = 7$. Similarly, $max_{lo}$ corresponds to the last $m$ value for which function $f_{max}^{lo}$ is greater than zero, and so $max_{lo} = 5$. Since $max_{lo} < min_{up}$, sufficient condition C3$_{[s]}$ is not satisfied.

It is easy to show that also the necessary and sufficient condition C3 does not hold. Recall that the necessary and sufficient condition states that convergence is ensured as long as $max_i \geq min_i$ for all encounter permutations $\pi_i$. However, this does not happen here. Consider (Figure 3) $f_{max}^{up}$ and $f_{min}^{lo}$, i.e., functions $f_{max}$ and $f_{min}$ in the best case. The first integer values of $m$ before the functions become negative determine the values of $max_{up}$ and $min_{lo}$. Since $max_{up} = 5$, from Corollary 1 we have that $max_i$ varies in the range $[5, 5]$, i.e., $max_i$ is always equal to 5 regardless of the permutation considered. This means that, for the permutation corresponding to the $min_i$ worst case, the source node will not be able in any case to send more than 5 copies with finite first-hop expected delay (while 7 are required). Hence, convergence cannot be achieved in this case.
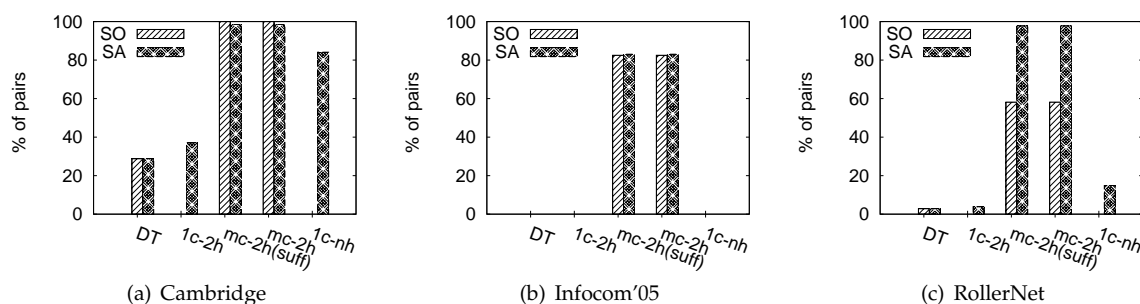
Fig. 4. Percentage of convergent pairs in traces.

## 8 CONVERGENCE IN MOBILITY TRACES

We conclude the paper by applying the convergence conditions derived in the previous sections to a set of contact traces that are frequently used in the literature: Cambridge [26], Infocom'05 [26], and RollerNet [27], composed respectively of 12, 41, and 62 nodes. Due to space limitations, we do not recall further properties of these traces and we move straight to the application of our convergence conditions. The Pareto exponents are estimated using Maximum Likelihood Estimation, setting $t_{min}$ equal to the sampling interval used for collecting the trace ($120s$ for both Infocom'05 and Cambridge, $15s$ for RollerNet). Under this configuration, after applying the Cramer-von Mises criterion, we obtain that the Pareto hypothesis cannot be rejected for $80\%$, $97\%$, and $85.5\%$ of pairs for Cambridge, Infocom'05, and RollerNet, respectively[7].

In Figure 4 we show the percentage of pairs for which the strategies identified in Section 4 achieve convergence, in both the social-oblivious and the social-aware case. We omit the $n$-hop schemes that share the convergence conditions of other strategies that consumes less resources, as discussed in Sections 5.3 and 6.3. However, we provide the results obtained by applying the sufficient convergence conditions in Corollaries 3 and 6, in order to highlight that in these cases they identify correctly the right set of pairs. In the Cambridge dataset (Figure 4(a)), the best performance in terms of convergence is delivered by the social-oblivious $m$-copy 2-hop scheme, with the social-aware $m$-copy 2-hop scheme just slightly behind. In the Infocom dataset (Figure 4(b)) there is basically a tie between the same two strategies, with the social-aware one performing slightly better in this case. In the RollerNet scenario (Figure 4(c)), the social-aware $m$-copy 2-hop scheme clearly outperforms the others.

From the Pareto exponent distribution point of view, the Cambridge dataset is the one more shifted towards higher values (min=1.39, median=1.68, max=2.86) with respect to Infocom (min=1.24, median=1.43, max=1.82) and RollerNet (min=1.27, median=1.58, max=3.27). This is reflected by the performance of the Direct Transmission scheme. In the Infocom dataset no pair meets with

an exponent higher than 2. This is a very unfortunate case from the convergence standpoint, but the $m$-copy 2-hop schemes are still able to overcome this limitation and to reach around $80\%$ of pairs with a convergent expected delay. Figure 4 also confirms the importance of taking into account heterogeneity when studying the convergence of forwarding strategies. In fact, applying the convergence condition for the homogeneous case derived in [25] ($\alpha > \frac{2}{N} + 1$, assuming that we set $m > \frac{N}{2}$) and using the estimated exponent for the aggregate intermeeting times ($\alpha = 1.44$), we would have obtained a "convergent" verdict for the $m$-copy 2-hop scheme *for all source-destination pairs* in the Infocom scenario, but this is not always the case, as shown in Figure 4(b).

## 9 CONCLUSIONS

Assuming heterogenous Pareto intermeeting times, in this paper we have derived the conditions on the Pareto exponents such that the expected delay of a large family of forwarding protocols is finite. Our main result for the social-oblivious case is that convergence is not improved by using more than two hops (and in some conditions direct transmission, with just one hop, is the most efficient choice). In the social-aware case, instead, allowing more than two hops can provide convergence when all other strategies fail. As for the comparison of single-copy and multi-copy schemes, we found that multi-copy strategies can, in some cases, outperform single-copy strategies in terms of convergence of the expected delay. The use of multiple copies, in fact, benefits from the parallel delivery of the message from different nodes, which may overcome the limitations of individual nodes in achieving a finite expected delay. Finally, comparing social-oblivious and social-aware multi-copy solutions we were able to prove mathematically that there is not a clear winner between the two, since either one can achieve convergence when the other fails depending on the underlying mobility scenario.

Besides the theoretical value of the above convergence model per se, we believe that such model has also important practical implications. For the majority of forwarding schemes, nodes would be able to evaluate online whether a policy can achieve convergence or not (hence they can decide which one is to be preferred). For example, convergence can be easily verified for 1-hop and 2-hop strategies, since it is perfectly reasonable

---

7. Different acceptance percentages in [27] are due to a different $t_{min}$ setting. We believe that it is more correct to set $t_{min}$ equal to the granularity of the trace because samples smaller than the granularity are a just a few and only emerge due to statistical fluctuations.

to assume that nodes can learn the Pareto exponents of their direct neighbours. The only relevant policy (from the convergence standpoint) for which nodes may not be able to verify the convergence online (because it would require the knowledge of the exponents of potentially distant nodes) is the social-aware $m$-copy $n$-hop scheme. One way to address this problem is to let the source node pick this strategy only as a last resort, i.e., only when it is not possible to collect the required exponents for testing its convergence and when two-hop schemes are not able to achieve it.

## REFERENCES

[1] T. Spyropoulos, K. Psounis, and C. Raghavendra, "Efficient routing in intermittently connected mobile networks: The single copy case," *IEEE/ACM Trans. on Netw.*, vol. 16, no. 1, pp. 63–76, 2008.

[2] M. Grossglauser and D. Tse, "Mobility increases the capacity of ad hoc wireless networks," *IEEE/ACM Trans. on Netw.*, vol. 10, no. 4, pp. 477–486, 2002.

[3] T. Spyropoulos, K. Psounis, and C. Raghavendra, "Efficient routing in intermittently connected mobile networks: The multiple-copy case," *IEEE/ACM Trans. on Netw.*, vol. 16, no. 1, pp. 77–90, 2008.

[4] C. Lee, "Heterogeneity in contact dynamics: helpful or harmful to forwarding algorithms in DTNs?" in *WiOPT'09*. IEEE, 2009, pp. 1–10.

[5] T. Spyropoulos, T. Turletti, and K. Obraczka, "Routing in Delay-Tolerant Networks Comprising Heterogeneous Node Populations," *IEEE Trans. Mobile Comput.*, pp. 1132–1147, 2009.

[6] A. Picu, T. Spyropoulos, and T. Hossmann, "An analysis of the information spreading delay in heterogeneous mobility dtns," in *IEEE WoWMoM*, 2012, pp. 1–10.

[7] Z. Haas and T. Small, "A new networking model for biological applications of ad hoc sensor networks," *IEEE/ACM Trans. on Netw.*, vol. 14, no. 1, pp. 27–40, 2006.

[8] A. Chaintreau, P. Hui, J. Crowcroft, C. Diot, R. Gass, and J. Scott, "Impact of human mobility on opportunistic forwarding algorithms," *IEEE Trans. Mobile Comput.*, pp. 606–620, 2007.

[9] H. Cai and D. Eun, "Crossing over the bounded domain: From exponential to power-law intermeeting time in mobile ad hoc networks," *IEEE/ACM Trans. on Netw.*, vol. 17, no. 5, pp. 1578–1591, 2009.

[10] T. Karagiannis, J.-Y. Le Boudec, and M. Vojnovic and, "Power law and exponential decay of intercontact times between mobile devices," *IEEE Trans. Mobile Comput.*, vol. 9, no. 10, pp. 1377–1390, 2010.

[11] A. Passarella and M. Conti, "Analysis of individual pair and aggregate inter-contact times in heterogeneous opportunistic networks," *IEEE Trans. Mobile Comput.*, vol. 12, no. 12, pp. 2483–2495, 2013.

[12] V. Conan, J. Leguay, and T. Friedman, "Characterizing pairwise inter-contact patterns in delay tolerant networks," in *Autonomics'07*, 2007.

[13] C. Boldrini, M. Conti, and A. Passarella, "Performance modelling of opportunistic forwarding under heterogenous mobility," *Computer Communications*, pp. 1–17, 2014.

[14] C. Lee and D. Eun, "Exploiting Heterogeneity in Mobile Opportunistic Networks: An Analytic Approach," in *IEEE SECON'10*. IEEE, 2010, pp. 1–9.

[15] A. Vahdat and D. Becker, "Epidemic routing for partially connected ad hoc networks," Tech. Rep., 2000.

[16] A. Lindgren, A. Doria, and O. Schelén, "Probabilistic routing in intermittently connected networks," *LNCS*, pp. 239–254, 2004.

[17] P. Hui, J. Crowcroft, and E. Yoneki, "Bubble rap: Social-based forwarding in delay tolerant networks," *IEEE Trans. Mobile Comput.*, vol. 10, no. 11, pp. 1576–1589, nov. 2011.

[18] E. Daly and M. Haahr, "Social network analysis for information flow in disconnected Delay-Tolerant MANETs," *IEEE Trans. Mobile Comput.*, pp. 606–621, 2008.

[19] C. Boldrini, M. Conti, and A. Passarella, "Exploiting users' social relations to forward data in opportunistic networks: The HiBOp solution," *Perv. and Mob. Comp.*, vol. 4, no. 5, pp. 633–657, 2008.

[20] P. Costa, C. Mascolo, M. Musolesi, and G. Picco, "Socially-aware routing for publish-subscribe in delay-tolerant mobile ad hoc networks," *IEEE J. Sel. Areas Commun.*, vol. 26, no. 5, pp. 748–760, 2008.

[21] C. Boldrini, M. Conti, and A. Passarella, "Less is more: long paths do not help the convergence of social-oblivious forwarding in opportunistic networks," *ACM/SIGMOBILE MobiOpp*, vol. 12, pp. 1–8, 2012.

[22] S. Burleigh, A. Hooke, L. Torgerson, K. Fall, V. Cerf, B. Durst, K. Scott, and H. Weiss, "Delay-tolerant networking: an approach to interplanetary internet," *IEEE Commun. Mag.*, vol. 41, no. 6, pp. 128–136, 2003.

[23] C. Boldrini, M. Conti, and A. Passarella, "From pareto inter-contact times to residuals," *IEEE Commun. Lett.*, vol. 15, no. 11, pp. 1256–1258, 2011.

[24] N. Johnson and S. Kotz, *Distributions in Statistics: Continuous Univariate Distributions: Vol. 1.* Houghton Mifflin New York, 1970.

[25] C. Boldrini, M. Conti, and A. Passarella, "The stability region of the delay in Pareto opportunistic networks," IIT-CNR, Tech. Rep. TR-13/2013, http://cnd.iit.cnr.it/chiara/pub/techrep/boldrini2013convergence_tr.pdf.

[26] P. Hui, A. Chaintreau, J. Scott, R. Gass, J. Crowcroft, and C. Diot, "Pocket switched networks and human mobility in conference environments," in *ACM CHANTS*. ACM, 2005, pp. 244–251.

[27] P.-U. Tournoux, J. Leguay, F. Benbadis, J. Whitbeck, V. Conan, and M. D. de Amorim, "Density-aware routing in highly dynamic DTNs: The rollernet case," *Mobile Computing, IEEE Transactions on*, vol. 10, no. 12, pp. 1755–1768, 2011.

**Chiara Boldrini** is a researcher at IIT-CNR. She received her M.Sc. degree in Computer Engineering and her Ph.D. in Information Engineering both from the University of Pisa, Italy, in 2006 and 2010, respectively. Her research interests are in the area of opportunistic and delay-tolerant networking. More specifically, she has been studying routing, content dissemination and energy saving issues, with a special focus on performance modelling using stochastic modelling techniques and complex network analysis. Alongside, she has been working on models for human mobility that can realistically represent the behaviour of mobile users.

**Andrea Passarella** (PhD in Comp. Eng. '05) is with IIT-CNR, Italy. He was a Research Associate at the Computer Laboratory, Cambridge, UK. He published 100+ papers on mobile social networks, opportunistic, ad hoc and sensor networks, receiving the best paper award at IFIP Networking 2011 and IEEE WoWMoM 2013. He was PC Co-Chair of IEEE WoWMoM 2011, Workshops Co-Chair of IEEE PerCom and WoWMom 2010, and Co-Chair of several IEEE and ACM workshops. He is in the Editorial Board of Elsevier Pervasive and Mobile Computing and Inderscience IJAACS. He was Guest Co-Editor of several special sections in ACM and Elsevier Journals. He is the Vice-Chair of the IFIP WG 6.3 "Performance of Communication Systems".

**Marco Conti** is a Research Director of the Italian National Research Council (CNR). He has published in journals and conference proceedings more than 300 research papers and four books related to design, modelling, and performance evaluation of computer networks, pervasive systems and social networks. He is Editor-in-Chief of Elsevier Computer Communications journal and Associate Editor-in-Chief of Elsevier Pervasive and Mobile Computing journal. He received the Best Paper Award at IFIP TC6 Networking 2011, IEEE ISCC 2012 and IEEE WoWMoM 2013. He has served as General and Program (Co)Chair for several conferences, including Networking 2002, IEEE WoWMoM 2005 and 2006, IEEE PerCom 2006 and 2010, ACM MobiHoc 2006, and IEEE MASS 2007.

# Optimal duty cycling in mobile opportunistic networks with end-to-end delay guarantees

Elisabetta Biondi, Chiara Boldrini, Andrea Passarella, and Marco Conti
IIT-CNR, Via G. Moruzzi 1, 56124 Pisa, Italy
Email: first.last@iit.cnr.it

*Abstract*—**Opportunistic communications have been recently proposed as a key strategy for offloading traffic from 3G/4G cellular networks, which is particularly beneficial in case of crowded areas where many users are interested in similar contents. To conserve energy, duty cycling schemes are typically applied, and therefore contacts between nodes may become intermittent and sporadic also in dense networks. It is thus of paramount importance to accurately tune the duty cycling policy in order to meet energy requirements without compromising the quality of communications. In this paper, building upon a model of duty cycling in opportunistic networks that we have validated in a previous work, we study how to optimise the value of the duty cycle in order to provide probabilistic guarantees on the delay experienced by messages. More specifically, for a broad range of end-to-end delay distributions, we provide closed-form approximated solutions for deriving the optimal duty cycle such that the probability that the delay is smaller than a target value $z$ is greater than or equal to a configurable probability $p$.**

## I. INTRODUCTION

Opportunistic networks have been conceived at the intersection between Mobile Ad hoc NETworks (MANET) and the Delay Tolerant (DTN) paradigm. In the conventional model, they exploit the movements of the nodes of the network (people with their smart, handheld devices like tablets and smartphones) in order to deliver messages to their destinations according to the store-carry-and-forward paradigm: nodes hold messages while they move and forward them to other nodes that are in radio contact, until messages reach their final destination. Opportunistic communications were initially seen as a standalone solution for those scenarios in which the nodes of the network were sparse and the infrastructure unavailable (disaster/emergency scenarios, developing countries, etc.). Recently, however, they have become one of the key strategies for mobile data offloading [1], whose main goal is to offload the traffic from cellular networks to other types of networks (e.g., WiFi infrastructured or MANET) in a synergic way, in order to address the overloading of the 3G/4G infrastructure.

In case of crowded environments (and thus dense networks) overloading may be even more critical, and opportunistic networking techniques can be usefully applied, as follows. Due to the typical Zipf-like shape of content interest, it is likely that large fractions of users in the crowd are interested in few, very popular contents (e.g., those mostly related to the area where the crowd gathers). Multicast can be a solution to reduce the traffic load only when content requests can be synchronised. When requests are generated dynamically by users, exploiting communications between users' devices is a more flexible solution, as content can be sent through the cellular network only to a few of them, exploiting opportunistic communications for the rest. The D2D technology addresses this goal to some extent, and is currently proposed in latest LTE releases. In this paper we focus on offloading through ad-hoc WiFi or Bluetooth technologies, as this approach permits to exploit additional portions of the spectrum (and, therefore, additional bandwidth) with respect to that allocated to cellular networks. A possible roadblock in this scenario is the fact that direct communications consume significant energy. To address this, nodes are typically operated in duty cycling mode, by letting their WiFi (or Bluetooth) interfaces ON only for a fraction of time. The joint effect of duty cycling and mobility is that, even if the network is dense, the resulting patterns in terms of communication opportunities is similar to that of conventional opportunistic networks, as devices are able to directly communicate with each other only when they come in one-hop radio range *and* both interfaces are ON.

The net effect of implementing a duty cycling scheme is thus the fact that some contacts between nodes are missed because the nodes are in power saving mode. Hence, detected intercontact times, defined as the time between two consecutive contact events during which a communication can take place for a pair of nodes, are longer than intercontact times determined only by mobility, when a duty-cycling policy is in place. This heavily affects the delay experienced by messages, since the main contribution to message delay is in fact due to the intercontact times. In our previous work [2], we have focused on exponentially distributed intercontact times and we have studied how these are modified by duty cycling, obtaining that intercontact times remain exponentially distributed but their rate is scaled by the inverse of the duty cycle (see Proposition 1, Section III). Building upon this result, we have then investigated how the first moments of the end to end delay vary with the duty cycle for a number of opportunistic forwarding schemes. In addition, we have found that energy saving and end-to-end delay both scale linearly with the duty cycling. Therefore, for a single message delivery, the same energy saved through duty cycling is spent because the network must stay alive longer. Thus, the main advantage of duty cycling is enabling the network to carry more messages by being alive longer (rather than improving the energy spent for each single delivery).

Our work in [2] assumed that the value of the duty cycle was given and studied its effects on important performance metrics such as the delay, the network lifetime, and the number of messages successfully delivered to their destination. More in general, the duty cycling can be seen as a parameter that can be configured, typically, based on some target performance metrics. To this aim, the main contribution of this paper

is a mathematical model that allows us to tune the duty cycle in order to meet a given target performance, expressed as a probabilistic guarantee (denoted as $p$) on the delay experienced by messages. Considering probabilistic, instead of hard, guarantees, allows us to cover a very broad range of application scenarios also beyond best-effort cases – all but those requiring real-time streaming. Specifically, we study the case of exponential, hyper-exponential and hypo-exponential delays (please recall that any distribution falls into one of these three cases, at least approximately [3]), deriving the optimal duty cycle for each of them. For the simple case of exponential delays we are able to provide an exact solution. For the other two cases, we derive an approximated solution and the conditions under which this approximation introduces a small fixed error $\varepsilon$ (which is always below $0.14$) on the target probability $p$. Specifically, in the worst case, the approximated duty cycle introduces an error on the target probability $p$ of about $0.1$ (hyper-exponential case) and $0.14$ (hypo-exponential case), while in the other cases the error is well below these thresholds.

The paper is organised as follows. In Section II we overview the literature on duty cycle optimisation for opportunistic networks. After having introduced the network and duty cycle model that we consider in this work (Section III) we derive in Section IV the optimal duty cycles for the case of exponential, hyper-exponential, and hypo-exponential delays. Then, in Section V, given a target performance for the delay, we discuss how the optimal duty cycle affects the volume of messages delivered during the network lifetime and we highlight that in the case of hyper-exponential delays it is possible to achieve a lower duty cycle than hypo-exponential delays for a given target performance. Finally, Section VI concludes the paper.

## II. RELATED WORK

There are not many contributions in the DTN literature studying the optimisation of the duty cycling policy. In [4], using a fixed duty cycle scheme, Wang et al. study the relationship between the probability of missing a contact and the associated energy consumption (considered inversely proportional to the contact probing interval). Building upon these results, [4] provides some heuristic algorithms to achieve an optimal contact probing. Differently from this work, in this paper we mathematically define the optimisation problem and we provide an analytical, closed form, result.

In [5], Gao and Li focus on the design of an adaptive duty cycle that minimises wakeups during intercontact times (which are useless, from a contact probing standpoint). Differently from [5], we have chosen to optimise the duty cycle directly, based on the performance goal that we want to achieve. While it is true that an optimisation based on intercontact times impacts directly on the delay performance, it is not straightforward how to control the one based on the other. With our model, instead, we can directly go from the requirements in term of probability of staying below a fixed delay threshold to the corresponding duty cycle value. In addition, differently from [5], we focus on a fixed duty cycle, similar to [6] [7] [4]. It is still an open research point which duty cycling strategy is to be preferred. However, preliminary results in [4] show that, under some assumptions, fixed duty cycle is the optimal strategy.

Another contribution focused on duty cycle optimisation is [8], in which Altman and Azad study the optimisation of node activation in DTN relying on a fluid approximation of the system dynamics. However, the problem analysed is different from the one studied in this paper, since in [8] nodes, once activated, remains active. In addition, this model is based on the assumption of i.i.d. intercontact times, while it has been shown that realistic intercontact times are intrinsically heterogeneous. For this reason, here we focus on heterogenous (but still independent) intercontact times.

## III. PRELIMINARIES

We assume that user mobile devices alternate between ON and OFF states, whose duration is fixed. We denote as duty cycle $\Delta$ the ratio between the duration of the ON and OFF states, and as $T$ their sum. We assume that when a node is in the ON state it is able to detect contacts with other nodes. Please refer to [2] for a discussion on how to apply this model to popular technologies such as Bluetooth and WiFi Direct. For the sake of simplicity, coarse synchronisation (e.g., controlled by the cellular infrastructure in the case of mobile data offloading) can be used to guarantee that ON intervals overlap between any pair of nodes, such that they can communicate during a contact if this overlaps with their ON phases. Under this assumption, in [2] we have investigated the effect of duty cycling on the detection of encounters between pairs of nodes. As discussed in Section I, this problem is extremely relevant to opportunistic networks, in which messages are delivered by means of consecutive exchanges between encountering nodes. In fact, the net effect of a duty cycling policy is to reduce the number of contacts that can be exploited for exchanging messages. More specifically, we have shown that, when intermeeting times follow an exponential distribution[1], the contact rate between a tagged node pair is approximately decreased by a factor $\Delta$. We summarise this result below.

*Proposition 1:* Considering a tagged pair of nodes $i$ and $j$ with exponential intercontact time of rate $\lambda_{ij}$, the detected intercontact time, i.e., the effective intercontact time when a duty cycling policy is in place, features approximately an exponential distribution with rate $\Delta\lambda_{ij}$, as long as $\lambda_{ij}T \ll 1$, where $T$ is the duty cycling period.

In [2] we have shown that the condition $\lambda_{ij}T \ll 1$ holds for the majority of contact traces available in the literature. Please note also that the above result has been obtained assuming that the duration of a contact is negligible with respect to the duration of the OFF period, which is reasonable (for example, results in [11] show that in absence of duty cycling the median contact duration is below 48s, while the period of typical duty cycling policies is in the order of several minutes).

Exploiting the result in Proposition 1, in our previous work [2] we have evaluated how intercontact times modified by the duty cycling policy affect the first two moments of the pairwise end-to-end delay for a set of representative (both social-oblivious and social-aware) opportunistic forwarding

---

[1]Exponential intercontact times are a popular assumption in the related literature [9] [10], even if a general consensus on the best probability distribution to approximate the realistic intercontact process has not been reached yet.

strategies. Specifically, we have derived the following properties, which we will use extensively throughout the paper:

**P1** The dependence of the coefficient of variation $c$ of the delay from $\Delta$ is negligible.

**P2** The expected delay when a duty cycling policy is in place (denoted as $E[D_\Delta]$) is approximately equal to the expected delay $E[D]$ with no duty cycle scaled by a factor $\frac{1}{\Delta}$, i.e, $E[D_\Delta] = \frac{E[D]}{\Delta}$.

**P3** The second moment of the delay when a duty cycling policy is in place (denoted as $E[D_\Delta^2]$) is approximately equal to the second moment of the delay $E[D^2]$ with no duty cycle scaled by a factor $\frac{1}{\Delta^2}$, i.e, $E[D_\Delta^2] = \frac{E[D^2]}{\Delta^2}$.

## IV. Setting the duty cycle for achieving a probabilistic guarantee on the delay

In this section we discuss how to derive the optimal duty cycle $\Delta_{opt}$ such that the delay of a tagged message remains, with a certain probability $p$, under a target fixed threshold $z$ or, in mathematical notation, $\Delta_{opt} = \min\{\Delta : P\{D_\Delta < z\} \geq p\}$. Since the delay increases with $\Delta$, the latter is equivalent to finding the solution to the following[2]:

$$\Delta_{opt} = \{\Delta : P\{D_\Delta < z\} = p\}. \tag{1}$$

Please note that in the following we will denote the CDF of $D_\Delta$ as $F_\Delta(x)$. In order to find the solution to Equation 1, the distribution of the delay $D_\Delta$ should be known. Although it is in general unfeasible to obtain an exact closed form for the distribution of $D_\Delta$ (except for some trivial cases, such as when the source node can only deliver the message to the destination directly), it is often possible to compute its moments, either exactly or approximately, under different distributions for intercontact times, as shown, e.g., in [9][12]. When the first two moments of the delay can be derived, it is possible to approximate its distribution with either a hypo-exponential or hyper-exponential random variable, using the moment matching approximation technique [3]. So, assuming that we have derived the first moment $E[D_\Delta]$ and the second moment $E[D_\Delta^2]$ of the delay using, e.g., the models in [9] [12], exploiting property P1, we can compute the coefficient of variation $c$ as $\sqrt{\frac{E[D^2]}{E[D]^2} - 1}$. Then, when $c$ is greater than one, $D_\Delta$ can be approximated using a 2-stages hyper-exponential distribution with the same moments of $D_\Delta$, as stated in the following Lemma.

*Lemma 1 (Hyper-exponential approximation):* The two moments matching approximation of $D_\Delta$ with coefficient of variation $c \geq 1$ is a 2-stages hyper-exponential distribution with parameters $(\lambda_1, p_1), (\lambda_2, p_2)$ given by the following:

$$\begin{cases} p_1 = \frac{1}{2}\left(1 + \sqrt{\frac{c^2-1}{c^2+1}}\right) \\ \lambda_1 = \frac{2p_1}{E[D_\Delta]} \end{cases} \quad \begin{cases} p_2 = 1 - p_1 \\ \lambda_2 = \frac{2p_2}{E[D_\Delta]} \end{cases} \tag{2}$$

Vice versa, when the coefficient of variation of the delay is smaller than 1 (but greater than $\frac{1}{\sqrt{2}}$ [13]), $D_\Delta$ can be approximated with an hypo-exponential distribution with CDF $F_X(x) = 1 - \frac{\mu_2}{\mu_2-\mu_1}e^{-\mu_1 x} + \frac{\mu_1}{\mu_2-\mu_1}e^{-\mu_2 x}$, for all $x \geq 0$, according to the following lemma.

---

[2]In the rest of the paper, for convenience of notation, we will drop subscript *opt* since all $\Delta$ we derive are the optimal ones.
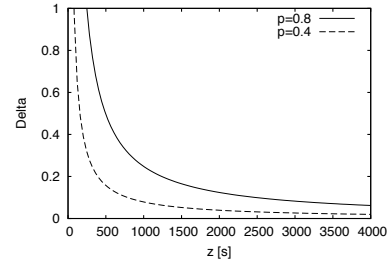


Fig. 1. $\Delta$ optimum for the exponential delay.

*Lemma 2 (Hypo-exponential approximation):* The two moments matching approximation of $D_\Delta$ with a coefficient of variation $c \in (\frac{1}{\sqrt{2}}, 1)$ is an hypo-exponential distribution with rates $\mu_1, \mu_2$ given by the following:

$$\begin{cases} \mu_1 = \frac{2}{E[D_\Delta]} \cdot \frac{1}{1+\sqrt{1+2(c^2-1)}} \\ \mu_2 = \frac{2}{E[D_\Delta]} \cdot \frac{1}{1-\sqrt{1+2(c^2-1)}} \end{cases} \tag{3}$$

In the rest of the section, we will analyse the optimisation problem in Equation 1 assuming that $D_\Delta$ features an exponential (Section IV-A), hyper-exponential (Section IV-B) or hypo-exponential distribution (Section IV-C). Please note that all three cases are possible starting from exponential intercontact times.

### A. The exponential case

The simplest case is when the delay features a coefficient of variation $c$ equal to one. In this hypothesis, the distribution of the delay is exponential with parameter $\lambda_\Delta = E[D_\Delta]^{-1}$. Then, it is straightforward to derive Theorem 1.

*Theorem 1:* The optimal duty cycle when $D_\Delta$ features an exponential distribution is given by the following:

$$\Delta = -\frac{\log(1-p)}{\lambda z}, \tag{4}$$

where we indicate with $\lambda$ the parameter of the exponential distribution obtained with $\Delta = 1$, i.e., $\lambda = E[D]^{-1}$.

*Proof:* We know that $\lambda_\Delta = \frac{1}{E[D_\Delta]}$, hence, since $E[D_\Delta] \sim \frac{E[D]}{\Delta}$ (Property P2), we have that $\lambda_\Delta = \lambda\Delta$. Thus, we can rewrite Equation 1 as $1 - e^{-\lambda\Delta z} = p$, from which $\Delta$ can be easily obtained. ∎

For the sake of example, in Figure 1 we plot $\Delta$ obtained from Theorem 1 setting $p = 0.8$. $E[D]$ is set to $154s$, which is the average expected delay obtained in [2] for a simple social-aware policy that selects the next relay of a message based on its contact rate with the destination and assuming the average contact rate equal to $4.07 \cdot 10^{-3}s^{-1}$ (the average contact rate measured in the RollerNet contact dataset [14]). Figure 1 shows that, as expected, when the target delay threshold is too small, it is impossible to achieve it with a probabilistic guarantee $p$, regardless of the value of the duty cycle. Instead, starting from $z = -\frac{\log(1-p)}{\lambda}$, $\Delta$ is inversely proportional to $z$.

Varying the parameters $z$ and $p$, Equation 1 describes a surface in $\mathbb{R}^3$, and more precisely the surface $K$ given by:

$$K = \{(z, p, \Delta) \in \mathbb{R} \times [0, 1] \times [0, 1] : P\{D_\Delta < z\} = p\}. \tag{5}$$
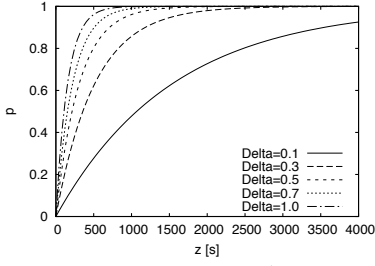
Fig. 2. Level set $K_\Delta$ for different values of $\Delta$ for the exponential delay

Given a certain duty cycle $\Delta \in (0,1]$, we can thus describe $K$ as the union of its level sets $K_\Delta$ or, in other terms, $K = \bigcup_{\Delta \in (0,1]} K_\Delta$ where:

$$K_\Delta = \{(z,p) \in \mathbb{R} \times [0,1] : P\{D_\Delta < z\} = p\}. \quad (6)$$

$K_\Delta$ is thus the set of pairs $(z,p)$ that can be obtained with a given duty cycling $\Delta$. It can be useful to plot $K_\Delta$ for different $\Delta$ in order to study whether it is possible to slightly compromise on the target performance in order to achieve a lower duty cycle. Assuming that we want $z = 250s$, in the exponential case (Figure 2) we can achieve it with a probability 0.8 with $\Delta = 1$ or with 0.68 with $\Delta = 0.7$, thus saving battery lifetime. Similarly, if we want to guarantee a target probability $p = 0.8$, with $\Delta = 1$ we obtain approximately $z = 250s$. If we are more flexible in terms of $z$, we can choose level set $K_{0.7}$ which gives $z = 350s$. This kind of analysis can be performed also for the hyper-exponential and hypo-exponential delays, with similar results.

### B. The hyper-exponential case

When the coefficient of variation of the delay is greater than one, the delay can be approximated with an hyper-exponential distribution as stated in Lemma 1. This means that Equation 1 becomes $1 - p_1 e^{-\lambda_1 z} - p_2 e^{-\lambda_2 z} = p$, where parameters $(\lambda_1, p_1), (\lambda_2, p_2)$ are given by Equation 2. From Equation 2, $\lambda_1$ and $\lambda_2$ depend on $\Delta$ (while $p_1$ and $p_2$ do not), thus, denoting with $\lambda_1^0$ and $\lambda_2^0$ the rates when $\Delta = 1$ and exploiting property P2, we can write Equation 1 as follows:

$$1 - p_1 e^{-\lambda_1^0 \Delta z} - p_2 e^{-\lambda_2^0 \Delta z} = p. \quad (7)$$

The exact solution $\Delta$ to this equation cannot be found analytically because Equation 7 cannot be inverted. However, in Theorem 2 below, we show how to obtain an approximated solution $\Delta_a$ that introduces a small error at most equal to $\varepsilon$.

*Theorem 2:* Let us $\lambda^0$ denote $E[D]^{-1}$ and $\lambda_1^0, \lambda_2^0$ the rates of the hyper-exponential delay (Equation 2) for $\Delta = 1$. When delay $D_\Delta$ has coefficient of variation greater than one, given a threshold $z$ of the delay and a target probability $p$, for every fixed $\varepsilon \geq \min\{\varepsilon_1, \varepsilon_2\}$ (whose definition is provided in the proof below), the duty cycle defined by:

$$\Delta_a = \begin{cases} \dfrac{1}{z}\left[-\dfrac{1-p-p_2}{\lambda_2^0 p_2} + \right. \\ \left. \quad + \dfrac{1}{\lambda_1^0}W\left(\dfrac{p_1^2}{p_2^2}e^{\frac{\lambda_1^0(1-p-p_2)}{\lambda_2^0 p_2}}\right)\right] & \text{if } \varepsilon_1 < \varepsilon_2 \\ \\ -\dfrac{\log 1-p}{\lambda^0 z} & \text{if } \varepsilon_1 \geq \varepsilon_2, \end{cases} \quad (8)$$

where $W$ is the Lambert function[3], verifies that $|F_{\Delta_a}(z) -$

$p| \leq \varepsilon$ and so it is a good approximation of the solution to Equation 7.

*Proof:* We will provide below an intuitive sketch of the proof whose detailed version can be found in [15]. The idea for finding an approximate solution to Equation 7 is to identify an approximation $\tilde{F}(z)$ that is close to $F_\Delta(z)$ under some conditions. So, we build a function $\tilde{F}$ for which it is possible to solve Equation 7 and for which $\min\{\varepsilon_1, \varepsilon_2\}$ is the error introduced (we will clarify this point below). Specifically, we have identified the following function:

$$\tilde{F}(z) = \begin{cases} 1 - p_1 e^{-\lambda_1^0 \Delta z} - p_2(1 - \lambda_2^0 \Delta z) & \text{if } \varepsilon_1 < \varepsilon_2 \\ 1 - e^{-\lambda^0 \Delta z} & \text{if } \varepsilon_1 \geq \varepsilon_2 \end{cases} \quad (9)$$

Let us denote with $\tilde{F}_1(z)$ and $\tilde{F}_2(z)$ the two parts of $\tilde{F}(z)$ in the above equation. In $\tilde{F}_1(z)$, we have approximated the third term on the left hand side of Equation 7 using the Taylor expansion, after noting that this term contributes to $F_\Delta(z)$ less and less as the coefficient of variation $c$ increases. Vice versa, the pure exponential behaviour ($\tilde{F}_2(z)$) dominates when $c$ is close to 1. Both $\tilde{F}_1(z)$ and $\tilde{F}_2(z)$ can be solved to find $\Delta$, from which Equation 8 follows.

The quality of these two approximations depends on the desired tolerance to the error that we inevitably introduce when we approximate $F_\Delta(z)$. If we tolerate a large error, either approximation can be chosen. Instead, if we want to achieve the smallest error, depending on the coefficient of variation of $D_\Delta$ we might have to prefer the one or the other. In the following we briefly discuss how to identify the minimum error introduced by $\tilde{F}_1(z)$ and $\tilde{F}_2(z)$, which we denote with $\varepsilon_1$ and $\varepsilon_2$ respectively. Let us start with $\tilde{F}_1(z)$. We want to find the region for which $|F_{\Delta_a}(z) - p| \leq \varepsilon$ or, equivalently, $|F_{\Delta_a}(z) - \tilde{F}_1^{(\Delta_a)}(z)| \leq \varepsilon$, where we denote with superscript $(\Delta_a)$ the fact that the CDF is computed using the approximated solution for $\Delta$. Solving the above inequality, we find that it holds for all $p < p_{max}$, where $p_{max}$ is a function of $c$ and $\varepsilon$ (due to lack of space, we do not report its formula here, please refer to [15] for details). Specifically, $p_{max}$ monotonically increases with $\varepsilon$. So, if we want to derive the minimum error for which inequality $|F_{\Delta_a}(z) - p| \leq \varepsilon$ holds for all $p$, we have to solve equation $p_{max}(c, \varepsilon) = 1$. We obtain the following:

$$\varepsilon_1 = \frac{(a-1)\left(-(a-1)W\left(\frac{(a+1)^2 e^{\frac{a+1}{a-1}}}{(a-1)^2}\right) + a + 1\right)^2}{4(a+1)^2}, \quad (10)$$

where again $W(x)$ denotes the Lambert function and $a$ is defined as $\sqrt{1 + 2(c^2 - 1)}$.

Let us now consider $\tilde{F}_2(z)$. We are able to prove that function $|F_{\Delta_a}(z) - p|$ has a maximum in $p^*$. We derive $p^*$ by finding the $p$ in which the derivative of $|F_{\Delta_a}(z) - p|$ becomes zero. Then, $\varepsilon_2$ can be computed as $\varepsilon_2 = |F_{\Delta_a}(z) - p^*|$, obtaining the following:

$$\varepsilon_2 = \frac{1}{2}(a+1)^{-2/a}\left(a\sqrt{2 - a^2} + 1\right)^{1/a} \cdot$$
$$\cdot \left(\frac{(a-1)(a+1)^2}{a\sqrt{2 - a^2} + 1} - \frac{a\sqrt{2 - a^2} + 1}{a+1} + 2\right), \quad (11)$$

---

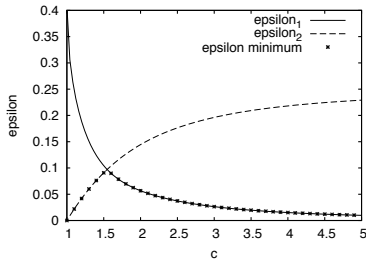[3]The Lambert function is defined as $W(x)e^{W(x)} = x$, for all $x \geq -\frac{1}{e}$

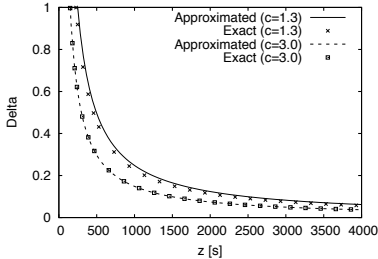Fig. 3. Error introduced by $\tilde{F}_1(z)$ and $\tilde{F}_2(z)$, varying $c$.



Fig. 4. $\Delta$ optimum (approximated vs exact) for the hyper-exponential delay with target probability $p = 0.8$ and varying $c$.

where again $a = \sqrt{1 + 2(c^2 - 1)}$. Thus, for both $\varepsilon_1$ and $\varepsilon_2$ we have derived a closed-form expression that tells us that the error that we make with our approximation is fixed for a given $c$. ∎

In Figure 3, we show how $\varepsilon_1$ and $\varepsilon_2$ vary with respect to the coefficient of variation $c$. As expected, for small $c$ (recall that we are in the hyper-exponential case, so $c > 1$ by definition) the exponential assumption $\tilde{F}_2$ allows us to achieve smaller errors. The opposite is true for large $c$. The worst case is reached for $c \sim 1.5$, when the minimum error is around $0.1$, which is still low. In Figure 4 we plot how the optimal duty cycle varies with $z$, setting the target probability to $p = 0.8$, for two values of coefficient of variation ($c = 1.3$ and $c = 3$). In both cases the approximation is good (the exact value is computed with standard numerical techniques to solve Equation 7). Specifically, when $c = 1.3$ the minimum error that can be achieved is $0.06$ and is provided by $\tilde{F}_2(z)$, hence confirming the predominance of the exponential behaviour for $c$ close to 1. Vice versa, when $c = 3$ the minimum error is $0.026$ and is provided by $\tilde{F}_1(z)$. It is also interesting to notice that smaller duty cycles can be achieved when $c$ increases, i.e., when the variability of the delay is higher. The importance of this result will be further discussed in Section V.

### C. The hypo-exponential case

When the coefficient of variation $c$ of the delay $D_\Delta$ is smaller than one, following Lemma 2, it is possible to approximate the delay with a hypo-exponential distribution. In particular, using property P2, if we denote with $\mu_1^0$ and $\mu_2^0$ the parameters obtained when $\Delta = 1$ in Equation 3, we can rewrite Equation 1 making explicit the dependence on $\Delta$:

$$1 - \frac{\mu_2^0}{\mu_2^0 - \mu_1^0}e^{-\mu_1^0 \Delta z} + \frac{\mu_1^0}{\mu_2^0 - \mu_1^0}e^{-\mu_2^0 \Delta z} = p. \quad (12)$$

As in the hyper-exponential case, this equation can not be directly inverted for finding $\Delta$, but it is possible to derive an
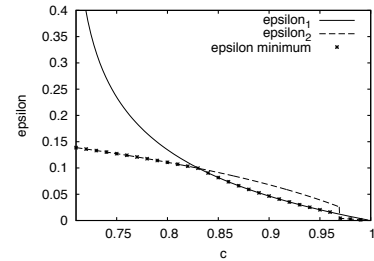


Fig. 5. Error introduced by $\tilde{F}_1(z)$ and $\tilde{F}_2(z)$, varying $c$.



(a) $c = 0.75$      (b) $c = 0.9$

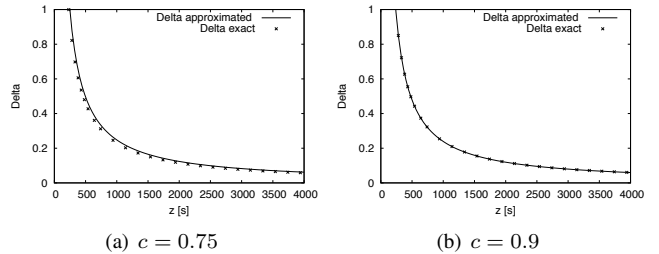Fig. 6. $\Delta$ optimum (approximated vs exact) for the hypo-exponential delay with target probability $p = 0.8$ and varying $c$.

approximate solution for which a small fixed (for a given $c$) error is introduced.

*Theorem 3:* Let $\mu_1^0$ and $\mu_2^0$ be the parameters given by Equation 3 with $\Delta = 1$. When the delay $D_\Delta$ has coefficient of variation smaller than one, the duty cycle defined by:

$$\Delta_a = \begin{cases} -\frac{1}{\mu_1^0 z} \log\left[(1-p) \cdot \frac{\mu_2^0 - \mu_1^0}{\mu_2^0}\right] & \text{if } \varepsilon_1 < \varepsilon_2 \\ -\frac{\log \frac{1-p}{\lambda^0 z}}{\lambda^0 z} & \text{if } \varepsilon_1 \geq \varepsilon_2, \end{cases} \quad (13)$$

verifies that $|F_{\Delta_a}(z) - p| \leq \varepsilon$ (with $\varepsilon \geq \min\{\varepsilon_1, \varepsilon_2\}$, see the proof in [15]), and so it is a good approximation of the solution to Equation 12.

Due to lack of space and since the rationale follows that of the proof for Theorem 2, we omit the proof of the above theorem, which can however be found in [15].

In Figure 5 we plot $\varepsilon_1$ and $\varepsilon_2$ varying $c$. When $c$ is close to one, both approximations are very good. For values of $c$ roughly in the interval $(0.83, 0.97)$, $\tilde{F}_1(z)$ provides better results, while, for low values of $c$, $\tilde{F}_2(z)$ is to be preferred. In Figures 6(a) and 6(b) we show how the optimal duty cycle varies with $z$, setting the target probability to $p = 0.8$, for two values of coefficient of variation ($c = 0.75$ and $c = 0.9$, respectively). In both cases the approximation and the exact value are very close. In Figure 6(a) the minimum error that can be achieved is $0.13$ and is provided by $\tilde{F}_2(z)$, while in Figure 6(b) the minimum error is $0.05$ and is provided by $\tilde{F}_1(z)$.

### V. OPTIMAL DUTY CYCLE AND TRAFFIC GAIN

In this section we investigate how the choice of the optimal duty cycle affects the volume of traffic carried by the network. As already discussed, the advantage of implementing a duty cycle policy is that device batteries are preserved and, as a consequence, the lifetime of the network increases. Specifically, with a duty cycle $\Delta$ and a baseline network lifetime $L$ (i.e., with $\Delta = 1$), the network lifetime when a duty cycling
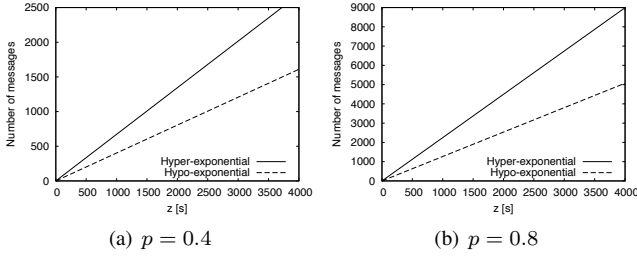
Fig. 7. $\mathcal{N}(\Delta)$ varying $z$ for different $p$ in the case of hyper-exponential ($c = 3$) and hypo-exponential ($c = 0.75$) delay.

policy is in place is given by $\frac{L}{\Delta}$. A longer network lifetime is very useful because it allows nodes to exchange messages for a longer time. If we assume, similarly to [2], that messages are generated according to a Poisson process with rate $\eta$, we can derive how the total number of messages $\mathcal{N}$ *delivered* by the nodes varies with $\Delta$. Due to lack of space, in the following we only consider the hypo-exponential and hyper-exponential cases. First, in the theorem below we recall the main results for $\mathcal{N}$ derived in [2].

*Theorem 4:* If the delay $D_\Delta$ has coefficient of variation $c$ greater than one, the volume $\mathcal{N}(\Delta)$ of messages delivered by the system under duty cycling $\Delta$ is given by:

$$\mathcal{N}(\Delta) = \frac{\eta L}{\Delta} - \eta E[D_\Delta]\left[1 - \tfrac{1}{2}e^{\frac{-L}{E[D_\Delta]\Delta}}\left(e^{\left(1+\sqrt{\frac{c^2-1}{c^2+1}}\right)} + e^{\left(1-\sqrt{\frac{c^2-1}{c^2+1}}\right)}\right)\right]. \tag{14}$$

Instead, if the delay $D_\Delta$ has coefficient of variation $c$ smaller than one, the volume $\mathcal{N}(\Delta)$ of messages delivered by the system under duty cycling $\Delta$ is given by:

$$\mathcal{N}(\Delta) = \frac{\eta L}{\Delta} - \eta E[D_\Delta]\cdot$$

$$\left[1 - \frac{1}{4\sqrt{1+2(c^2-1)}}\left(\left(1+\sqrt{1+2(c^2-1)}\right)^2 e^{-\left(\frac{2L}{\Delta E[D_\Delta]\left(1+\sqrt{1+2(c^2-1)}\right)}\right)} \right.\right.$$
$$\left.\left. -\left(1-\sqrt{1+2(c^2-1)}\right)^2 e^{-\left(\frac{2L}{\Delta E[D_\Delta](1-\sqrt{1+2(c^2-1)})}\right)}\right)\right]. \tag{15}$$

If we substitute in the above equations the optimal $\Delta$ derived in the previous section, we obtain how $\mathcal{N}$ varies as a function of the target performance $(z, p)$. In order to study this dependence, we set the network lifetime $L$ to $60000s$ and we assume that each node generates one message every ten minutes ($\eta = \frac{1}{600}s^{-1}$). In Figure 7(a) we set $p$ to the value 0.8 and we plot $\mathcal{N}$ varying $z$, while in Figure 7(b) we set $p = 0.4$. Besides the expected result that the less stringent the performance requirements (i.e., higher $p$) the higher the volume of traffic (because smaller duty cycles can be used), we observe an interesting difference between the two delay distributions. The traffic delivered under hyper-exponential delays is always higher than that exchanged under hypo-exponential delays. This is due to the fact that, as we have seen in Section IV-B, when $c$ increases we can achieve smaller optimal duty cycle for a given target performance $(z, p)$, hence saving more energy and increasing the lifetime of the network.

## VI. CONCLUSION

In this work we have studied how to optimise the duty cycle in order to guarantee, with probability $p$, that the delay of messages remains below a threshold $z$, assuming that inter-contact times are exponentially distributed. We have provided an exact solution for the case in which the delay follows an exponential distribution, and approximated solutions for the cases in which the coefficient of variation of the delay is greater than or smaller than 1. We have also demonstrated that the approximation of $\Delta$ introduces an error $\varepsilon$ whose formula we have provided and that is small and fixed for a given coefficient of variation $c$ of the delay. Finally we have focused on the volume of traffic delivered by the network when the optimal duty cycle is implemented, and we have discussed how the two parameters $z$ and $p$ impact on the number of messages delivered. Specifically, we have shown that the optimisation of the duty cycle is more efficient with hyper-exponential delays, as it achieves lower duty cycles and thus provides higher energy gains.

## REFERENCES

[1] J. Whitbeck, Y. Lopez, J. Leguay, V. Conan, and M. D. De Amorim, "Push-and-track: Saving infrastructure bandwidth through opportunistic forwarding," *Perv. and Mob. Comp.*, vol. 8, no. 5, pp. 682–697, 2012.

[2] E. Biondi, C. Boldrini, A. Passarella, and M. Conti, "Duty cycling in opportunistic networks: intercontact times and energy-delay tradeoff," IIT-CNR 22-2013, Tech. Rep. 22/2013, http://cnd.iit.cnr.it/chiara/pub/techrep/biondi2013duty_tr.pdf.

[3] H. Tijms, *A First Course in Stochastic Models*. Wiley, 2003.

[4] W. Wang, V. Srinivasan, and M. Motani, "Adaptive contact probing mechanisms for delay tolerant applications," in *MobiCom*. ACM, 2007, pp. 230–241.

[5] W. Gao and Q. Li, "Wakeup scheduling for energy-efficient communication in opportunistic mobile networks," in *IEEE INFOCOM*, 2013.

[6] H. Zhou, J. Chen, H. Zhao, W. Gao, and P. Cheng, "On exploiting contact patterns for data forwarding in duty-cycle opportunistic mobile networks," *IEEE Trans. on Vehic. Tech.*, pp. 1–1, 2013.

[7] O. Trullols-Cruces, J. Morillo-Pozo, J. M. Barcelo-Ordinas, and J. Garcia-Vidal, "Power saving trade-offs in delay/disruptive tolerant networks," in *WoWMoM*. IEEE, 2011, pp. 1–9.

[8] E. Altman, A. Azad, T. Başar, and F. De Pellegrini, "Combined optimal control of activation and transmission in delay-tolerant networks," *IEEE/ACM Trans. on Netw.*, vol. 21, no. 2, pp. 482–494, 2013.

[9] A. Picu, T. Spyropoulos, and T. Hossmann, "An analysis of the information spreading delay in heterogeneous mobility dtns," in *IEEE WoWMoM*, 2012, pp. 1–10.

[10] W. Gao and G. Cao, "User-centric data dissemination in disruption tolerant networks," in *IEEE INFOCOM*, 2011, pp. 3119–3127.

[11] S. Gaito, E. Pagani, and G. P. Rossi, "Strangers help friends to communicate in opportunistic networks," *Computer Networks*, vol. 55, no. 2, pp. 374–385, 2011.

[12] C. Boldrini, M. Conti, and A. Passarella, "Performance modelling of opportunistic forwarding under heterogenous mobility," IIT-CNR, Tech. Rep. TR-12/2013, http://cnd.iit.cnr.it/chiara/pub/techrep/boldrini2013heterogenous_tr.pdf.

[13] G. Bolch, S. Greiner, H. de Meer, and K. Trivedi, *Queueing Networks and Markov Chains: Modeling and Performance Evaluation with Computer Science Applications*. Wiley, 2006.

[14] P.-U. Tournoux, J. Leguay, F. Benbadis, J. Whitbeck, V. Conan, and M. D. de Amorim, "Density-aware routing in highly dynamic dtns: The rollernet case," *IEEE Trans. on Mob. Comp.*, vol. 10, no. 12, pp. 1755–1768, 2011.

[15] E. Biondi, C. Boldrini, A. Passarella, and M. Conti, "Optimal duty cycling in mobile opportunistic networks with end-to-end delay guarantees," IIT-CNR 2014, Tech. Rep., http://cnd.iit.cnr.it/chiara/pub/techrep/biondi2014optimisation_tr.pdf.

# Robust Adaptive Modulation and Coding (AMC) Selection in LTE Systems using Reinforcement Learning

Raffaele Bruno, Antonino Masaracchia, Andrea Passarella

Institute of Informatics and Telematics (IIT)

Italian National Research Council (CNR)

Via G. Moruzzi 1, Pisa, ITALY

E-mail: {r.bruno, a.masaracchia,a.passarella}@iit.cnr.it

*Abstract*—Adaptive Modulation and Coding (AMC) in LTE networks is commonly employed to improve system throughput by ensuring more reliable transmissions. Most of existing AMC methods select the modulation and coding scheme (MCS) using pre-computed mappings between MCS indexes and channel quality indicator (CQI) feedbacks that are periodically sent by the receivers. However, the effectiveness of this approach heavily depends on the assumed channel model. In addition CQI feedback delays may cause throughput losses. In this paper we design a new AMC scheme that exploits a reinforcement learning algorithm to adjust at run-time the MCS selection rules based on the knowledge of the effect of previous AMC decisions. The salient features of our proposed solution are: $i$) the low-dimensional space that the learner has to explore, and $ii$) the use of direct link throughput measurements to guide the decision process. Simulation results obtained using ns3 demonstrate the robustness of our AMC scheme that is capable of discovering the best MCS even if the CQI feedback provides a poor prediction of the channel performance.

*Index Terms*—LTE, channel quality, adaptive modulation and coding (AMC), reinforcement learning, performance evaluation.

## I. INTRODUCTION

The Long Term Evolution (LTE) is an acronym that refers to a series of cellular standards developed by 3GPP to meet the requirements of 4G systems. In particular, LTE has been designed to provide high data rates, low latency, and an improved spectral efficiency compared to previous cellular systems. To achieve these goals LTE adopts advanced physical layer technologies, such as OFDMA and multi-antenna techniques, and it supports new Radio Resource Management (RRM) functions for link adaptation [1]. In particular, adaptive modulation and coding (AMC) has been proposed for LTE, as well as many other wireless communication systems, to increase channel throughput [2]. In general, AMC techniques try to optimally select the channel coding and modulation scheme (MCS), while fulfilling a certain Block Error Rate (BLER) constraint[1] by taking into account the current channel conditions and the

---

[1]The BLER for a certain user is defined as the ratio between the number of erroneous resource blocks and the total number of resource blocks received by that user. In the LTE standard it is mandated that the selected MCS ensures an average BLER under the measured channel conditions lower than 10% [3].

receiver's characteristics (e.g., antenna configuration). For LTE downlink transmissions, traditional AMC schemes rely on the channel quality indicator (CQI) feedbacks that are periodically reported by the user terminals (UEs) to their base stations (eNBs) [3]. How CQI values should be computed by the UE using channel state information (e.g., SINR measurements) is implementation dependent. In principle, an eNB can use other information in addition to the CQI values reported by UEs, such as HARQ retransmissions, to determine the selected MCS. In practical implementations - as better explained in Section II - the UEs directly selects the MCS value that, if used by the eNB under the measured channel conditions, would achieve the maximum possible throughput by guaranteeing that the BLER is below 10%. This value is then mapped onto a CQI value and fed back to the eNB (that translates it back into the corresponding MCS value) [4], [5]. Therefore, the key focus of AMC algorithms is to define how UEs can compute MCS values that satisfy the BLER requirements.

Several technical challenges have to be addressed to design efficient AMC solutions for LTE systems. In particular, in practical LTE systems, the SINR values of multiple subcarriers are aggregated and translated into a one-dimensional link quality metric (LQM), since the same MCS must be assigned to all subcarriers assigned to each UE. Popular methods that are used in LTE to obtain a single effective SINR from a vector of physical-layer measurements related to subcarriers are the exponential effective SINR mapping (EESM) [6] or the mean mutual information per coded bit (MMIB) [7]. Once the LQM is found, AMC schemes typically exploit *static mappings* between these link quality metrics and the BLER performance of each MCS to select the best MCS (in terms of link throughput). In other words, for each MCS a range of LQM values is associated via a look-up table, over which that MCS maximises link throughput. Either link-level simulations or mathematical models can be used to generate such static BLER curves under a specific channel model. Unfortunately, past research has shown that it is difficult to derive accurate link performance predictors under realistic channel assumptions [5], [8]–[10]. Furthermore, a simulation-based approach to derive the mapping between LQM values

and BLER performance is not scalable since it is not feasible to exhaustively analyse all possible channel types or several possible sets of parameters [11]. The second main problem with table-based AMC solutions is that a delay of several transmission time intervals (TTIs) may exist between the time when a CQI report is generated and the time when that CQI feedback is used for channel adaptation. This is due to processing times but also to the need of increasing reporting frequency to reduce signalling overheads. This mismatch between the current channel state and its CQI representation, known as *CQI ageing*, can negatively affect the efficiency of AMC decisions [12], [13]

To deal with the above issues, in this paper we propose a new flexible AMC framework, called RL-AMC, that autonomously and at run-time decides upon the best MCS (in terms of maximum link-layer throughput) based on the knowledge of the outcomes of previous AMC decisions. To this end we exploit reinforcement learning techniques to allow each eNB to update its MCS selection rules taking into account past observations of achieved link-layer throughputs. Specifically, the purpose of the decision-making agent in our AMC scheme is to discover which is the correction factor that should be applied to CQI feedbacks in order to guide the transmitters in selecting more efficient MCSs. An important feature of our proposed scheme is the use of a low-dimensional state space, which ensures a robust and efficient learning even under time-varying channel conditions and mobility. Through simulations in ns3 we show that our AMC method can improve the LTE system throughput compared to other schemes that use static mappings between SINR and MCS both under pedestrian and vehicular network scenarios. Furthermore, our AMC is capable of discovering the best MCS even if the CQI feedback provides a poor prediction of the channel performance.

Before presenting our solution, it is important to point out that other studies [14]–[17] have proposed to use machine learning techniques to improve AMC in wireless systems. The main weakness of most of these solutions is to rely on machine learning algorithms (e.g., pattern classification [15] or SVM [14], [16]) that require large sets of training samples to build a model of the wireless channel dynamics. Similar to our work, the AMC scheme proposed in [17] exploits Q-learning algorithms to avoid the use of model-training phases. However, the MCS selection problem in [17] is defined over a continuous state space (i.e., received SINR), and even after discretisation a large number of states must be handled by the learning algorithm.

The remaining of this paper is organised as follows. Section II overviews existing proposals to implement AMC techniques in LTE networks. Section III introduces the principles of reinforcement learning, and introduces the Q-learning algorithm. Section IV describes our RL-AMC scheme. In Section V we report simulation results to demonstrate the performance improvements of the proposed scheme. Section VI concludes the paper with final remarks.
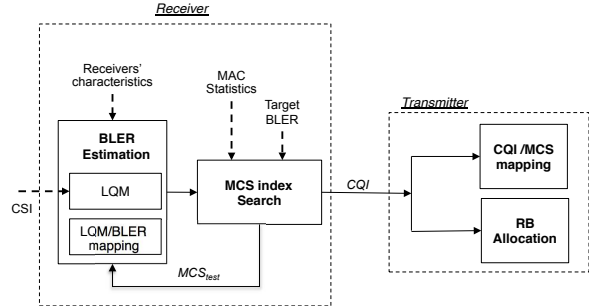


Fig. 1. AMC functional architecture.

## II. AMC IN LTE

For the sake of illustrative purposes, in Figure 1 we show a functional architecture for a practical AMC scheme for LTE systems. At the receiver's side, a first module is responsible for processing the channel state information (e.g., per-subcarrier received SINR values) to obtain a BLER estimation under the assumption of a specific channel model. Specifically, the receiver maps the channel measurements into a single link quality metric. Then, an offline look-up table is used to map this LQM to a BLER estimate for each MCS. These BLER curves are used to find the highest-rate MCS index that can satisfy a 10% BLER target. Finally, the selected MCS index is sent in the form of a CQI feedback to the transmitter. Based on such CQI feedback the transmitter performs resource scheduling and MCS selection.

Most of existing research on AMC schemes for LTE is focused on the problem of CQI calculation given a link quality metric. As mentioned in Section I a popular and sufficiently accurate method for LQM calculation is EESM. For instance, the authors in [18] study the MCS performance under an AWGN channel. Accurate packet error prediction for link adaptation via a Gaussian approximation of coding and decoding performance is proposed in [19]. A novel LQM metric for link adaptation based on raw bit-error-rate, effective SINR and mutual information is investigated in [20]. In [4] the authors proposed MCS selection based on packet-level effective SINR estimates rather than block-level SINR values, and they describe different averaging schemes to map BLER onto packet error rates. On the other hand, the authors in [5], [21] develops statistical models of the EESM under different channel models and use those models to analyse the throughput of EESM-based AMC for various CQI feedback schemes. A second group of paper studies channel predictors to deal with the CQI ageing. The authors in [12] derive closed-form expressions for the average throughput of an adaptive OFDMA system under the assumption of imperfect CQI knowledge. The performance of different CQI predictors, such as Kalman filtering or linear prediction with stochastic approximation, are evaluated in [13] and [22].

## III. BACKGROUND ON REINFORCEMENT LEARNING (RL)

Reinforcement Learning (RL) is a popular machine learning technique, which allows an agent to automatically determine the optimal behaviour to achieve a specific goal based on the positive or negative feedbacks it receives from the environment in which it operates after taking an action from a known set of admissible actions [23]. Typically, reinforcement learning problems are instances of the more general class of Markov Decision Processes (MDPs), which are formally defined through:

- a finite set $S = \{s_1, s_2, \ldots, s_n\}$ of the $n$ possible states in which the environment can be;
- a finite set $A(t) = \{a_1(t), a_2(t), \ldots, a_m(t)\}$ of the $m$ admissible actions that the agent may perform at time $t$;
- a transition matrix $P$ over the space $S$. The element $P(s, a, s')$ of the matrix provides the probability of making a transition to state $s' \in S$ when taking action $a \in A$ in state $s \in S$;
- a reward function $R$ that maps a state-action pair to a scalar value $r$, which represents the immediate payoff of taking action $a \in A$ in state $s \in S$.

The goal of a MDP is to find a *policy* $\pi$ for the decision agent, i.e., a function that specifies the action that the agent should choose when in state $s \in S$ to maximise its expected long-term reward. More formally, if an agent follows a policy $\pi$ starting from a certain state $s$ at time $t$ the policy value over an infinite time horizon, also called the value-state function, is simply given by

$$V^\pi(s) = \sum_{k=0}^{\infty} \gamma^k r_{t+k} , \qquad (1)$$

where $\gamma \in [0, 1]$ is a *discount factor* that weights future rewards. Then an *optimal* policy $\pi^*$ is, by definition, the one that maximise the value-state function. As a consequence, the policy that ensures the maximum possible expected reward, say $V^*(s)$, could be obtained by solving an optimisation problem $V^*(s) = \max_\pi V^\pi(s)$. If the transition matrix is known such optimisation problem can be expressed using a system of nonlinear equations by using techniques such as dynamic programming [23]. However, in most practical conditions it is hard, if not even impossible, to acquire such complete knowledge of the environment behaviour. In this case there are model-free learning methods that continuously update the probabilities to perform an action in a certain state by exploiting the observed rewards. Such methods adopt an alternative characterisation of policy goodness based on the state-action value function, or Q-function. Formally, the function $Q^\pi(s, a)$ computes the expected reward of taking an action $a$ in a starting state $s$ and then following the policy $\pi$ hereafter. Owing to the Bellman's optimality principle, it holds that a greedy policy (i.e., a policy that at each state selects the action with the largest Q-value) is the optimal policy. In other words, it holds that $V^*(s) = \max_{a \in A} Q^*(s, a)$ with $Q^*(s, a) = \max_\pi Q(s, a)$.

In this work we use a model-free solving technique for reinforcement learning problems known as *Q-learning* [24], which constructs the optimal policy by iteratively selecting the action with the highest value in each state. The core of this algorithm is an iterative value update rule that each time the agent selects an action and observes a reward makes a correction of the old Q-value for that state based on the new information. This updating rule is given by:

$$Q(s, a) = Q(s, a) + \alpha \left[ r(s, a) + \gamma \max_{a'} Q(s', a') - Q(s, a) \right] , \qquad (2)$$

where $\alpha \in [0, 1]$ is the learning rate. Basically, the $\alpha$ parameter determines the weight of the newly acquired information over state-action value information. In our AMC framework we use $\alpha = 0.5$.

The advantage of Q-learning is that it is guaranteed to converge to the optimal policy. On the negative side, the convergence speed may be slow if the state space is large due to the *exploration vs. exploitation dilemma* [23]. Basically, when in state $s$ the learning agent should exploit its accumulated knowledge of the best policy to obtain high rewards, but it must also explore actions that it has not selected before to find out a better strategy. To deal with this issue, various exploration strategies have been proposed in the literature, ranging from simple greedy methods to more sophisticated stochastic techniques, which assign a probabilistic value for each action $a$ in state $s$ according to the current estimation of $Q(s, a)$. In Section IV we discuss more in detail such exploration strategies.

## IV. AN RL-BASED AMC SCHEME (RL-AMC)

In order to apply the Q-learning approach to the MCS selection problem it is necessary to define: $i$) the state space of the problem, $ii$) the feedbacks that the decision agent receives from the LTE network, and $iii$) the admissible actions for the agent with the action selection strategy. In our RL-based AMC framework, the problem state consists of CQI feedbacks and their evolution trends. The reward is the instantaneous link throughput obtained by a user after each transmission. Finally, an action is the selection of a correction factor to be applied to each CQI feedback to identify the best MCS under the current channel conditions. In the following, we describe in details the operations of our proposed AMC algorithm.

First of all, it is important to clarify that the AMC decision agent interacts with the environment (i.e., the LTE network) at discrete time instants, called epochs. At each epoch the agent receives some representation of the LTE channel state and on that basis selects an action. In the subsequent epoch the agent receives a reward, and finds itself in a new state. In our AMC framework we assume that an epoch is the time when the UE receives a segment of data, either new or retransmitted. Without loss of generality we also assume that the decision agent is provided with a mapping rule that establishes a relationship between SINR values and MCS indexes. Note that our solution is not restricted to any specific BLER models but *an initial MCS value is only needed to bootstrap the*

*learning process* and to reduce the size of the state space. Thus, it is not necessary that this mapping is accurate nor adjusted to the unique characteristics of each communication channel. In Section V we will investigate the robustness of our AMC scheme to inaccurate CQI representation of channel performance.

Intuitively, a straightforward approach to define the state of the MCS selection problem would be to use the SINR values of received segments of data[2] as state variables, as in [17]. However, the SINR is a continuos variable and it should be discretised to be compatible with a discrete MDP formulation. The main drawback is that a fine discretisation leads to a large-dimensional state space, which increases convergence and exploration times. To avoid this problem, we directly use CQI-based metrics for the state representation. Specifically, we adopt a two-dimensional space $S = \{s_1, s_2\}$ to characterise the LTE communication channel. The first state variable represents the CQI value (called $CQI^m$) that the UE should select using the internal look-up table that associates BLER and MCS and received SINR. The second state variable represents the $\Delta CQI^m$ value, which is defined as the difference between the last two consecutive $CQI^m$ estimates. In other words, $\Delta CQI^m$ provides a rough indication of the trend in channel quality evolution. For instance, $\Delta CQI^m < 0$ implies that the channel quality is temporarily degrading.

Since the objective of the MCS selection procedure should be to maximise the link throughput it is a natural choice to define the reward function as the instantaneous link-layer throughput achieved when taking action $a$ (i.e., applying a correction factor to current CQI value taken from the mapping function) when in state $s$ (i.e., given the pair $\{CQI_t^m, \Delta CQI_t^m\}$). More precisely, we assume that the reward value of an erroneous downlink transmission is null. On the other hand, the reward for a successful downlink transmission is given by

$$R(s_{t_1}, a_{t_1}) = \frac{TB}{\#TTIs \text{ in } [t_1, t_2]} , \qquad (3)$$

where $TB$ is the MAC transport block size (i.e., the number of useful bits that could be carried in a certain number of RBs with a certain MCS), while the denominator is the time between the time $t_1$ when that segment of data was first scheduled and the time $t_2$ when it was successfully received[3].

The core of the Q-learning algorithm is represented by the set $A$ of admissible actions. In our learning model we assume that an action consists of applying a correction factor to the CQI value that is initially estimated by means of the internal look-up table. As discussed above, the mapping relationship between SINR values and MCS may be inaccurate and the correction factor allows the agent to identify the best

---

[2]We recall that LTE physical layer relies on the concept of resource blocks. A segment of data or transport block is basically a group of resource blocks with a common MCS that are allocated to a user. Typically, a packet coming from the upper layers of the protocol stack will be transmitted using multiple segments of data.

[3]A segment of data that is discarded after a maximum number of retransmissions has also a null reward.

modulation and coding scheme (in the sense of maximising the link throughput) for the given channel conditions. For instance, it may happen that the SINR-to-MCS mapping is too conservative for the current channel conditions and an MCS with an higher data rate can be used without violating the target BLER requirement. In this case the correction factor should be positive. Furthermore, a correction factor is also needed to compensate eventual errors due to CQI feedback delay. More formally, we assume that an action taken by the AMC decision agent at time $t$ is one possible choice of an integer number in the set $(-k, \ldots, -2, -1, 0, 1, 2, \ldots k)$, that we denote as $a_t$ in the following. This index is added to the original $CQI^m$ value to compute the CQI to be sent to the eNB, denoted as $CQI^f$. The line of reasoning for this adjustment is as follows. Let us assume that the agent state at time $t$ is $\{CQI_t^m, \Delta CQI_t\}$. We argue that if $\Delta CQI_t < 0$ we should prefer conservative MCS selections (and thus use values of $a_t$ lower than 0) because the channel trend is negative, while if $\Delta CQI_t \geq 0$ we can try to use MCSs offering higher data rates (and thus positive values for $a_t$). Recalling that the CQI is an integer between 0 and 15 [3], this can be expressed by writing that the CQI feedback, say $CQI_t^f$, that should be sent to the eNB by the UE to guide the selection of the MCS index for downlink transmissions at next epoch $t+1$ should be

$$CQI_t^f = \max\left[0, \min\left[CQI_t^m + a_t, 15\right]\right] , \qquad (4)$$

where $a \in [0, 1, 2, \ldots k]$ if $\Delta CQI_t \geq 0$ and $a \in [-k, \ldots, -2, -1, 0]$ otherwise. Thus, the set of admissible actions is different whether the channel-quality trend is negative or non-negative. Before proceeding it is useful to point out that the choice of the $k$ value determines how aggressively we want to explore the problem state space. In general, the selection of the $k$ value could take into account the CQI difference statistics, i.e., to what extent a current CQI may be different from the reported CQI after a feedback delay [10]. In Section V-C we will discuss this aspect more in detail.

A very important learning procedure is the action selection rule, i.e., the policy used to decide which specific action to select in the set of admissible actions. As discussed in Section III there is a tradeoff between exploitation (i.e., to select the action with the highest Q-value for the current channel state) and exploration (i.e., to select an action randomly). The simplest approach (called $\epsilon$-greedy [23]) would be to use a fixed probability $\epsilon$ to decide whether to exploit or explore. A more flexible policy (called *softmax* action-selection rule [23]) is to assign a probability to each action, basing on the current Q-value for that action. The most common softmax function used in reinforcement learning to convert Q-values into action probabilities $\pi(s, a)$ is the following [23]:

$$\pi(s, a) = \frac{e^{Q(s,a)/\tau}}{\sum_{a' \in \Omega_t} e^{Q(s,a')/\tau}} , \qquad (5)$$

where $\Omega_t$ is the set of admissible actions at time $t$. Note that for high $\tau$ values the actions tend to be all (nearly) equiprobable. On the other hand, if $\tau \to 0$ the softmax policy becomes the

same as a merely greedy action selection. In our experiments we have chosen $\tau = 0.5$.

## V. PERFORMANCE EVALUATION

In this section, we assess the performance of our proposed RL-AMC scheme in two different scenarios. In the first one a fixed CQI is fed back to the eNB by each UE. Without the use of reinforcement learning AMC necessarily selects a fixed MCS independently of the current channel conditions. Then, we demonstrate that our RL-based AMC is able to converge towards the best MCS even if the initial CQI estimate are totally wrong. In the second scenario we compare RL-AMC against the solution described in [25], which exploits spectral efficiency estimates to select MCS. Specifically, the spectral efficiency of user $i$ is approximated by $\log_2(1 + \gamma_i/\Gamma)$, where $\gamma_i$ is the effective SINR of user $i$ and $\Gamma$ is a scaling factor. Then, the mapping defined in the LTE standard [26] is used to convert spectral efficiency into MCS indexes and, then, into CQI feedbacks. In this case, we show that our reinforcement learning algorithm is able to improve the accuracy of the CQI mapping at run time.

### A. Simulation setup

All the following experiments have been carried out using the ns3 packet-level simulator, which includes a detailed implementation of the LTE radio protocol stack. As propagation environment, we assume an *Urban Macro* scenario, where path loss and shadowing are modelled according to the COST231-Hata model [27], which is widely accepted in the 3GPP community. The fast fading model is implemented using the Jakes model for Rayleigh fading [28]. To limit the computation complexity of the simulator pre-calculated fading traces are included in the LTE model that are based on the standard multipath delay profiles defined in [29]. In the following tests we have used the *Extended Typical Urban* fading propagation model with pedestrian (3 km/h) and vehicular (30 km/h) users' speeds. The main LTE physical parameters are summarised in Table I. Regarding the network topology, the considered scenario is composed by a single cell and a number of users, chosen in the range $[10, 100]$, which move according a Random Waypoint Model (RWM) [30] within the cell, if not otherwise stated. A downlink flow, modelled with an infinite buffer source, is assumed to be active for each UE. Finally, the eNode B adopts the resource allocation type 0, thus only allocating resource block groups (RBGs) to scheduled UEs. Given the downlink system bandwidth (see Table I) a RBG comprises two RBs [3]. RBGs are assigned to UEs following a Round Robin (RR) scheduler that divides equally the available RBGs to active flows. Then, all the RBs in the allocated RBGs used the MCS index that is signalled in the last received CQI feedback. Furthermore, the implemented version of RR algorithm is not adaptive, which implies that it maintains the same RBGs and MCS index when allocating retransmission attempts.

All results presented in the following graphs are averaged over five simulation runs with different network topologies.

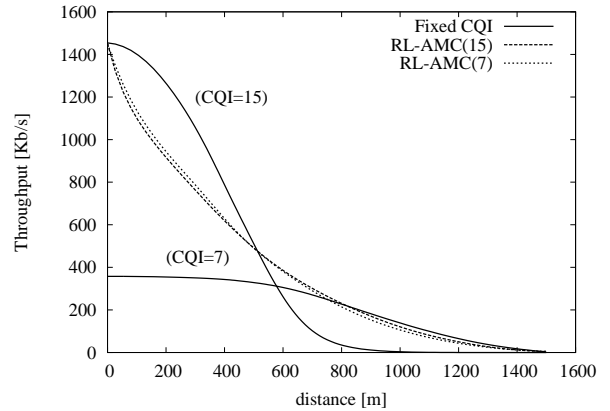| Parameter | Value |
|---|---|
| Carrier frequency | 2GHz |
| Bandwidth for downlink | 5 MHz |
| eNB power transmission | 43 dBm |
| Subcarrier for RB | 12 |
| SubFrame length | 1 ms |
| Subcarrier spacing | 15 KHz |
| Symbols for TTI | 14 |
| PDCCH & PCFICH (control ch.) | 3 symbols |
| PDSCH (data ch.) | 11 symbols |
| CQI reporting | periodic wideband |
| CQI processing time | 2 TTIs |
| CQI transmission delay | 4 TTIs |



Fig. 2. Average throughput as a function of the distance of the tagged user from the eNB in a pedestrian scenario.

Confidence intervals are very tight and are not shown in the figures. Each simulation run lasts 150 seconds.

### B. Results for fixed CQI

In this first set of simulations we assume that ten UEs are randomly deployed in the cell and they are static. Then an additional tagged user is moving with pedestrian speed from the center of the cell to its boundaries. However, independently of the UE position the CQI feedback is constant. Then, Figure 3 shows a comparison of the throughput achieved by the tagged user with and without reinforcement learning. This is obviously a limiting case which is analysed to assess the robustness of our RL-AMC scheme even when CQI provides a very poor prediction of channel performance. As expected with fixed MCS the user throughput is constant when the MCS is over provisioned, while it rapidly goes to zero after a critical distance. On the contrary, our RL-AMC is able to discover the correction factor that should be applied to the initial CQI to force the selection of a more efficient MCS. In addition, the performance of RL-AMC are almost independent of the initial CQI value. Note that in this case RL-AMC must explore the full range of CQI values and we set $k$ in (4) equal to 15.

### C. Results with adaptive CQI

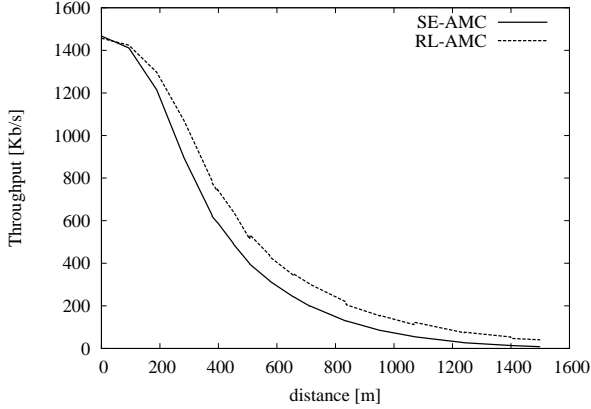In the following experiments we assume that each UE implements the SINR to CQI mapping described in [25]. First of all

Fig. 3. Average throughput as a function of the distance of the tagged user from the eNB in a pedestrian scenario.



Fig. 4. Average cell throughput as a function of the number of UEs in an urban vehicular scenario.

we consider the same network scenario as in Figure 2, i.e., ten static UEs randomly deployed and one tagged UE moving at pedestrian speed. Then, Figure 3 shows a comparison of the throughput achieved by the tagged user with both SE-AMC and RL-AMC schemes at different distances of the tagged UE from the eNB. We can observe that the MCS selection in SE-AMC is too conservative and this results in a throughput loss. On the contrary, RL-AMC method is able to discover the MCS configuration that can ensure a more efficient use of the available channel resources. This is more evident at intermediate distances from the eNB when short-term fading may lead to use more frequently low-rate MCSs. As shown in the figure, the throughput improvement varies between 20% and 55% in the range of distances between 200 meters and 800 meters.

In the second set of simulations we consider a more dynamic environment in which there is an increasing number of UEs in the cell, and all the UEs are moving according to RWM with speed 30 km/h and pause time equal to 5 seconds. Figure 4 shows a comparison of the aggregate cell throughput with both SE-AMC and RL-AMC schemes as a function of the network congestion (i.e., number of UEs). The results clearly indicate that the throughput improvement provided by RL-AMC is almost independent of the number of UEs and it is about 10%. We can also observe the the cell capacity initially increases when going from 10 to 20 UEs. This is due to two main reasons. First, RR is able to allocate RBs in a more efficient way when the number of UEs is higher. Second, the higher the number of UEs and the higher the probability that one of the UEs is close to the eNB and it can use high data-rate MCSs.

To investigate more in depth the behaviour of the considered AMC schemes, in Figure 5 we show the probability mass function of the number of retransmissions that are needed to successfully transmit a segment of data in a cell with 50 UEs moving as described above. We remind that the same MCS is used for both the first transmission attempt and the eventual subsequent retransmissions. We can observe that with
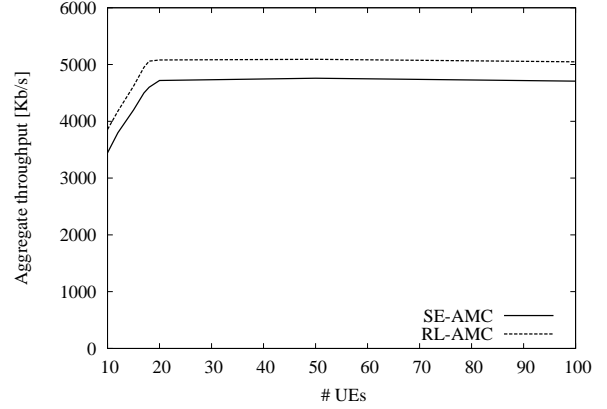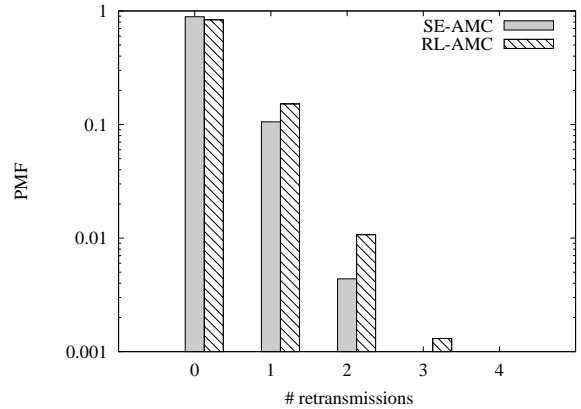


Fig. 5. Probability mass function of the number of retransmissions in an urban vehicular scenario with 50 UEs.

RL-AMC the probability to successfully transmit a segment of data at the first transmission attempt is slightly lower than with SE-AMC. However, the probability of successfully transmiting a segment of data after one or two retransmissions is higher with RL-AMC than with SE-AMC. This confirms our previous observation that the initial MCS selection of SE-MAC is more conservative. On the contrary, RL-AMC is able to also explore MCS with higher data rates when the channel conditions are more favourable and this is beneficial for the throughput performance. Note that this is achieved without violating the BLER requirements imposed by the LTE standard.

## VI. CONCLUSIONS

In this paper,we have presented a new AMC method for LTE networks that is based on reinforcement learning techniques, We have discussed how inaccurate feedbacks on channel qualities and the complexity of modelling link performance under realistic channel models may easily lead to inaccurate MCS selections. By exploiting reinforcement learning, we can significantly reduce the impact of channel prediction errors on the performance of link adaptation. As future work we plan to explore the use of SINR measurements for directly guiding

the MCS selection. In this case scale-spacing method have to be designed to reduce the state space. A critical extension of this work concerns the investigation of methods to reduce the (typically long) convergence delays of reinforcement learning. To this end recent advancements in RL theory, such as actor-critic methods, will be considered.

## REFERENCES

[1] A. Ghosh, R. Ratasuk, B. Mondal, N. Mangalvedhe, and T. Thomas, "LTE-advanced: next-generation wireless broadband technology [Invited Paper]," *IEEE Wireless Communications*, vol. 17, no. 3, pp. 10–22, June 2010.

[2] R. Fantacci, D. Marabissi, D. Tarchi, and I. Habib, "Adaptive modulation and coding techniques for OFDMA systems," *IEEE Transactions on Wireless Communications*, vol. 8, no. 9, pp. 4876–4883, Semptember 2009.

[3] 3GPP: Technical Specification Group Radio Access Network, "Evolved Universal Terrestrial Radio Access (E-UTRA); Physical layer procedures (Release 11)," 3GPP TS 36.213 V11.3.0, June 2013.

[4] J. Fan, Q. Yin, G. Li, B. Peng, and X. Zhu, "MCS Selection for Throughput Improvement in Downlink LTE Systems," in *PRoc. of IEEE ICCCN'11*, 2011, pp. 1–5.

[5] J. Francis and N. Mehta, "EESM-Based Link Adaptation in Point-to-Point and Multi-Cell OFDM Systems: Modeling and Analysis," *IEEE Transactions on Wireless Communications*, vol. 13, no. 1, pp. 407–417, January 2014.

[6] S. Tsai and A. Soong, "Effective-SNR mapping for modeling frame error rates in multiple-state channels," 3GPP, Tech. Rep. 3GPP2-C30-20030429-010, 2003.

[7] J. Olmos, S. Ruiz, M. García-Lozano, and D. Martín-Sacristán, "Link Abstraction Models Based on Mutual Information for LTE Downlink," COST 2100, Tech. Rep. 11052, June 2010.

[8] Y. Blankenship, P. Sartori, B. Classon, V. Desai, and K. Baum, "Link error prediction methods for multicarrier systems," in *Proc. of IEE VTC-Fall'04*, vol. 6, 2004, pp. 4175–4179.

[9] K. Brueninghaus, D. Astely, T. Salzer, S. Visuri, A. Alexiou, S. Karger, and G.-A. Seraji, "Link performance models for system level simulations of broadband radio access systems," in *Proc. of IEEE PIMRC'05*, vol. 4, 2005, pp. 2306–2311.

[10] M. Ni, X. Xu, and R. Mathar, "A channel feedback model with robust SINR prediction for LTE systems," in *Proc. of EuCAP'13*, 2013, pp. 1866–1870.

[11] J. Ikuno, S. Pendl, M. Simko, and M. Rupp, "Accurate SINR estimation model for system level simulation of LTE networks," in *Proc. of IEEE ICC'12*, 2012, pp. 1471–1475.

[12] A. Kuhne and A. Klein, "Throughput analysis of multi-user ofdma-systems using imperfect cqi feedback and diversity techniques," *IEEE Journal on Selected Areas in Communications*, vol. 26, no. 8, pp. 1440–1450, October 2008.

[13] R. Akl, S. Valentin, G. Wunder, and S. Stanczak, "Compensating for CQI Aging By Channel Prediction: The LTE Downlink," in *Proc. of IEEE GLOBECOM'12*, 2012, pp. 4821–4827.

[14] G. Xu and Y. Lu, "Channel and Modulation Selection Based on Support Vector Machines for Cognitive Radio," in *Proc. of WiCOM'06*, 2006, pp. 1–4.

[15] R. Daniels, C. Caramanis, and R. Heath, "Adaptation in Convolutionally Coded MIMO-OFDM Wireless Systems Through Supervised Learning and SNR Ordering," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 1, pp. 114–126, January 2010.

[16] R. Daniels and R. Heath, "Online adaptive modulation and coding with support vector machines," in *Proc. of EW'10*, 2010, pp. 718–724.

[17] J. Leite, P. H. De Carvalho, and R. Vieira, "A flexible framework based on reinforcement learning for adaptive modulation and coding in OFDM wireless systems," in *Proc. of IEEE WCNC'2012*, 2012, pp. 809–814.

[18] Z. He and F. Zhao, "Performance of HARQ with AMC Schemes in LTE Downlink," in *Proc. of IEEE CMC'10*, vol. 2, 2010, pp. 250–254.

[19] P. Tan, Y. Wu, and S. Sun, "Link adaptation based on adaptive modulation and coding for multiple-antenna ofdm system," *IEEE Journal on Selected Areas in Communications*, vol. 26, no. 8, pp. 1599–1606, October 2008.

[20] T. Jensen, S. Kant, J. Wehinger, and B. Fleury, "Fast link adaptation for mimo ofdm," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 8, pp. 3766–3778, October 2010.

[21] Donthi, S.N. and Mehta, N.B., "An Accurate Model for EESM and its Application to Analysis of CQI Feedback Schemes and Scheduling in LTE," *IEEE Transactions on Wireless Communications*, vol. 10, no. 10, pp. 3436–3448, October 2011.

[22] T. Tao and A. Czylwik, "Combined fast link adaptation algorithm in LTE systems," in *Proc. of ICST CHINACOM'11*, 2011, pp. 415–420.

[23] R. Sutton and A. Barto, *Reinforcement Learning: An Introduction*. The MIT Press, March 1998.

[24] C. Watkins and P. Dayan, "Q-Learning," *Machine Learning*, vol. 8, pp. 279–292, 1992.

[25] N. Baldo, M. Miozzo, M. Requena-Esteso, and J. Nin-Guerrero, "An Open Source Product-oriented LTE Network Simulator Based on Ns-3," in *Proc. of ACM MSWiM'11*, 2011, pp. 293–298.

[26] 3GPP: Technical Specification Group Radio Access Network, "Conveying MCS and TB size via PDCCH," TSG-RAN WG1 R1-081483, March 2008.

[27] COST Action 231, "Digital mobile radio future generation systems," Final Report - EUR 18957, 1999.

[28] W. Jakes, *Microwave Mobile Communications*. John Wiley & Sons Inc., 1975.

[29] 3GPP: Technical Specification Group Radio Access Network, "Evolved Universal Terrestrial Radio Access (E-UTRA); Base Station (BS) radio transmission and reception," 3GPP TS 36.104 V11.7.0, January 2014.

[30] D. Karamshuk, C. boldrini, M. Conti, and A. Passarella, "Human mobility models for opportunistic networks," *IEEE Communications Magazine*, vol. 49, no. 12, pp. 157–165, 2011.

# Offloading through Opportunistic Networks with Dynamic Content Requests

Raffaele Bruno, Antonino Masaracchia, and Andrea Passarella
IIT-CNR
Via G. Moruzzi 1, 56124 Pisa, Italy
{first.last}@iit.cnr.it

*Abstract*—Offloading is gaining momentum as a technique to overcome the cellular capacity crunch due to the surge of mobile data traffic demand. Multiple offloading techniques are currently under investigation, from modifications inside the cellular network architecture, to integration of multiple wireless broadband infrastructures, to exploiting direct communications between mobile devices. In this paper we focus on the latter type of offloading, and specifically on offloading through opportunistic networks. As opposed to most of the literature looking at this type of offloading, in this paper we consider the case where requests for content are *non-synchronised*, i.e. users request content at random points in time. We support this scenario through a very simple offloading scheme, whereby no epidemic dissemination occurs in the opportunistic network. Thus our scheme is minimally invasive for users' mobile devices, as it uses only minimally their resources. Then, we provide an analysis on the efficiency of our offloading mechanism (in terms of percentage of offloaded traffic) in representative vehicular settings, where content needs to be delivered to (subsets of the) users in specific geographical areas. Depending on various parameters, we show that a simple and resource-savvy offloading scheme can nevertheless offload a very large fraction of the traffic (up to more than 90%, and always more than 20%). We also highlight configurations where such a technique is less effective, and therefore a more aggressive use of mobile nodes resources would be needed.

## I. INTRODUCTION

In the last few years, we have observed a drastic surge of data traffic demand from mobile personal devices (smartphones and tablets) over cellular networks [1]. This has already generated famous collapses of 3G networks in the recent past, (e.g. [2]), showing that standard cellular technologies may not be enough to cope with this data demand. Even though significant improvement in cellular bandwidth provisioning are expected through LTE-Advanced systems, the overall situation is not expected to change significantly [3]. Besides personal mobile devices, the diffusion of M2M and IoT devices is expected to increase at an exponential pace (the share of M2M devices is predicted to increase 5x by 2018 [1]), which is likely to generate a corresponding increase in the demand for mobile traffic (11-fold increase by 2018 [1]).

*Offloading* part of the traffic from the cellular to another, complementary, network, is currently considered one of the most promising approaches to cope with this problem, with offloaded traffic being foreseen to account for at least 50% of the overall traffic in the coming years [1]. Among the various forms of offloading that are currently investigated (described in more detail in Section II), in this paper we consider offloading through opportunistic networks. Opportunistic networks [4] exploit physical proximity between mobile nodes to enable direct communication between them. They typically exploit ad hoc enabling technologies like WiFi-direct or Bluetooth, and support dissemination of messages through multi-hop space-time paths, i.e., multi-hop paths that develop both over space - as in conventional ad hoc multi-hop networks - and over time - by exploiting contact opportunities between nodes that become available over time due to their mobility. In offloading schemes based on opportunistic networks (e.g., [5], [6], [7]), content is initially seeded on a subset of the mobile nodes, and then it spreads through the opportunistic network to reach all interested users. Cellular bandwidth is thus used only to seed the network, possibly to add additional seeds over time (as in [6]), and to send content directly to interested users upon expiration of the deadline for its delivery. These schemes permit to significantly offload the cellular network, while at the same time guaranteeing bounded delays in content delivery.

Most of the literature on opportunistic-based offloading investigates the scenario where a specific piece of content is generated, and the set of users to whom it has to be delivered is known already at that time and does not change subsequently, i.e. requests are *synchronised*. While significant, this scenario only partially captures relevant use cases. In particular, it does not cover cases where content demand is *dynamic*, i.e. users' requests for the same piece of content can arrive at different time instants. In the latter scenario offloading can still be applied: upon a request, content can reach the requesting user either through the opportunistic network, exploiting an ongoing dissemination process, or through the cellular network, in case the opportunistic dissemination does not reach the user in time. Offloading may even be more needed in case of dynamic requests, as synchronised requests could in principle be served also through multicast transmissions (although [8] shows that offloading is beneficial also when multicast is applied).

In this paper we start investigating dynamic content requests, with a particular focus on vehicular scenarios. We deliberately use a very simple offloading scheme, described in Section III, whereby resources provided by mobile nodes are minimally used. Nodes interested in a content store it for a limited amount of time after receiving it. New requests from other users are satisfied either when the requesting user encounters another user storing a copy of the content, or through the cellular network upon expiration of the delivery deadline.

As opposed to most of the literature looking at offloading through opportunistic networks, in our scheme we do not use any epidemic dissemination mechanism. On the one hand, this allows us to test a minimally invasive offloading scheme from the mobile users' perspective. As additional resources spent by mobile devices are sometimes considered a possible roadblock for offloading, our results show the offloading efficiency when this additional burden is extremely low. On the other hand, this simple scheme allows us to stress the efficiency of offloading in a particularly unfavourable configuration, thus providing a worst-case analysis, all other conditions being equal.

We focus on two complementary scenarios. In the first one, users move in a given physical area, and *all* request a piece of content, though at different points in time. This scenario is representative of users moving inside a limited area, and accessing very popular content, though not particularly time critical (i.e., content that does not generate a surge of requests immediately when it is generated). In the second scenario, users enter and exit (after a short amount of time) a given geographical area, and request content after a random amount of time after they entered the area. This complementary scenario is thus representative of users traversing a geographical area, as opposed to roaming there. Finally, in this scenario we also consider the case where content is requested only with a certain probability, i.e., when content has different levels of popularity.

We analyse the offloading efficiency in these scenarios, defined as the fraction of nodes receiving content through the opportunistic network. We characterise efficiency as a function of key parameters such as the number of users, the deadline of content requests, the time after which users drop the content after having received it, the popularity of the content. As we show in Section IV, even with an unfavourable opportunistic dissemination scheme, we find that offloading can be very efficient, as it is possible to offload up to more than 90% of the traffic. In other configurations, we find that the considered offloading scheme is less efficient, resulting in an offloading of only about 20%. In such cases, however, there is ample room for improvement, by further leveraging opportunistic networking resources, e.g., through more aggressive content replication schemes.

## II. RELATED WORK

Offloading can take several forms. In some cases, traffic is offloaded by using modifications inside the cellular architecture (e.g. LIPA/SIPTO [9] or small cells [10]), or other wireless access infrastructures, primarily WiFi [11], [12].

In this paper we consider offloading that exploits direct communications between mobile devices. Also in this case there are several approaches. In the 3GPP area, the device-to-device (D2D) [13] architectural modification to LTE has been defined, that devotes part of the cellular resources to direct communication between devices under strict control of a common eNB. Instead, we focus on using opportunistic networks together with cellular networks, as previously proposed, e.g. in [5], [6], [7], [8]. In this case, offloading

exploits technologies (such as WiFi direct or Bluetooth) that do not interfere with cellular transmissions, and therefore no coordination is required with the eNB. In addition, mobile devices run self-organising networking algorithm to disseminate offloaded content without strict control of the eNBs or any other central controller.

The most common scenario where opportunistic offloading is used is content dissemination to a set of interested users. In most cases, it is assumed that the set of users interested in receiving a piece of content is known when the content is generated (or, alternatively, the content is implicitly requested by all interested users immediately when it is generated) and do not change over time. In addition, content is "seeded" through the cellular network on a subset of interested users, and then a dissemination process starts in the opportunistic network in order to reach the rest of the users [5]. Typically, epidemic dissemination is assumed [14]. In addition, some other papers (e.g. [6], [7]) consider that content must be delivered to users within a given deadline. To meet this deadline, content can be sent through the cellular network to additional seeds during the dissemination process, and is finally sent to users that are still missing it when the deadline is about to expire ("panic zone"). To know which users have received the content, a lightweight control channel is implemented through the cellular network, whereby users send an ACK to a central controller that tracks the status of the dissemination process, and determines when to seed additional copies of the content, and when to directly deliver content to the users in the panic zone.

With respect to this body of work, this paper differs in two main aspects. On the one hand, we release the assumption that users interested in a content request it simultaneously. In our scenarios content requests occur over time dynamically. On the other hand, we do not assume epidemic dissemination of content, but consider that content is exchanged in the opportunistic network only between users that have requested it, when they encounter directly. Therefore, our scenario covers more general cases with respect to strictly synchronised requests, and, in addition, provides a worst-case analysis of the potential of offloading, as we use the least possible aggressive form of dissemination in the opportunistic network.

To the best of our knowledge, the only other paper where content requests are not synchronised is [15]. That paper assumes that users become interested in the content after a random amount of time after its generation, and the goal of the proposed system is to maximise the probability that the user have already the content by then. This is very different from our scheme, which works reactively, *after* users generate requests.

Finally, offloading has been also proposed specifically in vehicular environments. In this case offloading schemes often assume the presence of RoadSide Units (RSU) [16] to support the dissemination process (e.g., by pre-fetching popular contents), which we do not assume here, to obtain a solution requiring no additional infrastructure development. Last but not least, offloading is proposed also for aggregating and uploading traffic generated by cars, e.g., in the context of Floating

Car Data (FCD) [17]. This is clearly a different application and offloading scenario with respect to that considered in this paper.

## III. Offloading mechanisms

As anticipated in Section I we deliberately consider a simple scheme that uses very little resources of mobile nodes to support the offloading process. In general, we support scenarios where content is requested by users at random points in time. Similarly to [6], we assume the existence of a Central Dissemination Manager (CDM), that can communicate with all nodes through the cellular network and keeps track of the dissemination process. Without loss of generality[1], in the following we focus on the dissemination of a single piece of content to the set of interested users. The offloading mechanism is defined by the actions taken by requesting nodes and by the CDM, as described by Algorithms 1 and 2, respectively.

Let us focus first on the actions taken by requesting nodes (Algorithm 1). When a request is generated at a node, the node sends it to the CDM via the cellular network (line 3). The node is guaranteed to receive the content within a given *content timeout*. During the timeout, the node tries to get the content from encountered nodes (lines 5-12). If the timeout expires, it receives it directly from the CDM (lines 13-16). Upon receiving the content, the node sends an ACK to the CDM (line 9 and, implicitly, line 14). In addition, it keeps the content for a *sharing timeout*, during which it can share the content with other encountered nodes (lines 18-20). After the expiration of the *sharing timeout* the content is deleted from the local cache. Note that requests and ACKs are supposed to be much shorter than the content size, and thus do not significantly load the cellular network.

Let us now focus on the actions taken by the CDM (Algorithm 2). Thanks to requests and ACKs, the CDM is always aware of the status of content availability in the network. Upon receiving a request, it checks whether some other node is already storing a copy of the content or not. In the latter case (lines 4-6) there is no chance that the user can get the content opportunistically through another node, and the CDM sends the content directly through the cellular network. In the former case (lines 7-21), it waits to receive an ACK during the *content timeout* (lines 8-15), indicating that the node has received the content. If this does not happen, it sends the content directly to the node (lines 16-20). Finally, upon expiration of the *sharing timeout* for a given node the CDM updates the view on the number of nodes with the content (lines 22-23)[2].

---

[1]Strictly, this is the case when congestion on the opportunistic network is low, and therefore the effect of multiple contents offloaded at the same time can be neglected. This is typically assumed in the literature on offloading through opportunistic networks.

[2]Note that the CDM implementation could be further simplified by allowing the nodes that select a content to send a message over the cellular network to inform the CDM. In this way, the CDM does not need to maintain separate timers for each of the nodes that have received the content. It is also reasonable to assume that such confirmation message would be a negligile overhead for the cellular network.

---

**Algorithm 1** Actions taken by requesting nodes

$\triangleright$ Run by a tagged node $k$
1: **Upon** request for content $C$
2: content_received = **false**
3: **Send** content_request to CDM
4: **if** $C$ not received immediately from CDM **then**
$\triangleright$ try with opportunistic contacts
5:     **while** content_timeout is not over **do**
6:         request $C$ to encountered nodes
7:         **if** content received **then**
8:             content_received = **true**
9:             **Send** ACK to CDM
10:             **break**
11:         **end if**
12:     **end while**
13:     **if** content_received == **false then**
14:         **Receive** $C$ from CDM
15:         content_received = **true**
16:     **end if**
17: **end if**
18: **while** sharing_timeout is not over **do**
$\triangleright$ available for opportunistic sharing
19:     **Send** $C$ to encountered nodes upon request
20: **end while**
21: **Cancel** content $C$

---

**Algorithm 2** Actions taken by CDM

$\triangleright$ Run by the CDM for content $C$
**Init** #nodes_with_$C$ = 0
1: **Upon** request from node $k$
2: $k$_served = **false**
3: **if** #nodes_with_$C$ == 0 **then**
4:     **Send** $C$ to $k$
5:     #nodes_with_$C$++
6:     **Set** sharing_timeout for node $k$
7: **else**
8:     **while** content_timeout is not over **do**
9:         **if** ACK received by $k$ **then**
10:             #nodes_with_$C$++
11:             $k$_served = **true**
12:             **Set** sharing_timeout for node $k$
13:             **break**
14:         **end if**
15:     **end while**
16:     **if** $k$_served = **false then**
17:         Send $C$ to $k$
18:         #nodes_with_$C$++
19:         **Set** sharing_timeout for node $k$
20:     **end if**
21: **end if**

22: **Upon** sharing_timeout for node $k$ over
23: #nodes_with_$C$ = #nodes_with_$C$-1

With respect to offloading mechanisms proposed for opportunistic networks (e.g., [5], [6]) our algorithms present several differences. First, there is no proactive seeding of the network. This is because requests arrive at the CDM dynamically, and there is no knowledge of which nodes will generate a request, and when. Therefore, we adopted a reactive policy, i.e. we wait for requests without doing any proactive seeding. Second, we want to use minimally mobile node resources in the opportunistic network. This is to make the offloading mechanism less intrusive as possible, as the additional mobile devices' resource usage brought about by offloading is often considered a possible severe drawback. Therefore, we do not use epidemic dissemination in the opportunistic network. For the same reasons, we assume that users drop content some time after receiving it. Still, our algorithms guarantee bounded delay, and impose similar overhead on the CDM as in previous proposals [6]. Clearly, Algorithms 1 and 2 can be easily modified to exploit additional resources of mobile devices (e.g., using more aggressive forms of dissemination or doing initial proactive seeding), if needed.

## IV. PERFORMANCE EVALUATION

### A. Scenarios and performance indices

We test the performance of the proposed offloading schemes in two different vehicular scenarios, hereafter denoted as Scenario A and B.

In Scenario A we capture cases where a group of vehicles move inside a geographical area covered by a cell, and roam always inside that cell. Vehicles move on a stretch of road crossing the cell, and come back when arriving at the boundary. The resulting traffic is therefore bidirectional. Nodes move with a speed randomly selected (with uniform distribution) in an interval $[v_{min}, v_{max}]$, and can exchange content directly while being within a maximum transmission range $T_{RX}$ from each other. We consider $N$ nodes in the simulations, which all request the content. Requests are generated from the beginning of the simulation sequentially, according to a Poisson process with rate $\lambda$ (i.e. two requests are spaced by an exponentially distributed time interval). Simulations lasts until all nodes have requested the content, and their *sharing timeouts* are all expired. In other words, we start from a condition where no nodes have any copy of the content, and we analyse the behaviour of the system until no copy of the content is available after all nodes have received it. While assuming vehicles go back and forth on a given road segment is a simplification, the scenario is still representative of movement patterns confined in a geographical area served by a cellular network, where a given content is very popular and thus requested by all users (though at different points in time). More in general, the scenario is representative of movement patterns whereby vehicles roam in such a geographical area, can move in opposite directions and can communicate with each other when being close enough, irrespective whether such movements occur on the same street or on different, nearby streets.

In Scenario B we capture cases where nodes are not necessarily staying in the same area, but there is a constant flux of vehicles entering and exiting the area. Again, we assume that vehicles move on a road and we focus on a road segment covered by a cell (we select speeds as in Scenario A). Traffic is again bidirectional, and we keep the number of nodes constant, and assume that a new vehicle enters the area when another one has left. When entering the area, vehicles become interested in the content with a given probability $p$. If they are interested, they generate a request after a time interval uniformly distributed between the time when they enter and the time when they reach the centre of the cell. Taking the same terminology of [6], we define a panic zone as the area of the cell $\Delta$ meters before the boundary. The *content timeout* is set so that the CDM sends the content directly when vehicles enter the panic zone. Finally, vehicles keep the content while being inside the cell. At the beginning of simulations, nodes are distributed randomly (with a uniform distribution) in the cell, are interested in content with probability $p$, and generate a request at a point in time uniformly distributed between the simulation start time and when they are midway towards the border of the cell. Simulations stop after 100 requests have been generated (50 in the case of low popularity content, without noticeable loss of statistical significance of the results), and the corresponding users have been all served. With this scenario we explore different cases with respect to Scenario A. After an initial transient phase, we are able to show a steady-state behaviour of offloading, in cases where vehicles enter and exit an area with a given flux and density. In other words, we can show how much offloading is efficient in making a given content "survive" in a geographical area, by only exploiting replicas available on vehicles of interested users passing through that area. This is an application of the basic floating content idea [18] to the case of vehicular networking environment in presence of offloading. In addition, only a fraction of the nodes can be interested into the content, i.e. the content can have different levels of popularity.

We ran simulations, using the NS3 with the LENA module for LTE[3], for various sets of parameters, as indicated in Table I. Specifically, we varied the number of nodes in both Scenarios, the request rate, the *content timeout* and the *sharing timeout* in Scenario A, and the content popularity in Scenario B. Request rate and popularity obtain similar effects in the two scenarios, as they modify the average number of nodes interested (and receiving) content at any point in time, and thus the density of nodes with a content replica. We performed at least 5 simulation runs for each set of parameters, using the independent replication method [19]. The main performance figure we consider is the offloading efficiency, defined as the fraction of content messages that reach the users through opportunisitc communications. For this index we computed the confidence intervals (with 95% confidence level) over the replications. To get a more precise idea on the dynamics of the offloading process over time, we also computed, on each 5s

---

[3]http://networks.cttc.es/mobile-networks/software-tools/lena/

time window, the average (across simulation replicas) number of copies of content stored on mobile nodes, and the average number of new content deliveries through the cellular and the opportunistic network, respectively.

| | Scenario A | Scenario B |
|---|---|---|
| speed (Km/h) | [80,120] | [90,110] |
| cell diameter (Km) | 4 | 1 |
| $N$ (nodes) | 20, 40 | 20, 40 |
| $T_{RX}$ (m) | 200 | 50 |
| $p$ | 1 | 0.5, 0.75, 1 |
| $\lambda$ (req/s) | 1, 0.5, 0.2 | – |
| *content timeout* (s) | 60, 90, 120 | – |
| *sharing timeout* (s) | 5,10,20,30,60,120 | – |
| $\Delta$ (m) | – | 50 |

*B. Analysis of scenario A*

We start by analysing the system performance in Scenario A. To this end, Figure 1 shows the offloading efficiency obtained in a wide set of different network configurations, in which we vary the node density, the content request rate, as well as the *content timeouts* and the *sharing timeouts*. Several general observations can be drawn from the shown results. First, the offloading efficiency increases with the node density. The main reason is that the higher the node density, the higher the contact rate between the mobile devices. Thus, there are more opportunities for opportunistic dissemination between interested users. As far as the impact of the request rate ($\lambda$) we observe two regimes. When the *sharing timeout* is low, higher request rates result in higher offloading. This is intuitive, because higher request rates results in requests being more concentrated in time. When nodes share the content only for very short amounts of time (see for example the case of 5s), concentrating the requests in time increases the probability of encountering other nodes sharing the content. Less intuitive is the behaviour for large sharing timeouts, where higher request rates results in *lower* offloading efficiency. The reason of this will be more clear when analysing the evolution of dissemination over time (Figure 2). Intuitively, when requests are more concentrated in time, *content timeouts* for nodes that do not get the content via the opportunistic network are also more concentrated. As we will discuss later, when a timeout expires and content is delivered via the cellular network, this kicks off a fast increase in the dissemination of content via the opportunistic network in the region of the node whose *content timeout* has expired. When expirations are less concentrated in time (i.e., when request rates are lower), the opportunistic diffusion process has more time to spread content, and therefore the offloading efficiency increases.

A second interesting observation is related to the impact of the *sharing timeout* on the offloading efficiency. Our results indicate that if the content is sufficiently persistent in the network (e.g., *sharing timeout* $\geq$ *content timeout*) then the impact of the *sharing timeout* on the offloading efficiency is negligible. On the other hand, if the content is

volatile, i.e., it is cached in the local memory of interested users only for few seconds, then the number of copies of that content in the environment may be too small to allow an efficient opportunistic dissemination. For instance, with 20 mobile devices and a content request rate of 0.2 req/s the offloading efficiency can be as low as 20% (this degradation of the offloading efficiency is less remarkable in denser networks). Interestingly, if the content request rate is high (i.e., $\lambda = 1$ req/s) then even a *sharing timeout* = 5s can still provide an offloading efficiency up to 60% in a cell with 20 mobile devices. A last observation is related to the effect of the *content timeout*. As shown in Figure 1c and Figure 1d an increase in the *content timeout* results into an increase of the offloading efficiency. This is more noticeable for large *sharing timeouts*, i.e. when content stays available on nodes for opportunistic dissemination longer. This is basically a joint effect of the fact that (i) content is available longer in the opportunistic networks (longer *sharing timeouts*) and (ii) interested nodes wait longer for requesting it via the cellular network (longer *content timeouts*).

To get a deeper understanding of the offloading dynamics, plots in Figure 2 show the temporal evolution of (i) the total number of mobile devices that have received the content via the cellular and the opportunistic network, respectively, and (ii) the number of copies of the content available in the network (i.e., the number of nodes that are storing and sharing a copy of the content at that time). Note that plots are typically shown until just after the time when the last requesting node has received the content. The system evolution after that time is not particularly interesting: nodes progressively drop the content when their *sharing timeouts* expire. We show plots for extreme values of the considered parameters. Specifically, in Figures 2a and 2b we focus on two extreme values of the *sharing timeout* parameter, for the case of 40 nodes and 1 req/s (*content timeout* is always 60s). As expected, the main difference is the number of copies of the content available in the network, which is much higher in Figure 2b, resulting in a higher offloading efficiency. It is very interesting to observe the behaviour of the system after 60s, i.e. when the *content timeouts* for the first nodes generating requests expire. For larger *sharing timeouts* (larger than the *content timeout*), before that time, only 1 node (the first one requesting the content) can receive the content via the cellular network, as it is clear from Figure 2b. This is due to the behaviour of the CDM explained in Section III, that sends immediately a content to the first requesting node (as no other node stores the content yet), and then waits the *content timeouts* (i.e. 60s) for the next requests before taking any action. In other words, it is impossible to have more than one delivery via the cellular network in the first 60s, due to the CDM algorithm. When the *sharing timeout* is short, more copies of the content can be sent via the cellular network also before the first *content timeout* expires. This happens whenever a new request is generated and all nodes that have previously received the content have already dropped it (due to expiration of the *sharing timeout*). In both cases, after 60s from the start of the simulation, *content timeouts* start expiring, and
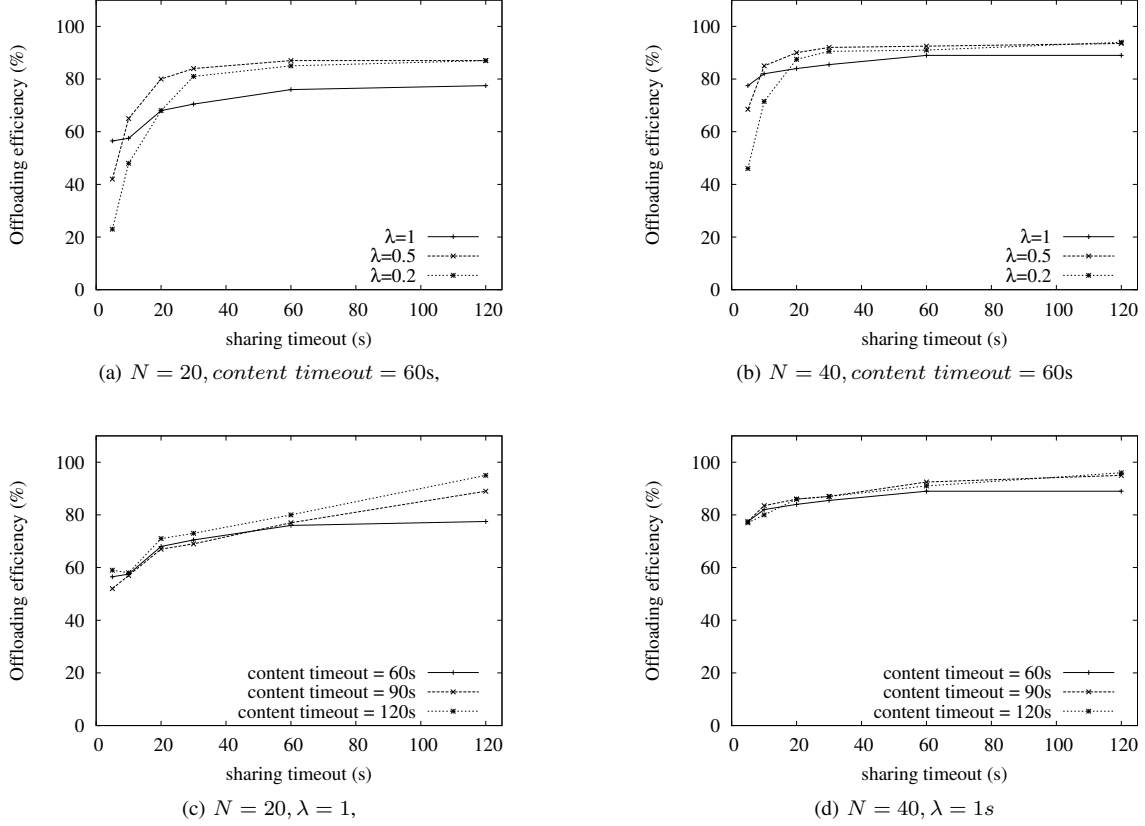
Fig. 1.   Scenario A: offloading efficiency for varying request rates, *content timeouts* and the *sharing timeouts*.

new copies of the content are sent via the cellular network. This generates a burst of dissemination in the opportunistic network, that is noticed by the steep increase of the curve related to opportunistic deliveries around that time. Note, in particular, that after 60s in both cases the rate of increase of the delivery via the opportunistic network is higher than the rate of increase of cellular deliveries. This means that each delivery via the cellular network is significantly amplified by deliveries in the opportunistic network.

Figures 2c and 2d show the evolution over time in the least favourable conditions for offloading, i.e. for low request rates ($\lambda = 0.2$) and very short *sharing timeout* (5s). Curves confirm the behaviour described before. In particular at low densities ($N = 20$) the *sharing timeout* is not long enough to sustain significant dissemination over the opportunistic network. The situation improves for denser networks ($N = 40$), but still the *sharing timeout* makes nodes drop content too fast with respect to the rate of arriving requests (anyway, the offloading efficiency is still between 20% and 40% even in these cases).

### C. Analysis of scenario B

Figure 3 shows the offloading efficiency in Scenario B for the two considered densities of nodes and the different content popularities ($p$). Results basically confirm previous observations. This is nevertheless important, as Scenario B is more representative of a "steady state" behaviour of the offloading

system, as nodes constantly enter and exit the cell at a given rate, and continuously generate requests (with a given probability). Again, denser networks ($N = 40$) achieve higher offloading efficiency. The effect of the popularity parameter is similar to that of the request rate in Scenario A: the higher the popularity, the higher the number of nodes sharing content, the higher the offloading efficiency. It is interesting to note, however, that, due to the mobility of the nodes, they stay within the cell only for about 30s in total, and, on average, stay in the cell for about 22s after having generated a request. This is the "useful time window" during which they can receive content via opportunistic dissemination. Even though this time window is rather short, offloading is very efficient, even at quite low popularities ($p = 0.2$).
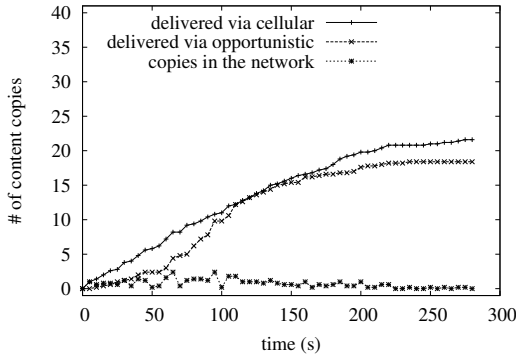
Finally, Figures 4a and 4b show the evolution over time for $N = 40$ nodes at the extreme popularity values. Besides confirming the general behaviour observed also in Scenario A, it is interesting to note that at high popularity the opportunistic dissemination alone is sufficient to keep enough copies of the content in the cell so that requesting nodes can find at least one before exiting. This is shown by the fact that the curve of delivery via the cellular network flattens out after an initial "seeding" interval. Instead, in case of less popular contents, there are cases where nodes do not encounter other nodes sharing a copy of the content before getting out of the
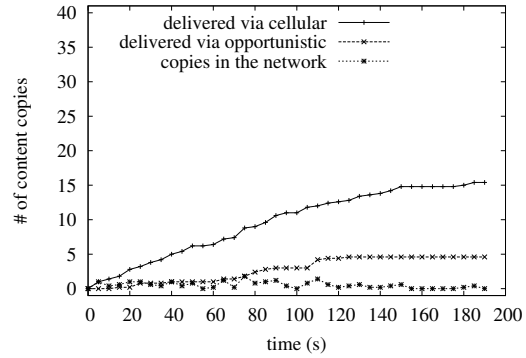
(a) $N = 40, \lambda = 1$ req/s, $sharing\ timeout = 5$s,

(b) $N = 40, \lambda = 1$ req/s, $sharing\ timeout = 120$s

(c) $N = 40, \lambda = 0.2$ req/s, $sharing\ timeout = 5$s,

(d) $N = 20, \lambda = 0.2$ req/s, $sharing\ timeout = 5$s

Fig. 2. Scenario A: temporal evolution of the number of content copies and served content requests in different network scenarios.
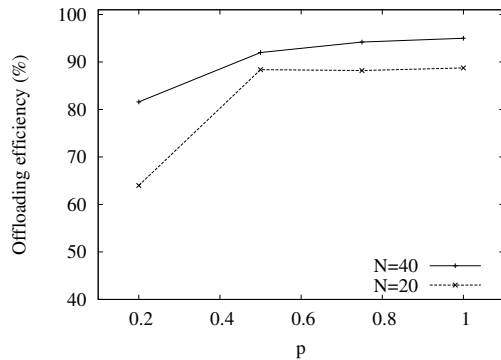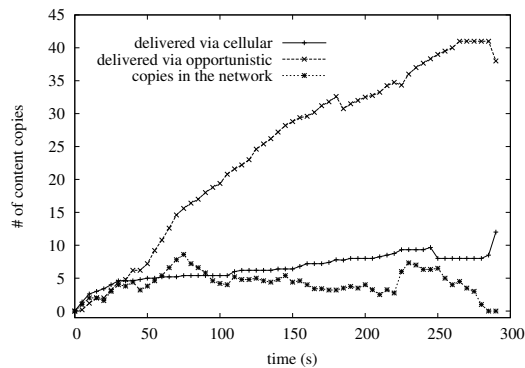


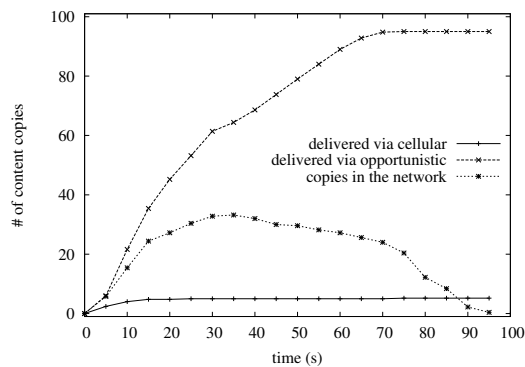Fig. 3. Scenario B: offloading efficiency for different content popularities.

cell, and therefore the CDM needs to serve them through the cellular network. Fluctuations in the number of copies stored in the network are mainly due to statistical fluctuations in the contacts and requests events. In addition, the curves drop towards the end of the simulation when only few requests need to be satisfied and no new requests are generated (remember that simulations stop when a maximum number of requests is reached).

## V. CONCLUSIONS

In this paper we have started to study the performance of offloading through opportunistic networks, in cases where content request are generated dynamically, and are not all synchronised at the moment when content becomes available. This general scenario is still to be satisfactorily addressed in the literature, and represents a large number of more specific scenarios. Interestingly, in such cases no support from cellular multicast mechanisms can be used, therefore offloading is even more critical. We have defined offloading mechanisms that guarantee bounded delays in content delivery, but, differently from existing literature, use as little as possible resources of mobile users' devices. This is also a critical point, as additional consumption of mobile devices' resources (storage, battery, etc.) is a drawback of offloading with opportunistic networks, that could limit its practical applicability. By considering minimal use of mobile devices' resources, we show that offloading can still be able to drastically reduce the traffic over the cellular network, also in a configuration that is unfavourable for its efficiency. Specifically, we tested the performance of our offloading schemes in vehicular environments, considering different densities of nodes, different popularity of content, and different parameters of the offloading protocols. Our results show that offloading can be very efficient also when using very limited resources of mobile devices, achieving

(a) low popularity ($p = 0.2$)



(b) high popularity ($p = 1$)

Fig. 4. Scenario B: temporal evolution of the number of content copies and served content requests in a network with $N = 40$ users.

offloading ratios up to more than 90%. Moreover, our results also highlight configurations of the protocols and parameters of the investigated scenarios where offloading is less efficient, and therefore would benefit from more aggressive policies, using additional resources of mobile devices. It is worth noticing, however, that offloading efficiency never drops below 20% in the considered cases, even though they may be quite challenging for the considered offloading mechanisms.

## REFERENCES

[1] Cisco, "Cisco visual networking index: Global mobile data traffic forecast, update, 2013–2018," Cisco, Tech. Rep., http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/white_paper_c11-520862.html, 2013.

[2] J. Wortham, "Customers Angered as iPhones Overload AT&T," *NY Times*, 2009, http://www.nytimes.com/2009/09/03/technology/companies/03att.html?_r=0.

[3] "Growing data demands are trouble for verizon, lte capacity nearing limits." [Online]. Available: http://www.talkandroid.com/97125-growing-data-demands-are-trouble-for-verizon-lte-capacity-nearing-limits/

[4] L. Pelusi, A. Passarella, and M. Conti, "Opportunistic networking: data forwarding in disconnected mobile ad hoc networks," *Communications Magazine, IEEE*, vol. 44, no. 11, pp. 134–141, November 2006.

[5] M. V. Barbera, A. C. Viana, M. D. de Amorim, and J. Stefa, "Data offloading in social mobile networks through {VIP} delegation," *Ad Hoc Net.*, vol. 19, no. 0, pp. 92 – 110, 2014.

[6] J. Whitbeck, Y. Lopez, J. Leguay, V. Conan, and M. D. de Amorim, "Push-and-track: Saving infrastructure bandwidth through opportunistic forwarding," *Pervasive and Mob. Comp.*, vol. 8, no. 5, pp. 682 – 697, 2012.

[7] L. Valerio, R. Bruno, and A. Passarella, "Adaptive data offloading in opportunistic networks through an actor-critic learning method," in *ACM CHANTS*, 2014.

[8] F. Rebecchi, M. D. de Amorim, and V. Conan, "Flooding data in a cell: Is cellular multicast better than device-to-device communications?" in *ACM CHANTS*, 2014.

[9] K. Samdanis, T. Taleb, and S. Schmid, "Traffic offload enhancements for eutran," *Comm. Surveys Tutorials, IEEE*, vol. 14, no. 3, pp. 884–896, Third 2012.

[10] S. Singh and J. G. Andrews, "Joint resource partitioning and offloading in heterogeneous cellular networks," *IEEE TWC*, vol. 13, no. 2, pp. 888–901, February 2014.

[11] F. Mehmeti and T. Spyropoulos, "Performance analysis of on-the-spot mobile data offloading," in *IEEE GLOBECOM*, Dec 2013.

[12] K. Lee, J. Lee, Y. Yi, I. Rhee, and S. Chong, "Mobile data offloading: How much can wifi deliver?" *IEEE/ACM Trans. Netw.*, vol. 21, no. 2, pp. 536–550, Apr. 2013.

[13] F. Malandrino, C. E. Casetti, C. F. Chiasserini, and Z. Limani, "Fast resource scheduling in hetnets with d2d support," in *IEEE INFOCOM 2014*.

[14] A. Vahdat and D. Becker, "Epidemic routing for partially connected ad hoc networks," Duke University, Tech. Rep. CS-200006, 2006.

[15] X. Wang, M. Chen, Z. Han, D. Wu, and T. Kwon, "Toss: Traffic offloading by social network service-based opportunistic sharing in mobile social networks," in *IEEE INFOCOM*, 2014.

[16] F. Malandrino, C. Casetti, C. Chiasserini, and M. Fiore, "Content download in vehicular networks in presence of noisy mobility prediction," *IEEE TMC*, vol. 13, no. 5, pp. 1007–1021, May 2014.

[17] R. Stanica, M. Fiore, and F. Malandrino, "Offloading floating car data," in *IEEE WoWMoM*, June 2013, pp. 1–9.

[18] J. Ott, E. Hyytiä, P. Lassila, J. Kangasharju, and S. Santra, "Floating content for probabilistic information sharing," *Pervasive and Mobile Computing*, vol. 7, no. 6, pp. 671 – 689, 2011.

[19] A. M. Law and W. D. Kelton, *Simulation Modeling and Analysis*. McGraw-Hill, 2000.

# Flooding Data in a Cell: Is Cellular Multicast Better than Device-to-Device Communications?

Filippo Rebecchi
UPMC Sorbonne Universités
Thales Communications &
Security
filippo.rebecchi@lip6.fr

Marcelo Dias de Amorim
LIP6/CNRS – UPMC
Sorbonne Universités
marcelo.amorim@lip6.fr

Vania Conan
Thales Communications &
Security
vania.conan@thalesgroup.com

## ABSTRACT

A natural method to disseminate popular data on cellular networks is to use multicast. Despite having clear advantages over unicast, multicast does not offer any kind of reliability and could result costly in terms of cellular resources in the case at least one of the destinations is at the edge of the cell (i.e., with poor radio conditions). In this paper, we show that, when content dissemination tolerates some delay, providing device-to-device communications over an orthogonal channel increases the efficiency of multicast, concurring also to offload part of the traffic from the infrastructure. Our evaluation simulates an LTE macro-cell with mobile receivers and reveals that the joint utilization of device-to-device communications and multicasting brings significant resource savings while increasing the cellular throughput.

## Categories and Subject Descriptors

C.2.1 [**Network Architecture and Design**]: Store and forward networks; Wireless communication;

## Keywords

Cellular multicast; mobile data offloading; hybrid networks; delay-tolerant networks.

## 1. INTRODUCTION

With the deployment of increasingly performing cellular technologies, such as the 3GPP Long Term Evolution (LTE) and LTE-A, mobile networks will provide ever higher data rates (up to 100 Mbps for LTE, and 500 Mbps for LTE-A) [6]. Operators will exploit these opportunities to offer ubiquitous access to next generation services to their customers, such as multimedia applications, leading users to consume content anywhere, anytime. As a result, the expected mobile traffic will be very problematic to handle during peak times [5]. Among these multimedia services, some involve delivering the same piece of data to a community of interested users. Examples that fit this use case are software updates, on-demand

videos, and road traffic information. When a multitude of co-located users are interested in the same content, two possible approaches could help operators to relieve their cellular infrastructures: *multicast* and *mobile data offloading.*

Multicast makes use of a single unidirectional link, shared among several users inside the radio cell, allowing, in principle, a more efficient use of network resources with respect to the case where each user is reached through dedicated bearers. Note that a more precise terminology would be "multicast/broadcast", because only a subset of nodes is concerned by the content (multicast), and the shared nature of the wireless medium (broadcast) is exploited to transmit data. For the sake of readability, in the following we will only employ the term "multicast". To ensure coexistence between multicast and unicast services, operators must reserve a fixed amount of resources for multicast transmissions. Lately, field trials for video service during crowded sport events like the superbowl have tested the effectiveness of multicast [8]. Despite its attractive features, multicast presents intrinsic and still unresolved issues that limit its exploitation due to the difficult adaptation to radio channel conditions. Section 2 will provide an example of these inefficiencies.

Mobile data offloading is an alternative low cost solution to reduce the burden on the infrastructure network [7, 3, 13]. Direct device-to-device (D2D) communications may be employed to lower the load on the infrastructure. The increase in the density of mobile users gives rise to an abundance of contact opportunities and represents a strong argument to support opportunistic offloading strategies. Not surprisingly, this has been identified as one of the key enabling technologies for future cellular network architecture [1]. In order to encourage subscribers to offer their battery and storage resources to this end, mobile providers may offer monetary incentives and pricing discounts. As a counterpart, users should accept a delayed content reception.

In this paper, *we explore the combination of opportunistic traffic offloading with multicasting.* As we will see later, this strategy allows significant reduction in the load on the access part of the cellular network. As standard multicast is not intended for retransmissions, performance suffers and resources are wasted in the case of a single bad channel user inside the cell, due to trade-offs in coverage and efficiency. By including D2D communications into the picture, we obtain additional performance gains in terms of radio resources. Well-positioned users participate in mitigating the inefficiencies of multicast, by sharing their short-range resources to hand over content to users in bad cellular channel conditions. Depending on the number of participants requesting data,

we find a break-even point that achieves a good trade-off in terms of covered users and reception delay.

To assess the performance of this joint multicast/D2D approach it is necessary to evaluate the amount of radio resources consumed at the base station. This leads us to introduce a finer model of radio resource consumption than previous works in the offloading literature. Existing proposals do not consider heterogeneous channel conditions and assume that delivering a given amount of data to different users has always the same cost. Such an assumption does not hold in reality, as radio resources vary according to the channel condition experienced by each user. In other words, transmitting the same piece of content to users with different channel conditions do lead to uneven costs at the base station. To the best of our knowledge, we are the first to evaluate this aspect in the context of data offloading.

As a summary, the main contributions of this paper are:

- **Joint offloading strategy.** Our strategy employs direct D2D transmissions to assist the cellular distribution via multicast, permitting to consistently save resources at the cellular base stations.

- **Fine-grained resource consumption analysis.** We evaluate resource consumption employing the smallest radio resource unit that can be assigned to users for data transmission. This analysis shows that existing macroscopic techniques fail to capture actual system behaviors.

The remainder of the paper is organized as follows. We first present the motivation of our work in Section 2. The proposed joint offloading architecture and operation is described in Section 3. We evaluate the proposed system using realistic mobility traces in Section 4. We push the related work to Section 5 so that the reader has enough material to capture our original contribution. We finally conclude the paper and identify topics for future research in Section 6.

## 2. MOTIVATIONAL EXAMPLE

LTE proposes an optimized broadcast/multicast service through eMBMS (*enhanced Multimedia Broadcast Multimedia Service*), a point-to-multipoint specification to transmit control/data information from the cellular base station (eNB) to a group of user entities (UEs) [10].

Cellular UEs can use different modulation and coding schemes (MCS) to deal with variable channel characteristics. Each UE experiences different radio conditions, depending on path loss, interference from other cells, and wireless fading. UEs that are closer to the base station are able to decode data at a higher rate, while others located near the edge of the cell have to reduce their data rate and use a degraded MCS. This heterogeneity (time-varying and user-dependent) reduces the effectiveness of multicast because the eNB uses a single MCS to multicast downlink data. Usually, the selected MCS should be robust enough to ensure the successful reception and decoding of the data-frame for each recipient inside the cell. Thus, the worst channel among all the receivers dictates performance. An increase in the number of UEs boosts the probability that at least one UE experiences bad channel conditions, degrading the overall throughput [4].

To quantify this effect, we simulate a $500 \times 500$ m$^2$ single LTE cell with an increasing number of randomly located receivers using the ns-3 simulator [12]. Fig. 1 presents the
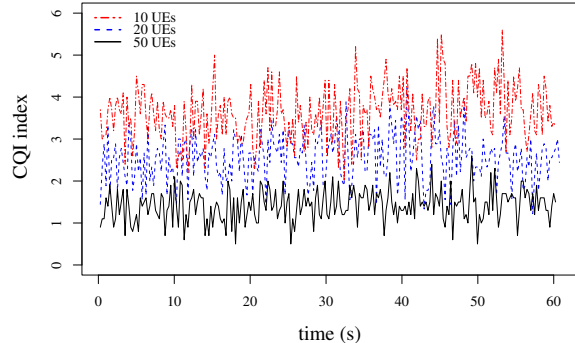


**Figure 1: Minimum CQI for different multicast group sizes. 100 runs, confidence intervals are tight and not shown in figure.**

**Table 1: CQI / MCS Table for LTE [2].**

| CQI index | Modulation schema | code rate x 1024 | Spectral Efficiency [bit/s/Hz] |
|---|---|---|---|
| 0 | | out of range | |
| 1 | QPSK | 78 | 0.1523 |
| 2 | QPSK | 120 | 0.2344 |
| 3 | QPSK | 193 | 0.3770 |
| 4 | QPSK | 308 | 0.6016 |
| 5 | QPSK | 449 | 0.8770 |
| 6 | QPSK | 602 | 1.1758 |
| 7 | 16-QAM | 378 | 1.4766 |
| 8 | 16-QAM | 490 | 1.9141 |
| 9 | 16-QAM | 616 | 2.4063 |
| 10 | 64-QAM | 466 | 2.7305 |
| 11 | 64-QAM | 567 | 3.3223 |
| 12 | 64-QAM | 666 | 3.9023 |
| 13 | 64-QAM | 772 | 4.5234 |
| 14 | 64-QAM | 873 | 5.1152 |
| 15 | 64-QAM | 948 | 5.5547 |

average minimum channel quality, in terms of CQI (Channel Quality Indicator), reported at the eNB by UEs (static). The reported CQI is a number between 0 (worst) and 15 (best) as listed in Table 1. The CQI indicates the most efficient MCS giving a Block Error Rate (BLER) of 10% or less. We realize that the average minimum CQI value decreases as the number of users in the multicast group increases. The result is that augmenting the number of multicast receivers clearly impacts the attainable cell throughput. Table 1 shows that a UE with the best CQI could theoretically receive 37 times the throughput of a UE with the lowest index.

This greatly motivates us to investigate methods to cope with the inefficiencies of multicast. We exploit the presence of alternative direct connectivity options available at UEs to relieve the cellular infrastructure load, while reducing the influence of UEs experiencing poor radio conditions.

## 3. JOINT D2D/MULTICAST OFFLOADING

We address the distribution of popular content to a set of $N$ mobile UEs inside a single LTE cell. Each UE is a multi-homed device that embeds both an LTE interface and a short range wireless technology that allows D2D communications

(we consider IEEE 802.11g in the paper). We want to transmit data to each UE with a guaranteed maximum *service delay D* at the minimum cost for the cellular infrastructure. In order to increase efficiency, we exploit D2D connectivity and store-and-carry forwarding. The challenging issue is that such a strategy is, by definition, unreliable, as it depends on many factors that are difficult to control (e.g., cellular channel quality, variable density of opportunistic neighbors, or interference on the D2D channel). To achieve guaranteed delivery, we consider an acknowledgment mechanism, and *panic zone* retransmissions similarly to [15]. When the service delay reaches its maximum value *D*, the eNB pushes all the missing data to uninfected nodes using unicast transmissions.

## 3.1 Cost function

What emerges from the analysis in Section 2 is that a UE with good channel quality can obtain higher bit-rates with the same amount of resource blocks (RBs), while bad channel users consume more RBs in order to transmit the same amount of data. To capture the allocation expenditure, we define the cost of transmitting to UE $v_i$ as:

$$c(v_i, k, t) = \left\lceil \frac{s_k}{\mathcal{T}_{v_i(t)}} \right\rceil, \tag{1}$$

where $s_k$ is the size in bytes of the $k$-th data block to be transmitted and $\mathcal{T}_{v_i(t)}$ is the transport block size (TBS) decided by the eNB. The cost function $c(v_i, k, t)$ measures the number of RBs needed to transmit a packet of length $s_k$ to $v_i$ at time $t$. In order to assign the MCS, and consequently $\mathcal{T}_{v_i(t)}$ (Table I, Tables 7.1.7.1-1, and 7.1.7.2.1-1 from [2]), the eNB uses the channel quality information obtained from the CQI messages that each mobile UE periodically transmits to the base station.[1]

## 3.2 Offloading strategy

The principles behind the joint multicast/D2D approach are: (1) at initial time, the eNB sends data to the $I_0$ UEs with the best radio conditions through a single multicast emission; (2) the UEs that have received the data ($I_0$ or less) start disseminating it in a D2D (epidemic) fashion; (3) before the maximum *service delay* D, we define a time interval, a *panic zone* where all the nodes that have not yet retrieved the content (either with the initial broadcast emission or in D2D fashion) receive it through unicast LTE emissions.

The proposed scheme allows all UEs to receive data by the deadline (as long as the panic zone is sufficiently large). It adapts to different *service delays* – the larger ones allowing for more D2D dissemination. Its performance relies essentially on one key parameter ($I_0$) that characterizes the number of UEs reached by the initial multicast transmission. Indeed, the eNB maintains a dynamic ranking of the UEs according to their instantaneous $c(v_i, k, t)$ values. By transmitting data with the MCS of the $I_0$-th ranked UE, the algorithm aims at reaching the best $I_0$ UEs in terms of channel quality. This immediately improves the usage of resources at the eNB, because it excludes the $N - I_0$ worst-channel UEs.

Fig. 2 offers a representative example of the proposed strategy with 6 UEs in the cell. Setting $I_0$=3, the eNB employs a MCS of 12 for the initial multicast emission. Thus, it reaches nodes with MCS of 12 and above, but leaves the

---

[1]The periodicity of CQI reports is comprised in the range $[2 - 160]$ ms in real LTE deployment.
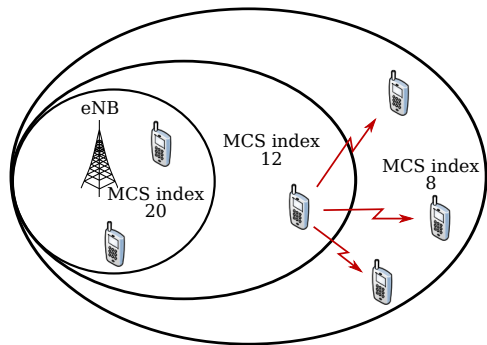


Figure 2: **UEs can decode data with a maximum modulation schema depending on their position in the cell. The eNB may decide to multicast at higher rate (E.g., MCS index 12). UEs unable to decode data are reached through out-of-band D2D links.**

three farther ones in outage (their MCS is 8). In the D2D dissemination phase, these *outaged* UEs benefit from nearby nodes, fetching data directly from them through out-of-band D2D transmissions. This cooperative strategy is by far more efficient in terms of cellular resource consumption than multicast alone, given that the transmission rate increases and the D2D links typically exploit a much larger bandwidth than cellular communications.

Here resides the novelty of our approach: *the eNB trades off the set of recipients that minimizes the multicast cost on the cellular network, while guaranteeing full coverage through D2D communications and panic re-injections when needed*. Next, we will determine the best $I_0$ values with the aid of simulations for different scenarios of utilization.

## 4. PERFORMANCE EVALUATION

### 4.1 Methodology and parameters

We compare the performance of the proposed joint distribution system with the one achieved by the classic cellular multicast alone. All the results presented in this section are averages over 25 independent simulation runs. Standard multicast implementation transmits data to all the UEs inside the cell using the MCS allowed by the lowest reported CQI value. Even in that case, UEs have no assurance of reception. The radio channel could suddenly degrade during data reception (e.g., due to fast fading or mobility), preventing certain users to correctly decode data. For this reason, we consider an additional resilience layer in the form of panic zone retransmissions, which guarantee full dissemination at the cost of much higher resource consumption.

For now, we consider a static number of UEs within the cell for each simulation run, to prove the validity of the concept. Future work will tackle the case where UEs can enter and exit the distribution area. Node mobility is implemented according to the random way-point model with speed fixed at 27 m/s and pause-time set at 0.5 s. We simulate UDP constant bit-rate downlink flows, each one with packet size $s_k = 2048$ bytes and a total load of 8 Mb. We implemented our joint D2D/multicast strategy in the ns-3. Since ns-3 does not natively support cellular multicast, we implemented an additional module that interacts with the packet scheduler to emulate single-cell multicast. The multicast module

| Parameter | Value |
|---|---|
| Cellular layout | Isolated cell, 1-sector |
| LTE downlink bandwidth | 5 MHz (25 RBs) |
| Frequency band | 1865 MHz (Band 3) |
| CQI scheme | Full Bandwidth |
| eNb TX-power | 41 dBm |
| Pathloss | Cost 231 |
| BS station height | 30 m |
| UE station height | 1.5 m |
| Fast fading | Extended Vehicular A (EVA) model |
| Multicast group size $N$ | 10, 25, 50 UEs |
| Service delay $D$ | 10, 30, 60, 90 s |
| % of direct recipients $I_0$ | 100 %, 70 %, 50 %, 30 % |

Table 2: ns-3 simulation parameters.

receives the CQI reports of UEs and decides the transmission rate following the steps explained in Section 3. We fix the bandwidth allocated for the multicast service at 5 MHz. 3GPP standard recommends not to reserve more than 60% of RBs to multicast [10], so the 5 MHz value could represent respectively the 50% or the 25% of RBs in a typical 10 or 20 MHz deployment. Other simulation parameters for the LTE cell are listed in Table 2.

Additionally, we implemented store-carry-forward routing mechanism at UEs to allow data forwarding on the WiFi interface. Regardless of its reception method, an unexpired packet can be forwarded on the WiFi interface upon meeting with neighbors. Neighbor discovery is implemented through a beaconing protocol. UEs periodically broadcast beacon messages containing their identifier and the list of buffered packets. Upon beacon reception UEs update their vicinity information and can decide to transmit a packet.

**Implementation assumption:** In simulation we make the following simplification:

- HARQ-level retransmissions and RLC-level feedback are disabled in multicast. This is a reasonable assumption: otherwise the eNB should merge the *ack/nack* messages received from all the UEs, and decide which is the best retransmission strategy. We guarantee data reception with *panic zone* retransmissions.

- The PUCCH channel is employed to acknowledge data reception towards the eNB. Panic zone retransmissions are then triggered looking at the list of received acknowledgments.

### 4.2 Reference strategies

*No D2D* is the basic strategy, where UEs have no direct connectivity options, and multicasting through the cellular infrastructure is the only means of distributing content. We compare this base case to our joint D2D/multicast strategy. We assess the performance for three different values of $N$ – the number of users inside the cell – respectively 10, 25, and 50, so to evaluate performance under different loads. We also consider various values for the parameter $I_0$ – the number of direct multicast recipients. In order to be consistent with the notation, we evaluate this value as a percentage of $N$.

### 4.3 Evaluation

**Reception Methods.** UEs may receive packets concurrently on two interfaces, using three different reception methods: multicast and unicast on the cellular interface, D2D on the WiFi interface. Fig. 3 provides the fraction of packets partitioned by their reception method. For now, we focus only on their relative weight. As expected, the fraction of packets delivered through multicast follows $I_0$. The fraction of *panic zone* and D2D messages strongly depends on the parameters $D$ and $N$. Tight service delays leave less time to opportunistic distribution to reach outaged UEs, resulting in a more intense use of panic retransmissions.

We can find a small amount of packet retransmitted during the *panic zone* even in the *No D2D* strategy. These are packets incorrectly decoded by UEs during the initial multicast emission. In the other strategies, D2D allows not to make use of retransmissions where possible, because UEs can retrieve missing packets from other UEs. For instance, the strategies *No D2D* and *100%* have the same fraction of multicast reception, but differ on the amount of panic and D2D messages. We note also that for sufficiently long *service delays*, *panic zone* is never triggered, and D2D transmissions meet the goal of guaranteeing total data diffusion. As we will show later this brings a lot of resource saving.

**Cellular Resource Analysis.** Mobile operators are primarily concerned about radio resource usage. Fig. 4 gives hints on the actual amount of RBs devoted to distribute data in the considered scenarios. Unlike previous figure, here we focus on the amount of consumed radio resources at the eNB.

The parameter $N$ strongly affects the number of employed resources. This is even more evident if we consider very short *service delays*. While the amount of resources devoted to multicast only slightly increases with the number of UEs, the impact of unicast re-injections heavily depends on the number of UEs in the cell. This happens because $N$ has a multiplier effect on unicast transmissions, and because with large probability uninfected UEs are the ones with the worst channel conditions. If we compare Fig. 3 and Fig. 4, it is impressive to note how in some cases a small fraction of unicast transmissions could translate into such a great resource usage. When $N$ is large, the choice of good values of $I_0$ becomes fundamental in order to avoid congesting the cell with too many panic retransmissions.

Another interesting result is that for any possible value of $N$ and $D$, we may always find a joint D2D/multicast strategy that offers better results than *No D2D*. For low values of $N$ and short delivery times, it is not possible to consider a great amount of outaged UEs, since opportunistic contacts between UEs are scarce, disallowing complete data dissemination before the expiration of the deadline. On the other hand, if we consider longer service times, and/or many UEs in the cell, simulation results tell us that it is possible to allow up to 70% of UEs in *outage*, reaching up to 3 times better resource efficiency. Note also that redundant multicast strategies with repeated retransmissions would never be optimal, since the amount of consumed resources would be more than double, without ensuring 100% data reception.

**Cellular Rate** The use of D2D communications increases the achievable data rate at the eNB. The proposed approach excludes from the set of direct multicast receivers the UEs in bad channel quality. These *outaged* users are reached opportunistically using D2D transmissions or through panic zone re-injections. In Fig. 5, we can evaluate the gain in terms of transport block size, with respect to the baseline *No D2D* strategy. An increase in the average TBS means that with the
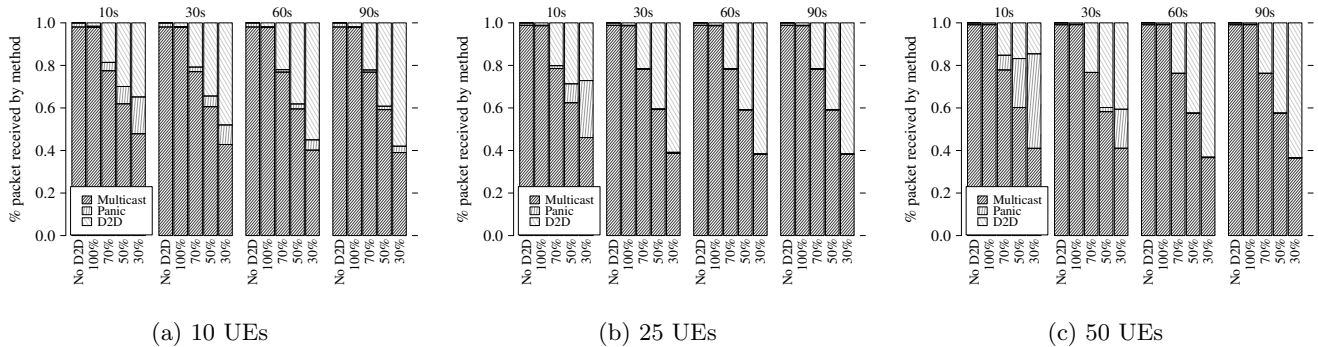
(a) 10 UEs      (b) 25 UEs      (c) 50 UEs

**Figure 3: Data packet ranked by reception method.** *Multicast* and *Panic* flows through the cellular infrastructure, *D2D* is on the WiFi channel.



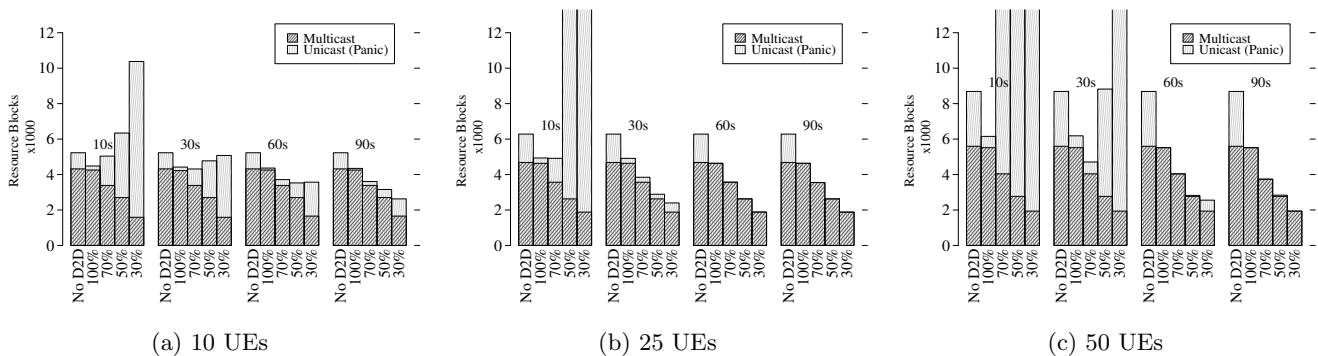(a) 10 UEs      (b) 25 UEs      (c) 50 UEs

**Figure 4: Average resource blocks employed at eNB to reach 100% dissemination. Note that even few panic zone retransmissions (in unicast) result very costly in resources.**

same amount of radio resources, the eNB can transmit more data. Average TBS is a mean to understand the quality of the cellular link connecting the eNB and the UEs. In general, the average TBS increases as $I_0$ decreases, since our strategy always serves the best placed UEs. This beneficial effect emerges if we look at the dashed curves (representing the TBS for multicast emissions only). However, the gain brought by the joint D2D/multicast strategy is often mitigated by the heavy use of panic zone re-injections to guarantee 100% data reception. This is the reason why the average TBS tends to saturate at the multicast value when the amount of panic retransmissions falls. For tight *service delays*, the opportunistic diffusion has not enough time to transfer all data packets to each UE. This forces the eNB to resort to panic unicast re-injections. The probability to have a bad channel is higher for the UEs that have not received the content, lowering the average TBS. For the *No D2D* and 100% strategies the penalty due to panic zone is negligible and never impacts the already low TBS in a noticeable manner. For larger maximum *service delay* the increase with respect to the conservative multicast-only strategy could be in the order of 2–3 times.

## 5. RELATED WORK

**Mobile data offloading.** D2D communications have been the target of intensive studies as a method to relieve the pressure on the cellular infrastructure. Typically only unicast

transmissions are considered. For instance, Han et al. identified the opportunity to save infrastructure data exploiting the social ties between users, proposing a subset selection mechanism based on contact history [7]. Similarly, Li et al. analytically formulated the problem of traffic offloading of multiple contents in a mobile environment. Under the assumption of Poisson contact, the optimal subset selection problem is solved under multiple constraints [11]. Barbera et al. analyzed contacts between end-nodes in order to select a subset of socially important VIP users, which are turned into data forwarders [3]. We proposed a simple re-injection based scheme that takes into account the evolution of the opportunistic dissemination [13]. In all these works the principal metric is the amount of data (or messages) saved on the infrastructure link. While this is an influential driver for evaluation, it does not fully represent the real amount of saved resources at the base station.

**D2D-aided multicast.** Bhatia et al., proposed the use of D2D communications to improve performance of multicast in 3G cellular networks [4]. A multihop ad hoc network is modeled analytically. A near-optimal discovery algorithm selects the best data forwarder for receivers with poor channel quality. The authors in [16] devised an algorithm to figure out the optimal number of relays inside the cluster. The paper focuses on in-band D2D communications, such that considered in [1]. Similarly, in [14], only the cluster head receives the content and is in charge of D2D retransmission inside its cluster. No hints are given on how clusters are created
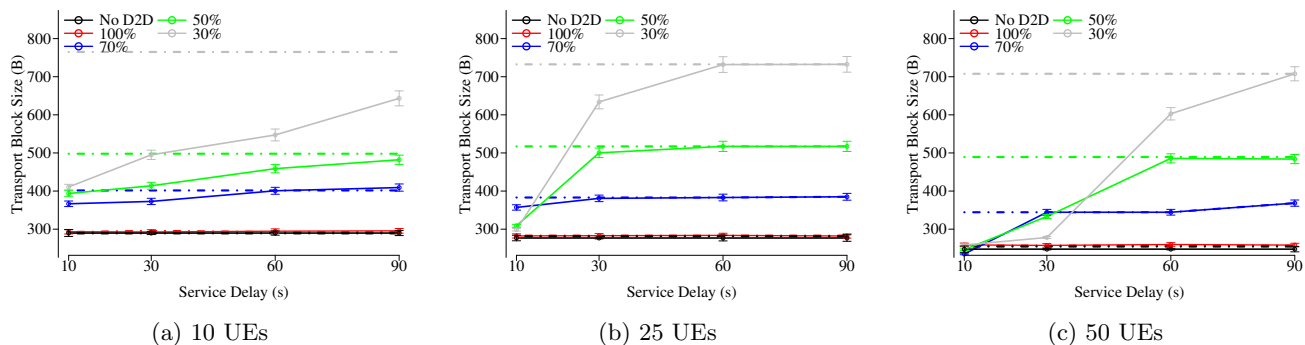
Figure 5: **Average transport block size for different** *service delays*. **Solid lines show the average (multicast and unicast), dashed lines display only multicast.**

## 6. DISCUSSION AND PERSPECTIVES

In this work, we have presented a hybrid distribution system for popular content with guaranteed delays. Multicast is a valuable option to distribute popular data into a cellular network. However, performance is limited by the channel quality of the worst UE in the cell. We proposed a framework that exploits D2D capabilities at UEs to counter the inefficiencies of cellular multicast. We evaluated the performance of a joint D2D/multicast strategy by varying the number of UEs in the cell and the maximum reception deadline. Simulation results prove that the use of D2D communications allows increasing the multicast transmission rate, saving resources and improving the overall cell throughput.

Future work will focus on the development of analytical models for epidemic data diffusion to aid the choice of which UEs to insert in the set of direct multicast recipients. Moreover, we will evaluate the scenario where multiple neighboring cells are active, and UEs can roam between them.

### Acknowledgment

## 7. REFERENCES

[1] 3GPP. TSG SA: Feasibility study for proximity services (ProSe) (release 12), 2012.

[2] 3GPP. TS 36.213 V11.2.0 Rel.11: Evolved universal terrestrial radio access (e-utra); physical layer procedures, 2013.

[3] M. V. Barbera, A. C. Viana, M. D. de Amorim, and J. Stefa. Data offloading in social mobile networks through VIP delegation. *Ad Hoc Networks*, 2014.

[4] R. Bhatia, L. Li, L. Haiyun, and R. Ramjee. ICAM: integrated cellular and ad hoc multicast. *IEEE Trans. on Mobile Computing*, 5(8):1004–1015, Aug 2006.

[5] Cisco. Cisco visual networking index: Global mobile data traffic forecast update (2012 − 2017), 2013.

[6] E. Dahlman, S. Parkvall, and J. Skold. *4G: LTE/LTE-advanced for mobile broadband*. Academic Press, 2013.

[7] B. Han, P. Hui, V. S. A. Kumar, M. V. Marathe, J. Shao, and A. Srinivasan. Mobile data offloading through opportunistic communications and social participation. *IEEE Trans. on Mobile Computing*, 11(5):821–834, May 2012.

[8] Z. Honig. Verizon demonstrating LTE Multicast during Super Bowl XLVIII (hands-on video). http://www.engadget.com/2014/01/29/verizon-lte-multicast/.

[9] F. Hou, L. Cai, P.-H. Ho, X. Shen, and J. Zhang. A cooperative multicast scheduling scheme for multimedia services in ieee 802.16 networks. *IEEE Trans. on Wireless Communications*, 8(3):1508–1519, Mar. 2009.

[10] D. Lecompte and F. Gabin. Evolved multimedia broadcast/multicast service (eMBMS) in LTE-advanced: overview and rel-11 enhancements. *IEEE Comm. Mag.*, 50(11):68–74, Nov. 2012.

[11] Y. Li, G. Su, P. Hui, D. Jin, L. Su, and L. Zeng. Multiple mobile data offloading through delay tolerant networks. In *ACM CHANTS*, Las Vegas, NV, USA, Sept. 2011.

[12] NS-3. Network simulator. http://www.nsnam.org.

[13] F. Rebecchi, M. D. de Amorim, and V. Conan. DROiD: Adapting to individual mobility pays off in mobile data offloading. In *IFIP Networking*, Trondheim, Norway, June 2014.

[14] S. Spinella, G. Araniti, A. Iera, and A. Molinaro. Integration of ad-hoc networks with infrastructured systems for multicast services provisioning. In *IEEE ICUMT*, pages 1–6, St. Petersburg, Oct. 2009.

[15] J. Whitbeck, Y. Lopez, J. Leguay, V. Conan, and M. D. de Amorim. Push-and-track: Saving infrastructure bandwidth through opportunistic forwarding. *Pervasive and Mobile Computing*, 8(5):682–697, Oct. 2012.

[16] B. Zhou, H.Hu, S.Huang, and H.Chen. Intracluster device-to-device relay algorithm with optimal resource utilization. *IEEE Trans. on Vehicular Tech.*, 62(5):2315–2326, June 2013.

The text mentions earlier: and discovered. Huo et al., proposed a cooperative multicast scheduling for 802.16 networks. A two phase schema is proposed, and all successful recipients of multicast participate in data retransmission using in-band D2D links [9].

# DISCLAIMER

*The information in this document is provided "as is", and no guarantee or warranty is given that the information is fit for any particular purpose. The above referenced consortium members shall have no liability for damages of any kind including without limitation direct, special, indirect, or consequential damages that may result from the use of these materials subject to any liability which is mandatory due to applicable law.*