**FLAVIA**
*FLexible Architecture*
*for Virtualizable wireless future Internet Access*

# Specific Targeted Research Project

# FLAVIA
## *FLexible Architecture for Virtualizable wireless future Internet Access*

# Deliverable Report

## D5.3 Analysis and design of solutions for scheduled technology enhancements

| | |
|---|---|
| Deliverable title | Analysis and design of solutions for scheduled technology enhancements |
| Version | 1.0 (final) |
| Due date of deliverable (month) | 31 Dec 2012 |
| Actual submission date of the deliverable (dd/mm/yyyy) | 15/02/2013 |
| Start date of project (dd/mm/yyyy) | 01 JUL 2010 |
| Duration of the project | 36 months |
| Work Package | WP5 |
| Tasks | 5.2 and 5.3 |
| Leader for this deliverable | BGU |
| Other contributing partners | CNIT, IMDEA, ALV, NEC, TID, IITP |
| Authors | O. Gurewitz, E. Biton, V. Mancuzo, P. Rost, P. Pillegi, J. Alonso |
| Deliverable reviewers | Giuseppe Bianchi, Ilenia Tinnirello, Maria Cristina Brugnoli (CNIT) |
| Deliverable abstract | The Document describes solutions developed in the FLAVIA project for enhancement of scheduled access technology. |
| Keywords | 802.16, LTE, throughput enhancement, Radio resource management, power control, power save, scheduling |

**FLAVIA**
*FLexible Architecture*
*for Virtualizable wireless future Internet Access*

Grant Agreement: FP7 - 257263

| Project co-funded by the European Commission within the Seventh Framework Programme | | |
|---|---|---|
| **DISSEMINATION LEVEL** | | |
| **PU** | Public | **X** |
| **PP** | Restricted to other programme participants (including the Commission Services) | |
| **RE** | Restricted to a group specified by the consortium (including the Commission Services) | |
| **CO** | Confidential, only for members of the consortium (including the Commission Services) | |

**REVISION HISTORY**

| Revision | Date | Author | Organisation | Description |
|---|---|---|---|---|
| 0.1 | 20 Dec 2012 | O. Gurewitz, E. Biton, | BGU | First draft |
| 0.2 | 15 Jan 2013 | V. Mancuso | IMDEA | Description of work carried out on opportunistic clustering and ICIC |
| 0.3 | 17 Jan 2013 | E. Biton | BGU | Introduction, outline |
| 0.4 | 23 Jan, 2013 | P. Rost | NEC | NEC contribution |
| 1.0 | 15 Feb 2013 | Giuseppe Bianchi, Ilenia Tinnirello, Maria Cristina Brugnoli | CNIT | Revisions, final editing and submission to the EC |

**FLAVIA**
*FLexible Architecture*
*for Virtualizable wireless future Internet Access*

Grant Agreement: FP7 - 257263

# TABLE OF CONTENT

# Executive summary

In this report we summarize the research activities on designing, analyzing and evaluating novel approaches to enhance the performance of existing schedule-based OFDMA systems. We provide thorough examination, how each of the proposed ideas and enhancements is supported by the flexibility and programmability of FLAVIA architecture. In particular, we examine the robustness of FLAVIA's architecture and examine its flexibility to support a variety of innovative solutions which improve the performance of schedule-based wireless communication in different aspects, and are not supported by existing platforms.

After a brief introduction, which provides a brief overview of FLAVIA scheduled MAC functional architecture and an overview of the topics tackled in this document, we provide a summarized description of each of the research topics explored by the FLAVIA partners, in the context of schedule-based OFDMA systems. Each topic is followed by a thorough discussion on FLAVIA support.

The report is split into three different research areas, each corresponding to a section:

(i) Section 2 focuses on radio resource management, in which we study both intra-cell resource allocation, i.e., how resources should be distributed within a cell (between the MSs), as well as inter-cell resource allocation, i.e., how resources should be distributed between the cells (between the BSs).

(ii) Cellular architecture and scenarios is the topic of Section 3, in which we examine new paradigms for organizing cellular systems which benefit from future cellular deployments. Among others we examine paradigms which allow more flexibility in cell size (possibly mixing different sizes), coordination between neighboring cells, incorporating different technologies, etc.

(iii) Analysis, modeling and system simulations are addressed in Section 4, in which we develop a set of tools for the validation of the proposed enhancements. The tools presented in this section are an important building block in the design and evaluation of FLAVIA's architecture, as they not only provide the platform to evaluate the innovative solutions suggested throughout this document and others, but also provide a tool to verify and analyze the flexibility of FLAVIA to support such enhancement for existing and potential technologies and standards.

# 1 Introduction

Work Package 5 (WP5) focuses on the design of novel approaches to enhance the performance of schedule based OFDMA systems, regardless of the employed technology. The suggested novel solutions should exploit FLAVIA's architecture, and prove its flexibility.

As detailed in Deliverable D3.1.1 and D3.1.2, FLAVIA's high level MAC architecture for scheduled systems is composed of modules, each representing a service that a scheduled access technology is expected to provide. Services can have common functions, i.e., a function can be reused by multiple services. The high-level FLAVIA functionalities for scheduled access systems (e.g., 802.16 and LTE) are depicted in Figure **1**. In particular, the figure below shows the main services needed for the implementation of scheduled access.



**Figure 1: FLAVIA scheduled MAC functional architecture.**

In the first part of the project, we started development and performance evaluation of novel approaches and solutions which exploit the flexibility of the FLAVIA architecture. The results of this work were reported in the previous deliverable D5.2.

In the second part of the project, we continued the work on improvement and further analysis of the approaches proposed in D5.2 and also developed and evaluated a number of new solutions.

We divide the research results into three main aspects:

I. radio resource management – focuses on enhancements of one or more services that are in the heart of the radio resource management layer

II. cellular architecture and scenarios – suggests new architectures and scenarios that are made possible by FLAVIA's flexible architecture

III. Analysis, Modeling and System simulations – provides tools to analyze and simulate the proposed enhancements

## Radio Resource Managements

Resource management plays a significant role in communication networks and particularly in scheduled wireless cellular networks. Under this research field we study both intra cell resource allocation, i.e., how resources should be distributed within a cell (between the MSs) as well as inter-cell resource allocation, i.e., how resources should be distributed between the cells (between the BSs). In particular, for the intra-cell case, we study (i) uplink scheduling and user pairing in collaborative MIMO, (ii) power save enhancements in IEEE 802.16e networks, (iii) prioritization of CQI/Sounding resources, and (iv) distributed opportunistic Scheduling. For the inter-cell case we study (i) joint scheduling and power control with noise rise constraints and (ii) Inter-cell interference management via base station scheduling.

### Collaborative MIMO user pairing

Our research on collaborative MIMO (CoMIMO) focuses on how to provide the capability in which two user terminals transmit simultaneously to one base station at the same UL time & frequency resources using a single antenna per MS. In such a configuration there is a high probability that pairing will work where a MIMO (matrix B) from a single MS will fail. Using CoMIMO, if average station rate is kept, then cell capacity can be increased without increasing the power per station (so more users can be served). The research focused on the development of the Co-MIMO pairing algorithm. The proposed algorithm selects pairing candidates in different decision levels: first PHY coupling metric are checked, then Maximum MCS gap between allocations, then Joint Interference budget, then Power control per MS, then Minimum allocation size (bytes) and finally Resource utilization degree.

### Power-save in IEEE 802.16e networks

The mobility of the IEEE 802.16 (WiMAX) cellular network implies limited accessibility to power, hence, power-saving techniques are crucial for conserving the power of the mobile terminals. In this work we provide a comprehensive measurements study of power-save mechanisms implemented

in current deployed networks. Particularly, we examine the power consumption of the various transmission and reception modes, and compare them with the power consumption during Idle and Sleep modes. We show both experimentally as well as analytically, that even though theoretically power consumption can be dramatically reduced by employing an efficient algorithm which alternates between power-save and active modes, lack of cross layer coordination between the applications, the operating system and the network card prevents efficient power save implementation. We suggest and implement a simple proactive buffering solution that delays the sporadic traffic generated by the upper layers, when the device is in Idle mode, and show that such simple enhancement can dramatically reduce the device power consumption.

We present a different approach for power-save operation termed *Intra Frame Power Save (IFPS)*, which does not require any cross layer coordination. We show how IFPS can dramatically reduce the power consumption even while the device is in operational mode, and suggest ways for further reducing power consumption by taking into consideration mobile devices supporting IFPS mode while performing the schedule by the base station.

## Prioritization of CQI/Sounding resources

Our research on feedbacks considers the allocation of CQI and sounding resources. Specifically, it considers the prioritization of CQI/Sounding resources when resource shortage is encountered within the BS. The approach was (i) to filter MS's with negligible traffic from having sounding resources, and (ii) to maximize sector throughput by prioritizing (near) full loading MS's.

## Distributed Opportunistic Scheduling in Non-Homogeneous Networks

In this work, we design novel distributed scheduling algorithms for multi-user Multiple Input Multiple Output (MIMO) systems and evaluate the resulting system capacity analytically. In particular, we consider algorithms which do not require sending channel state information to a central processing unit, nor do they require communication between the users themselves, yet, the resulting capacity closely approximates that of a centrally-controlled system, which is able to schedule the strongest user in each time-slot. In other words, multi-user diversity is achieved in a distributed fashion. Our analysis is based on a novel application of the Point-Process approximation. This technique, besides tackling previously suggested models successfully, allows an analytical examination of new models, such as non-homogeneous cases (non-identically distributed users) or various Quality of Service (QoS) considerations. This results in exact expressions for the capacity of the system under these schemes, solving analytically problems which to date had been open. Possible

applications include, but are not limited to, modern 4G networks such as 3GPP LTE, or random access protocols.

## Joint scheduling and power control with noise rise constraints

Consider the problem of joint uplink scheduling and power allocation. Being inherent to almost any wireless system, this resource allocation problem has received extensive attention. Yet, most common techniques either adopt classical power control, in which mobile stations are received with the same Signal-to-Interference-plus-Noise Ratio, or use centralized schemes, in which base stations coordinate their allocations.

In this work, we suggest a novel scheduling approach in which each base station, besides allocating the time and frequency according to given constraints, also manages its uplink power budget such that the aggregate interference, ``Noise Rise'', caused by its subscribers at the neighboring cells is bounded. Our suggested scheme is distributed, requiring neither coordination nor message exchange.

We rigorously define the allocation problem under noise rise constraints, give the optimal solution and derive an efficient iterative algorithm to achieve it. We then discus a relaxed problem, where the noise rise is constrained separately for each sub-channel or resource unit. While sub-optimal, this view renders the scheduling and power allocation problems separate, yielding an even simpler and more efficient solution, while the essence of the scheme is kept. Via extensive simulations, we show that the suggested approach increases overall performance dramatically, with the same level of fairness and power consumption.

## Inter-cell interference management via base station scheduling

The continuously increasing demand for higher data rates results in increasing network density, so that inter-cell interference is becoming the most serious obstacle towards spectral efficiency. Considering that radio resources are limited and expensive, new techniques are required for the next generation of cellular networks, to enable a more efficient way to allocate and use radio resources. In this framework, we target the design of a frequency reuse 1 scheme, which exploits the coordination between base stations as a tool to mitigate intercell interference by separating the scheduling from resource allocation. While common approaches proposed in the literature focus on the optimal user scheduling, we tackle the problem from a different angle. In particular, we formulate a base station scheduling problem to decide whether a base station is allowed to transmit to any of its users in a given subframe, without causing excessive interference to any of the users of other scheduled base stations. To this aim, we show that finding the optimal base station

scheduling is NP-hard, and formulate the BASICS (BAse Station Inter-Cell Scheduling) algorithm, a novel heuristic to approximate the optimal solution at low complexity cost. By means of numerical and packet-level simulations, we prove the effectiveness and reliability of the proposed solution as compared to the state of the art of inter-cell interference mitigation schemes.

## Cellular Architecture and Scenarios

Currently deployed cellular networks still rely on a macro- and micro-cellular architecture without advanced cooperation and coordination possibilities. However, this changes with the development and deployment of next-generation mobile communication systems such as IEEE 802.16m and 3GPP LTE-A. These systems rely on the support of relay nodes, femto-cells, as well as multi-cell cooperation. Furthermore, future deployments will be characterized by a tremendous increased cell density. This requires new approaches to organize cellular systems in order to benefit from the future cellular deployments. These approaches will heavily rely on novel MAC protocols and coordination techniques supporting the operation of cellular networks. Among others, those protocols will allow for more inter-BS coordination and information-exchange. Furthermore, they provide more flexibility which is a pre-requisite to implement future-proof algorithms. The FLAVIA framework shows one way to provide the means to implement scalable and flexible algorithm for future cellular networks. This section presents algorithms that exploit FLAVIA's flexibility and provide significant benefits in future cellular networks.

### Opportunistic scheduling with clusters

Opportunistic scheduling was initially proposed to exploit user channel diversity for network capacity enhancement. However, the achievable gain of opportunistic schedulers is generally restrained due to fairness considerations which impose a tradeoff between fairness and throughput. In this study, we show via analysis and simulation that opportunistic scheduling not only increases network throughput dramatically, but also can be fair to the users when they cooperate, in particular by forming clusters. We propose to leverage smartphone's dual-radio interface capabilities to form clusters among mobile users, and we design simple and scalable cluster-based opportunistic scheduling strategies which would incentivize mobile users to form clusters. We use a coalitional game theory approach to analyze the cluster formation mechanism, and show that proportional fair-based intra-cluster payoff distribution would bring significant incentive to all mobile users regardless of their channel quality.

### Dynamic assignment of UL/DL sub-frames in TDD systems

A prerequisite to obtain full spatial reuse is to cancel or mitigate inter-cell interference, which is limiting the cell throughput. In order to derive new strategies for cellular interference-mitigation, this section analyzes achievable uplink-downlink data rates for different inter-cell interference scenarios. Among others, this section discusses an asymmetric protocol exploiting cross-uplink-downlink interference, i.e. two adjacent cells do not operate simultaneously in uplink or downlink but only one of both is active in uplink and one is active in downlink. Using analytical results, it is shown that under specific conditions this approach provides performance gains over conventional approaches and close to multi-cell MIMO. Although this approach is not able to improve the performance under all channel conditions, it provides a new degree of freedom which might be exploited if inter-cell interference significantly impairs the performance.

### Coded unicast downstream traffic in a wireless network

In this study, we design, analyze and implement a network coding based scheme for the problem of transmitting multiple unicast streams from a single access point to multiple receivers. In particular, we consider the scenario in which an access point has access to infinite streams of data to be distributed to their intended receivers. After each time slot, the access point receives acknowledgments on previous transmissions. Based on the acknowledgements, it decides on the structure of a coded or uncoded packet to be broadcast to all receivers in the next slot. The goal of the access point is to maximize the cumulative throughput or discounted cumulative throughput in the system.

We first rigorously model the relevant coding problem and the information available to the access point and the receivers. We then formulate the problem using a Markov Decision Process with an infinite horizon, analyze the value function under the uncoded and coded policies and, despite the exponential number of states, devise greedy and semi-greedy policies with a running time which is polynomial with high probability. We then analyze the two users case in more detail and show the optimality of the semi-greedy policy in that case.

## Analysis, modelling and system simulations

### Analysis of fundamental trade-offs in scheduling and resource allocation

This research derives an analytical framework to evaluate achievable rates in a multi-user OFDMA system. Opportunistic schedulers and OFDMA allow for exploiting multiuser diversity using a flexible and simplified resource assignment. However, the granularity of resource assignments determines the

gains in spectral efficiency and requires to consider signaling overhead and finite channel coherence. Additionally to the signaling and pilot overhead, practical systems suffer from partial CSI, i.e., neither scheduler nor receiver have exact knowledge of the channel state. This work investigates the trade-off between achievable net-rates and channel characteristics as well as system parameterization. An analytical framework for the expected net rates in a multi-user system with partial CSI at scheduler and receiver is derived.

## Traffic-Centric modeling methodology

The Traffic-Centric Modelling solution technique was successfully applied to develop a model solution of LTE technology. The result is an initial flexible extensible framework that can be used to develop traffic-centric model solutions of different LTE scenarios. The features of the implemented model solution is very modular and can be programmed to mimick the FLAVIA architecture, allowing developers to test new protocols and configurations before actually applying it to their FLAVIA-certified network.

## System level simulation

Finally, in order to demonstrate the applicability of the investigated algorithms, impacts of the software interfaces, and the performance benefits in a 3GPP LTE system, established a credible system level simulator. The system level simulator has been calibrated with respect to large scale and small scale parameters.

# 2 Radio Resource Management

## 2.1 Collaborative MIMO user pairing (ALV)

### 2.1.1 Problem statement and motivation

In this research work we targeted the possibility to obtain MIMO using two user terminals that would transmit simultaneously to one base station at the same UL time & frequency resources (see Figure 2).



*Figure 2: Collaborative MIMO scenario*

This makes it possible to implement MIMO (matrix B) with one Tx antenna per MS.

### 2.1.2 Proposed solution

In the proposed algorithm we have tried to preserve the existing iterative scheduling-frame building flow. We have tried to keep it as simple as possible (i.e. non optimal) targeting random pairing matching. Other design constraints were to have minimum (frame building) rollback phases and to generate seeds for future extension to more efficient pairing matching algorithm.

In the I\iterative operation Scheduling-frame building phase the operation to get the next quota provides the full resource allocations (HARQ, maps, data resources, etc). The algorithm performs one time allocation per connection including unified DL & UL scheduling in which map reservations are done on the fly. Concerning Co-MIMO pairing, the algorithm looks for the higher channel spatial separation that would provide better pairing performance. It provides the same resource allocation size, joint interference budget and power control per coupled MS.

Figure 3 describes the main blocks of Co-MIMO algorithm. It calculates the next quota in the uplink and checks the possibilities to obtain MIMO pairing in 3 stages.

In the first stage the new UL quota is checked for MIMO candidacy first verifying the Co-MIMO Capabilities (SBC - TLV 177), then the QoS type (only BE is considered), then the allocation size demand in bytes (post limitations, pre-MRT), the radio conditions and finally filtering multi-burst. It can be decided that the quota will be SIMO-handle or, if the verifications were passed to go on with the Co-MIMO candidacy in a 2nd stage the algorithm.

In the Algorithm's 2nd stage we try to build Co-MIMO pairing set candidacy obtaining pairing candidates subset among already (quasi) scheduled unpaired Co-MIMO candidates. Due to current ("one shot") UL implementation and UL Co-MIMO candidacy conditions, in most cases one BE allocation will occupy the entire UL subframe data resources (i.e. one candidate). For every unpaired Co-MIMO candidate, two sequential pairing matching processing applies, the first does PHY coupling metric calculation and the second is related to coupling metrics advisory. The PHY coupling metric reflects the channel separation of pairing candidates by checking the MCS couple fixed margins and CQI measurements. Regarding the coupling metrics, the MRT inbound parameters include both pairing candidates parameters, PHY coupling metric and the number of subchannels (the minimum number of subchannels as existing unpaired Co-MIMO candidate and the maximum number of subchannels up to UL subframe resource limitations). The MRT Co-MIMO pairing matching criterions include PHY coupling metric, maximum MCS gap between allocations, Joint Interference budget, power control per MS, the minimum allocation size (bytes) and finally the resource utilization degree that can be reflected within 3rd algorithm stage. If PHY coupling metric is higher than a given threshold, further pairing processing is forwarded to MRT, otherwise the pairing is failed. The MRT output includes allocations attributes of pairing candidates. SIMO scheduling metrics are calculated as well as input for 3rd comparison stage

*Figure 3: Co-MIMO algorithm*

The algorithm 3rd stage is that of scheduling selection. In this stage, we compare among potential pairing candidates and SIMO taking the one with highest throughput gain (i.e. max. aggregated bytes). In the case of Co-MIMO pairing we need to add the MAP IE delta size reservation (Co-MIMO IE – SIMO MAP IE) and to finalize Co-MIMO including bursts allocations, the notification of final satisfied quotas to observers (scheduler, HARQ,…) and the burst descriptors. In the case of Quasi SIMO we have to mark resources/allocation as Co-MIMO candidate and to map IE reservation as for SIMO. Regarding the frame building procedures we have 2 options. The first consists of fully frame building procedures which in case of future Co-MIMO pairing, rollback for allocation attributes (MCS, reduced number of bytes, increased number of subchannels, etc…) and burst rebuild as well as final quota observers notifications will be required Second option consisting of partial frame building

procedures which include data resources & HARQ buffering reservation, burst allocation and notifications of satisfied quotas to observers (that are postponed up to final pairing/non-pairing result) and, at the end of frame building phase, unpaired Co-MIMO candidate allocations which should be finalized as SIMO allocations.

During scheduling-frame building phase free resources should be managed both for frame resources and Co-MIMO resources. Finalizing phase is required at the end of frame building phase in case of existing unpaired Co-MIMO candidates.

Some frame building aspects need to be taken into account. There will be separate HARQ region definitions for SIMO and Co-MIMO allocations pointed from different HARQ MAP IEs as can be seen in Figure 4. UL MIMO HARQ Subburst MAP IE is used for Co-MIMO allocations. HARQ MAP IEs reservations shall take into account all overheads in case of multi-HARQ regions.



*Figure 4: Separate HARQ regions for SIMO & Co-MIMO*

HARQ retransmissions for initially Co-MIMO pairs are done using SIMO allocations. HARQ budget for scheduled for unpaired Co-MIMO candidates should be reserved and finalized at the end of pairing decision.

### 2.1.2.1 FLAVIA support

This work can be taken as an example of how the enhanced modularity provided by FLAVIA's architecture for scheduled systems can be used to easily implement this approach to get MIMO based on the collaboration of MSs. Our approach tries to benefit from a feature usually not widely exploited in scheduled systems like it is Collaborative MIMO. In that sense, FLAVIA's

approach adds a lot of flexibility to exploit MIMO capabilities of scheduled systems.

### 2.1.3  Some insights on resulting advantage and benefits

The system was tested with a setup according to the scheme depicted in Figure 5. The demo run in an indoor conducted setup with 4 CPE's. Each MS was connected directly with RF cable to a different Tx/Rx antenna.



*Figure 5: Co-MIMO testing setup*

The channel plan used carrier frequency 2.6635GHz being the Bandwidth 10MHz. The BS Tx Power was configured to 38dBm. The TDD ratio was 29/18 DL/UL symbols in frame. Regarding the Air Link condition, all the CPE's had RSSI between (-35) to (-45) dBm and CINR: higher than 33dB.

UDP traffic was generated and sent over each of the UL service flows. The traffic rate was 5MBps for each MS. Data Packets size was constant and of 512 bytes. DL traffic was sent simultaneously for each MS. The system was configured with Tx rate 64QAM 5/6 and to generate no NACKS.

Based on such configuration, the expected results were that the throughput per MS would be approximately 3.9Mb/sec. and that the aggregated throughput would be approximately 15.5Mb/sec. This prediction fits with the obtained results that can be observed in Figure 6.

***Figure 6: Co-MIMO testing results***

The proposed collaborative MIMO scenario provides with high probability that a pairing will work where a MIMO from a single MS will fail. In such deployments, channel is pretty much constant and pairing can be assigned statically in advance. If deployment is not interference limited (e.g., directional antennas, high reuse, etc.), then capacity can be doubled, paying the price of increased power per station due to MIMO losses (pairings that are less than optimal, channel estimation losses), around 0.2 dB

If SNR is not an issue because the highest possible station rate has been reached anyway, then the deployment capacity can be effectively doubled

This way of getting MIMO has also some disadvantages. First of all transmitters can't take advantage of a known channel (pre-coding is not possible). The scheduler is more complex  because stations have to be paired and additional HARQ issues apply.

In the receiver there could be modem issues, such as different frequency shift for each MS. Finally, interference rejection performance (when applied) is degraded due to a decrease in the number of degrees of freedom.

Future work can include the enhancement of Co-MIMO matching algorithm using a scheduling subset approach and the inclusion of PHY coupling metric based on CQI reports.

## *2.2 Power-save in IEEE 802.16e network (BGU)*

### *2.2.1 Problem statement and motivation*

Next generation cellular technology (3.5G going on 4G) is emerging rapidly. For example Clearwire, a leading provider of 4G wireless broadband services in the U.S., has recently re-ported that 2011 ended with approximately 10.4 million subscribers, which is a 140% growth year over year from 4.3 million. UQ Communications announced reaching two million subscribers to UQ WiMAX in Japan. One of the biggest advantages of such technologies is their ability to provide widespread coverage providing connectivity all over. Consequently, users, mobile or static, can expect high speed mobile broadband services everywhere, even when distant from Access Point (AP) and not necessarily in the proximity of power supply. In addition, due to the portability, devices are expected to be a carry on, hence smaller and lighter with small batteries. Accordingly, power saving mechanisms are crucial techniques for conserving the power of the mobile terminals.

In this study we conduct a thorough measurement investigation of the power-save mechanisms implemented in current deployed networks, and examine their suitability to various applications. We investigate a new approach for power-save operation termed Intra Frame Power-save (IFPS), in which devices can enter power-save mode for much shorter intervals than the common Idle/Sleep modes (i.e., intra-frame resolution). We suggest ways for utilizing IFPS for reducing power consumption.

In particular, the contributions of the study are as follows: We examine the efficiency of power saving modes on an operational WiMAX network with widely deployed WiMAX network cards. We show that due to lack of cross layer coordination between the application, the operating system and the network card, the efficiency of the power save mechanisms as defined in the standard is relatively low. Particularly, keep alive messages of always on background applications prevent the card from activating power saving mode more frequently and for longer periods. We measure and characterize the power consumption of three commercial WiMAX devices, while performing Rx and Tx for various kinds of signals (e.g., data reception and transmission, map reception, passive listening, etc.). We identify an *Intra-Frame Power-Save* (*IFPS*) mechanism. Based on IFPS devices stay tuned to receive the maps and can participate in ordinary operations such as transmitting control and feedback signals. Nonetheless, when not transmitting or receiving data a device reduces its power consumption by switching off its transceivers. We show that the power consumption with IFPS is comparable with Sleep mode, but is free of the complexity to schedule and maintain the Sleep cycles.

Obviously, IFPS requires no latencies in entering or exiting power-save mode. To further utilize IFPS power saving, we suggest a scheduling mechanism that considers IFPS and maximizes the overall Intra-frame power saving potential both in the downlink and the uplink allocations. While in the downlink we show that a simple algorithm obtains the optimal allocation, in the uplink, due to a device's maximal transmission power, we show that the IFPS scheduling problem is NP-hard, and we suggest a heuristic and provide a tight lower bound.

Even though our measurements were performed on widely deployed devices running IEEE 802.16e-2009, the results and conclusions brought in this study not only apply to pre 4G WiMAX networks, but also to 4G networks based on 802.16m and the Long Term Evolution (LTE) of the 3GPP. This is because in both 4G technologies, there are similar protocols, (though not identical) for power-save; Sleep and Idle (with small changes in 802.16m), and the mode of discontinuous reception (DRX) in LTE.

### 2.2.2 Cross-layer traffic buffering for enhanced Idle mode performance

In this study we show that the efficiency of the power-save mechanisms is limited due to sparse traffic generated at the MAC layer as well as at higher layers. We further show that the efficiency of these power-save mechanisms can dramatically degrade as a function of the number of 'always-on' applications running in the background. Lack of coordination between the MAC layer and the upper layers as well as the operating system can cause an application to generate a "keep alive" message shortly after the device enters Idle mode, thus reducing the device's power-save efficiency dramatically.

Addressing this problem, we aim at better synchronizing the applications' background traffic and the MAC Idle intervals. This is done by buffering the packets arriving on the Uplink direction while the device is in Idle mode, and by that prevents the WiMAX card from exiting Idle allowing a longer power save interval.

#### 2.2.2.1 FLAVIA support

Implementing a simple buffer at the MAC layer, as suggested in this research, is made very easy with FLAVIA. Specifically, it can be easily integrated into the Data Transport service. However, traffic buffering could also be implemented in traditional protocol stack design, though it might require more effort. Nevertheless, the proposed traffic buffering is just a simple solution that aims at illustrating the potential gain in synchronization of the data transport with the MAC layer and with the Power Saving service in particular. This is only supported by a flexible MAC such as designed in FLAVIA. Specifically, the

unique FLAVIA interfaces between the Data Transport service, the QoS Strategy and the Power Saving service, allows a more advanced mechanism that synchronizes 'keep alive' messages (of all applications) and management messages with the power save mode windows, such that the device stays longer in power save mode and activates power save more frequently. On the other hand, when the Data Transport and the QoS Strategy services identify urgent messages, it informs the Power Saving service to resume to normal operation.

### 2.2.3  Some insights on resulting advantage and benefits (or if applicable trade-offs)

We evaluate the suggested solution by modifying the WiMAX driver such that, when the card reports Idle mode activation, the driver starts buffering all the packets in a special buffer for a predefined duration. We term this predefined duration 'hold-time', and we examine different hold times of 8, 16, 28 and 40seconds. We repeat the measurements for Skype and Dropbox, while keeping the baseline of having the device connected with no application in the background. For each hold time, we verified that the buffer does not affect the application's "keep alive" mechanism, e.g., throughout the tests the status of the Skype user was examined both locally (at the device itself) and remotely (as observed by remote users), and showed the user always online. Similar verification was done for Dropbox operation.

Note that as soon as the application becomes active, e.g., initiation of a Skype call, the WiMAX card deactivates Idle mode and the buffer is disabled. Accordingly, while the operation of active applications is not affected, the buffer hold time may impact the time an application becomes active (e.g., the time to initiate a Skype call).

Figure 2 summarizes the results for the 8 second hold time. The left bars relate to the base case without a buffer and the right bars depicts the behavior with the buffering mechanism. The figure clearly shows the increase in the time the system is in Idle mode, for example, 64% instead of 29% with Skype.

Our findings clearly show that the operating systems as well as the application layers may implement better mechanisms to utilize the MAC layer power save capabilities. For instance, an efficient buffer would synchronize packets from different applications while classifying the packets by their sizes and importance in order to identify when the device shall be awakened from Idle mode, and what is the maximal delayed period per packet in the buffer.

| | No Application | | | Dropbox | | | Skype | |
|---|---|---|---|---|---|---|---|---|
| Idle | 71.90% | 85.43% | | 22.19% | 35.86% | | 29.85% | 64.67% |
| Sleep | 11.07% | 7.18% | | 57.45% | 44.16% | | 48.62% | 22.45% |
| Online | 17.03% | 7.39% | | 20.36% | 19.99% | | 21.53% | 12.88% |

*Figure 7:   Power States with a buffer solution - Online, Sleep and Idle ratios without a buffer and with a buffer hold-time of 8 seconds*

More details are available in [1].

### 2.2.4  Intra-Frame Power Save Scheduling

WiMAX scheduling or resource allocation relates to the task in which the base station decides how to divide its time and frequency resources per all the users it serve. At each frame, the base station solves the scheduling problem and allocates its resources to the users. Each allocation defines the transmission power, the Modulation and Coding Scheme (MCS), and the allocated resource units of the frame. The process of user selection and resource allocation is influenced by many parameters such as user requirements (e.g., bandwidth requirements), application requirements (e.g., delay and jitter limitations), channel condition, scheduling criteria (e.g., fairness, max throughput), and more. Typically, the scheduling process comprises two processes: (i) scheduling and link adaptation and (ii) frame building. The scheduling and link adaptation is responsible to select the users to transmit/receive, the allocation size, the transmission power and the MCS.  The frame builder is responsible for the logical positioning of the allocation within the frame.  In this study we deal

only with the frame builder and only with the objective of power saving. Specifically, we propose two schemes, one for the downlink and one for the uplink, that optimize the frame allocation in terms of mobile power save.

We consider the two dimensional (2D) frame (the upper parts of Fig. 4.3 describe a typical frame) with $K$ representing the time and $M$ the frequency. The 2D basic resource units are called slots. We consider the set of allocations $A = \{a_i\}$, the output of the scheduler, serving as input to the frame builder. The notation $a_i$ represents the size of the allocation assigned to user $i$, in units of slots. We assume that the scheduler ensured that there is a feasible allocation for which $\sum_{i=1}^{n} a\_i \leq M \cdot K$.

The output of the frame builder consists of an $M \times K$ matrix $\boldsymbol{R}$, which represents the assignment of resource units to the users, i.e., $R_{m,k} = i, \; i \in \{0,1,2,\dots,N\}$ iff user $i$ is assigned with sub-channel $m$ on time slot $k$, where 0 represents that sub-channel $m$ on time slot $k$ was not assigned to any user. Finally, we denote by $e_i$ the (downlink/uplink) transmission completion time in time slots of the $i$'th user's allocation $e_i(\boldsymbol{R}) = \{\max k: \exists m, R_{m,k} = i\}$.

In contrast to the downlink behavior, in which the MS power consumption increases at the beginning of the downlink sub-frame, on the uplink, since the information on the uplink allocation is available in advance (e.g., in 802.16e the allocation of the current frame uplink transmission is transmitted at the UL-MAP of the previous frame), the device can stay on low power just prior to transmission time and reduce the power right after transmission. In addition, despite the fact that at least seemingly a device transmission power is more or less a linear function of the resources assigned to the device, e.g., two resource units will consume twice as much power as one resource unit, surprisingly, one can see that transmitting on one subchannel (one time slot) consumes almost the same power as transmitting on all subchannels of a one time slot. Both these observations suggest that, in contrast to power consumption during the downlink sub-frame which is determined by the end of the device downlink allocation, the uplink power consumption is determined only by how many time slots the allocation spans. This latter assertion implies that uplink allocations should be assigned according to frequency first (vertical) allocation.

Despite the aforementioned motivation of granting narrow time allocations to MSs, it is important to note that each MS is constrained in its maximal transmission power. Note that this max-power constraint determines for each device the maximal number of subchannels it can be assigned in each time slot. It is important to note, that in order to avoid crossing this max-power constraint, IEEE 802.16e specifies uplink allocations that meander horizontally in a time-first manner, as opposed to the downlink frequency-first allocation. Accordingly, the uplink allocation is spread for a longer time over fewer

frequencies. Clearly, even though this approach makes the tasks of both the scheduler and frame builder much easier, it is deficient regarding saving MSs power.

Denoting by $l_i$ the uplink allocation start time in time slots of the user $i$ allocation,

$l_i(\boldsymbol{R}) = \{\min k: \exists m, R_{m,k} = i\}$. The frame builder aims at minimizing the sum of all user transmission time spans to allow maximum intra-frame power save for all devices, while keeping the max-power constraint.

Formally, the uplink Intra-frame power-save scheduling problem is as follows:

$$\underset{\boldsymbol{R}}{\text{minimize}} \sum_{i=1}^{N} e_i(\boldsymbol{R}) - l_i(\boldsymbol{R})$$

$$s.t. \sum_{m=1}^{M} \sum_{k=1}^{K} \delta_{m,k}^{(i)} = a_i \forall i; \sum_{m=1}^{M} \delta_{m,k}^{(i)} \leq \quad (i) \ \forall i, k$$

where $H(i)$ denotes the maximal number of subchannels user $i$ can be assigned in each time slot.

The above allocation problem is NP-hard.

Accordingly, we propose a sub-optimal heuristic that solves the uplink power-save allocation problem. The idea behind the algorithm is to split the uplink sub-frame into horizontal (frequency) band of maximal size such that for each band there exists a set of devices (the union of all the sets contains all devices), such that: (i) the total allocations assigned to these devices by the scheduler are greater or equal the number of resources constrained in the band (ii) the max-power constraint of each of the devices in the set, is greater or equal to the width (number of sub-channels) in the band. Note that, under these constraints, each set of devices can be scheduled according to meandering in a frequency-first manner within its respective band.

More specifically, the IFPS uplink frame builder algorithm sorts the allocations $a_i \in A$ according to their max-power constraint, and creates a sorted list $s_i \in S \ s.t. H(s_1) \leq H(s_2) \leq \cdots \leq H(s_N)$. Next, it seeks the maximal size frequency band $\mu \leq M$ such that $\exists \eta; s.t. \sum_{j=\eta}^{N} a_{s_j} \leq \mu \cdot K$ and $H(s_i) \leq \mu; \forall s_i \geq \eta$. Then, it allocates the ordered $s_i$ each at a time in a time-first meandering pattern within the frequency band. The process is returned until all devices are scheduled.

A pseudo code and more details are given in [1]

### 2.2.4.1 FLAVIA support

FLAVIA scheduled based system architecture distinguishes between two scheduling services, namely (i) Scheduling Strategy and (ii) MAC scheduler,
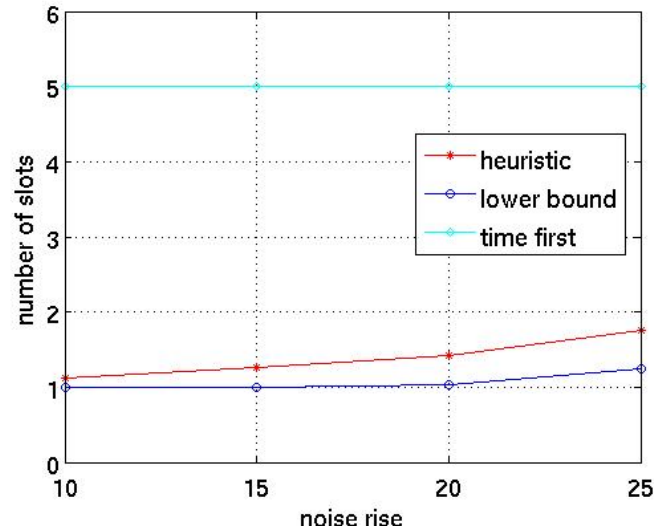
and defines the interfaces between them. In this power save enhancement we suggest a modification of the frame building process which is part of the MAC scheduler service. Accordingly, FLAVIA modularity enables us to modify one of the functions of the MAC services without the hassle of changing any other services. Accordingly, the implementation of the suggested enhancement becomes immediate and quite easy.

### 2.2.5 Some insights on resulting advantage and benefits

In order to examine the performance of the suggested heuristic in a realistic scenarios, we simulated it for a typical WiMAX uplink subframe sized 5x35 tiles. In more details, we simulate a two tier hexagonal deployment with 19 sites, each containing three cells (sectors). Each cell takes its scheduling decisions independently. MSs are located uniformly across the deployment with 20 MSs in each cell. Statistics are collected from all cells. We assume the WINNER II stochastic channel modeling for the urban macro-cell scenario, where typically MSs are located outdoors at street level and BSs are fixed clearly above surrounding building heights at a distance of 500m between them. The MS and the BS antennas gain are 0 and 17 dBi, respectively, and the assumed noise figure at the MS and at the BS are 7 and 5 dB, respectively. We considered the proportional fair scheduling, where the $k$-th tile is allocated to the user $i^*$ which maximizes $\frac{r_i}{T_i(k)}$, where $r_i$ is the user $i$ instantaneous rate and $T_i(k)$ the user $i$ average throughput. $T_i(k)$ is given by $T_i(k) = \beta T_i(k-1) + (1-\beta)B_i(k-1)$ with $\beta$ is the decay window and $B_i(k-1)$ is the number of bits delivered to MS $i$ at tile $k-1$ (0 if MS $i$ was not scheduled at tile $k-1$).

The user transmission power is derived by the noise rise constraint approach (see Subsection 0), where the maximal transmission power is set to 24 dBm.

Figure 3 depicts the average length of the allocation in time slots as a function of the noise rise for the traditional time first allocation as well as based on our heuristics. Additionally, the figure depicts a lower bound for allocation length in time. One can see that as the noise rise increases the allocation time also increases. Here the increase in noise rise requires higher transmission power, thus users reach their power limitation with fewer sub-channels. Thus the same allocation spans over less sub-channels and over more time slots.

***Figure 8: length of allocation in time slots as a function of the noise rise***

Figure 4 depicts the average length of the allocation in time slots as a function of the minimal allocation size. Here, increasing the allocation size requires increase in the allocation span in time slots.

In both figures, the advantages of the suggested heuristic compared to the traditional time first allocation is clearly presented. Furthermore, it is not far from the optimal (at case unachievable) lower bound.

*Figure 9: length of the allocation as a function of the minimal allocation size*

## 2.3 Prioritization of CQI/Sounding resources (ALV)

### 2.3.1 Problem statement and motivation

We would like to prioritize CQI/Sounding resources when resource shortage is encountered within the BS.

The most efficient way to support periodic DL CINR reports is over the CQI channel (CQUICH), therefore whenever new MS joins the BS a CQICH is allocated for it. The CQICH has a periodicity of 8 frames, an infinite duration (reports are stopped only when CQICH is de allocated).

MSs that support only matrix A MIMO require one CQICH but MSs with support for matrix B MIMO require two CQICHs. When working with two CQICHs, MS is supported with one main CQICH all the time and with additional second CQICH when needed. There must be a mechanism to provide the second CQICH to the MSs that need it and also that enables the MSs that are left without CQICH, to receive one when it is possible.

Channel sounding is a signaling mechanism where a mobile station (MS) transmits known waveforms on the uplink to enable the base station (BS) to estimate the BS-to-MS channel response. The sounding transmissions are periodic and are sent within a dedicated UL "zone" which is one or more symbol and is shared by several MSs. This is what we call sounding resources.

The channel sounding manager is responsible for allocating and de allocating sounding resources. It decides which MS will be allocated with sounding and how often (periodicity) and it takes care the sounding resources are used properly and not wasted. Without sounding allocations an MS can not be served with beam forming since antenna weighting is based on channel estimation.

### 2.3.2 Proposed solution

#### 2.3.2.1 Prioritization of sounding resources

Our prioritization targets are, first filtering negligible traffic MSs from having sounding resources and second maximizing sector throughput by prioritizing (near) full loading MSs. The periodicity allocated and de-allocates sounding resources based on MS priority. The priority metric for MS will be based on slots consumption of MSs within time interval (number of frames) given by:

$$\textbf{P(MS)= DL Data Slots Interval}$$

DL data slots are the number of slots that were allocated to MS within DL Sub-frame for data transmissions during time interval. Data slots to be counted are total allocated slots including all retransmissions.

The sounding prioritization mechanism is determined to be a slow process in which the system can locate the significant MSs candidate set for sounding allocation/de-allocation. The interval length (number of frames) is being determined by a configuration parameter *bsSoundingPriorityInterval*. Interval length value range is defined as minimum of 250ms up to 5000ms.

At the end of time interval, the MSs without sounding resources having the highest number of allocated slots and have an average effective MCS above a configured threshold *bsSoundingAllocationRateTh* should be replaced with MSs having sounding resources with lowest number of allocated slots as long as the number of slots for non-sounding resource MSs is bigger than the lowest ones for MSs with sounding resources.

When the above condition is not satisfied, we would still like to have allocation/de-allocation resources exchanges in order to compensate of the priority mechanism especially for ET scheduling. In that case the number of sounding resource allocation/de-allocation couples for interval should be limited by a hard coded parameter. Default value for the parameter should be 4.

By supporting MS with sounding and BF, we improve its DL SNR, therefore it can work at higher MCS and use less slots for its allocations. That way, more MSs can be supported in the same bandwidth and system capacity increases.

Choosing the MS with the highest slot consumption allows saving more slots and therefore improves the system capacity more than if other MS would have been chosen.

For Equal Rate scheduling, the MS with the highest number of slots has low DL SNR and may be benefit to get higher DL SNR improvement from having sounding. Since in ER, users with large number of slots are the dominant resource consumers, improving these MSs with contribute to higher aggregate sector throughput. Moreover, low number of consumed slots reflects low traffic users in which we do not want to prioritize for sounding or to be removed from sounding with high priority.

For Equal time scheduling, the number of allocated slots will help to filter the MSs with low traffic, however the drawback will be not prioritizing the MSs with same number of consumed slots however with different SNR improvement potential.

### 2.3.2.2 Prioritization of CQI channels

MS prioritizing is needed for times when CQICH is fully occupied. In such times, there are not enough CQICH allocations for all MSs and there will be

active MSs that are left without CQICH or MSs that do not have the additional CQICH for MIMO support (when MIMO required two CQICH).

The priority mechanism makes sure that those MSs that will work without CQICH or served in a lower quality than they required are the MSs that their influence on the BS throughput because of the degradation is the smallest influence.

MSs which need second CQICH for MIMO support or MSs which work without any CQICH may damage throughput more radically because:

1. Without the additional CQICH, BS will not allow MS to work with MIMO matrix B and to double the data rate
2. MSs that work without CQICH are "doomed" to work at MCS of QPSK 1/2 (or any other basic rate) therefore their user experience may be damaged. MSs without CQICH also require more BW at the expense of BW for other MSs

Based on the described above, the chosen criteria for MSs prioritizing is throughput maximizing/ spectral efficiency maximizing.

The MS that should work without CQICH is the MS that will have the less added BW as a result of working with wrong CQICH level of service. The prioritizing mechanism looks at the additional amount of data due to CQICH degradation and chooses the MS that will add the minimal amount of slot per frame after degradation.

With the amount of allocations "saved" by choosing the right MS to work without CQICH, there is a better and correct use of BW and spectral efficiency is reached

The MS priority reflects the addition amount of allocation per frame needed for the MS. It is based on two factors:

1. Avg bits per slot before degradation and after degradation
2. Avg slots per frame which reflects MS demand for data, QoS, scheduler prioritizing for this MS in compare to other MSs in the system.

Why bits per slot should influence the calculation? When taking the only CQICH from the MS or taking its additional CQICH, its demand for data allocations will remain the same and will not influence the amount of allocations which are given to it (unless scheduler is proportional fair).

In most cases (excluding the scenario where MS already worked at minimal rate) without the one CQICH or the additional CQICH (when needed) the number of bits per slot would decrees and so, MS would need more slots allocation in order to send the same amount of average data it used to send.

Assuming both MSs have the same average amount of data to send and one of them works at higher average rate (more bits per slot) than the other, and assuming CQICH degradation leads to the same new amount of bits per slots

for both, then the MS with the higher rate (bits per slot) would need more allocations per frame after degradation than the MS with the lower amount of bits per slot would need if it suffers degradation. If MS already work at minimal rate (in average) than working without CQICH will not influence it at all therefore its CQICH should be taken from it in case of shortage.

Why Avg slots per frame should influence the calculation? Assume two MSs have the same rate (bits per slot) before and after degradation. One of them needs big number of slots per frame while the other needs few slots per frame. After removing the CQICH to one of them, they both require same Avg amount of allocations, but the MS that has high BW demand will add a lot more additional slots per frame, than the other MSs that has less data demand.

In most cases we compare between MSs with different BW allocations and different bits per slots, as well as different influence of CQICH degradation. Therefore it is only the multiply of both that can provide a suitable answer to the question which MS should work without CQICH at the current state.

It should be notice that priority mechanism is completed with the algorithm of returning back MS without CQICH or with less CQICH than it should have to have the amount of CQICH it needs, according to same criteria that maximize the BS throughput.

### 2.3.2.3 FLAVIA support

FLAVIA's architecture provides us with the possibility to use dynamically different prioritization criteria for CQI/sounding resources. In the next section we describe a simulation in which different parameters were tested in order to discover the best one in a scenario in which CQI/sounding resources were lover than the number of registered MSs. Such a situation should be more deeply analyzed but probably the optimal criteria can be different depending on the circumstances and it would be good that the prioritization mechanism/criteria can be dynamically changed.

### 2.3.3 Some insights on resulting advantage and benefits

Due to budget constraints the prioritization algorithms were not implemented. However we performed a simulation in order to evaluate the goodness of the approach we took under stressing conditions, specifically when the number of CQI resources is lower than the number of registered MSs in a given BS.

The simulation environment is described in Table 1.

| Simulation Environment |
|---|
| DL:UL ratio: 29:18 |
| 30 DL sub-channels |

| |
|---|
| SIMO |
| HARQ |
| 56 CQI resources |
| 70 active MSs – 14 MSs with no CQI resources |
| Interval length – 50-400 frames |
| Up to 4-10 switches at the end of time interval |
| PF – Beta 0.99/ET – Beta 0.001 |
| Packet size 1500 bytes |
| Basic rate QPSK ½ rep. 6 |
| 25% CQI resources oversubscription eligible active full loaded MSs |
| Immediate switching of MSs at the end of interval |
| Upgraded/Downgraded rate is updated immediately |

*Table 1: Simulation environment for CQI/sounding prioritization algorithms*

In the simulation a prioritization metric is being used based on fixed time intervals. Statistics are gathered within time interval and then used as raw data for prioritization decisions at the end of time interval. At the end of such time interval, MSs without sounding resource having metric value higher than MSs with sounding resource are being switched (up to maximum switches per interval). The Prioritization Metrics that we have considered are listed in Table 2.

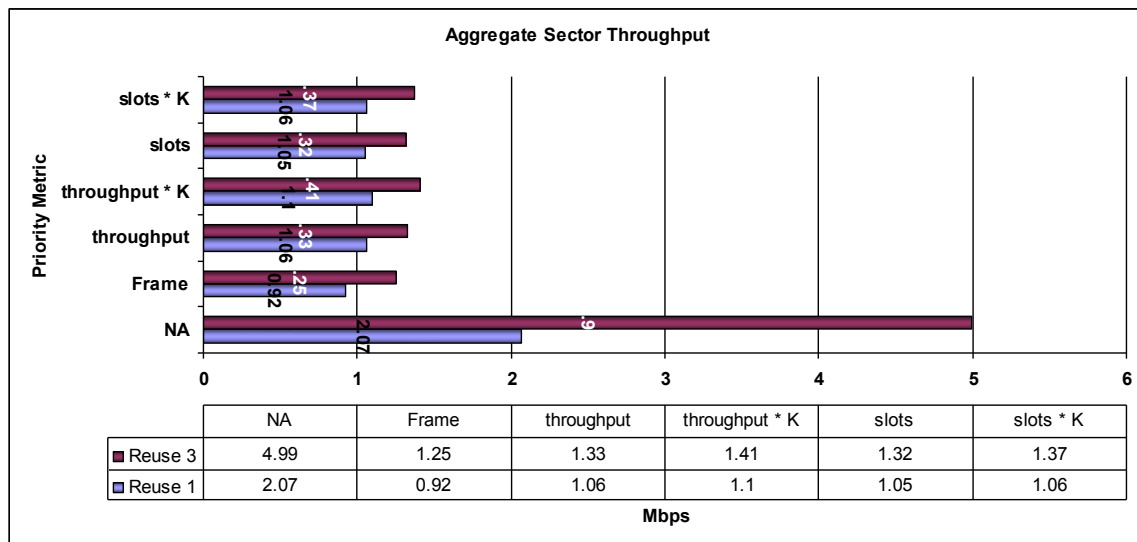| Prioritization Metrics |
|---|
| **Slots** : Total number of slots allocated to MS (including retransmissions) within time interval |
| **Slots * K factor**: Total number of slots allocated to MS (including retransmissions) within time interval multiplied by K factor where <br> K factor = \|MCS no resource - MCS resource\| |
| **Throughput**: Total number of bytes sent to MS within time interval (no retransmissions are considered) |
| **Throughput * K factor**: Total number of bytes sent to MS within time interval (no retransmissions) multiplied by K factor where <br> K factor = \|MCS no resource - MCS resource\| |
| **Frame**: Last scheduled frame number |

*Table 2: Prioritization metrics*

The simulation provided different results for Equal Rate (ER) scheduling strategy and Proportional Fair (PF) scheduling strategy.

In Figure 10 we can see the results obtained for ER. The NA entry refers to the situation in which the number of CQI resources was equal to the number of registered MSs (50).

We can observe that for ER scheduling and for the situation in which we have 50 CQI resources and 74 MSs, the sector throughput is decreased dramatically by around 50% for reuse 1 and 75% for Reuse 3.  This is due to the fact that MSs with no sounding resource (using basic rate) consumes most of BW to achieve their data rate. Then, all priority metrics have limited capability due to the significant impact of basic rate MSs.

**Aggregate Sector Throughput**

|         | NA   | Frame | throughput | throughput * K | slots | slots * K |
|---------|------|-------|------------|----------------|-------|-----------|
| Reuse 3 | 4.99 | 1.25  | 1.33       | 1.41           | 1.32  | 1.37      |
| Reuse 1 | 2.07 | 0.92  | 1.06       | 1.1            | 1.05  | 1.06      |

**Mbps**

*Figure 10: Prioritization performance with Equal Rate scheduling strategy*

ER scheduling suffers tremendous degradation for all loading schemes and priority metrics. Fully loaded MSs with no CQI resource have tremendous impact of aggregate sector throughput (and per MS throughput). Every used metric gain is very limited, since MSs with no CQI resource cannibalize resources. 'Frame' metric is the worst among all priority metrics (20% below highest metric).

For partial loading system, all priority metrics except 'frame' metric screens out the negligible MSs from getting sounding resource and the first to be removed from having one. The effect for wrong decision is harmful as we could observe in the case of 'frame' metric.

Finally, we observed in the simulation that longer time intervals increase performance especially for fully loaded scheme. For instance time interval around 400 frames increase ER performance by additional 10%.



**Aggregate Sector Throughput**

|  | NA | Frame | throughput | throughput * K | slots | slots * K |
|---|---|---|---|---|---|---|
| Reuse 3 | 7.35 | 6.24 | 5.98 | 6.35 | 6.07 | 6.82 |
| Reuse 1 | 4.26 | 3.75 | 3.51 | 3.98 | 3.61 | 3.92 |

**Mbps**

*Figure 11: Prioritization performance with Proportional Fair Rate scheduling strategy*

The results obtained in the simulation for PF scheduling can be seen in Figure 11.

For Reuse 1 the 'throughput * K' metric achieves best performance (7% below NA scheme) while as for Reuse 3 the 'slots * K' metric achieves highest gain (7% below NA scheme). It is interesting to observe that the PF is fed by changes of MCS levels of MSs. The 'throughput/slots * K' metric ensures maximum changes of MCS levels of MSs. We can observe that using no K factor reduces performance by additional 10%  (metric 'throughput').

For both reuse schemes 'throughput' metric is the worst since no switching takes place. Then, there is a need to switch out of priority users. All metrics except 'frame' metric filter the negligible MSs of having sounding resource and first to release their sounding resource.

Interval length has impact on PF performance. For PF, the more the interval is short, the better opportunism for PF.

For partial loaded system all metrics introduce same performance. Using K factor within system is not trivial, however it can boost performance by up to 10% on fully loaded system under shortage

Finally the last observation is that ET and PF metrics gain results are correlated

## 2.4 Distributed Opportunistic Scheduling in Non-Homogeneous Networks (BGU)

### 2.4.1 Problem statement and motivation

The desire to cope with the growing demand for wireless access networks which are the last hop connecting users wirelessly to the high speed backbone network, stimulates new innovations in coding schemes, resource allocation, transmission techniques, etc. One such emerging technology in the recent decade is multi-user schemes which exploit the spatial diversity of the wireless channel. Nonetheless, even though analytical schemes such as multiple access channel codes and Dirty Paper Coding (DPC) for Gaussian broadcast channel have been shown to dramatically boost network capacity, e.g., [2], in practice these techniques are not used as they involve complex computations, as well as transmit to too many users simultaneously. Accordingly, most worldwide deployments still use traditional mechanisms, via either scheduled or random access protocols, which ensure that only a single user (or a small group of users) is active at any given time, forfeiting the huge capacity gain potential due to channel diversity.

Although the current literature does include numerous results on scheduling and capacity for wireless networks for the case of homogeneous users, i.e., for scenarios in which one user is selected in each slot and the channel experienced by the users is independent and identically distributed (i.i.d.) (e.g., [3, 4, 5, 6, 7]), our current understanding of how to successfully schedule users (or groups of users) in various non-homogeneous, spatially correlated or slow fading models, and the ensuing system performance is remarkably limited. Yet, it is clear that these models, albeit complicated, are necessary to address the diverse conditions in current and future networks.

In this study, we design novel distributed channel access algorithms for Multiple Input Multiple Output (MIMO) systems and analyze their capacity under diverse, non-homogeneous models. In particular, our contribution is twofold.

First, we suggest a novel technique, based on the Point Process approximation, to analyze the expected capacity of scheduled multi-user MIMO systems. We show that this approximation not only allows us to derive results which analyze the asymptotic (in the number of users) capacity for various threshold based algorithms (i.e., algorithms in which at each time slot only users with anticipated capacity greater than a given threshold transmit) in the case of homogeneous users, but also enable us to extend these existing results to analyze the capacity of similar systems for the case of non-homogeneous users, i.e., for the case in which the channel experience of different users is not identically distributed. We further broaden the results for the heterogeneous users, from cases in which users are inherently non-uniform, e.g., due to spatial location differences, to cases in which the users are deliberately not identical, e.g., due to Quality of Service (QoS) constrains or when fairness considerations are added. To date, these scenarios have not been subjected to rigorous analysis.

Second, we leverage our understanding of the system in non-homogeneous channel conditions to design and analyze a novel distributed algorithm, which achieves a constant factor of the maximal multi-user diversity without centralized processing or communication among the users. Our algorithm is based on the Mini Slotted Alternating Priority (MSAP) protocol [8], in which the data transmission period is preceded by a short reservation interval for which the user transmitting in the following data interval is selected. The suggested algorithm is completely distributed, and does not require any information collection as to which users are backlogged or what is the channel condition of each backlogged user. Our reservation mechanism is as previously a threshold based algorithm, in which at each mini time slot only the users with capacity greater than a given threshold transmit in the following data interval. The near optimal capacity is achieved by choosing thresholds for each mini slot in a way that minimizes the unutilized air time of the system, i.e., minimizes the number of unutilized data intervals. We validate all the analytical results suggested via thorough simulation results.

### 2.4.2  Preliminaries

#### 2.4.2.1 MIMO Capacity
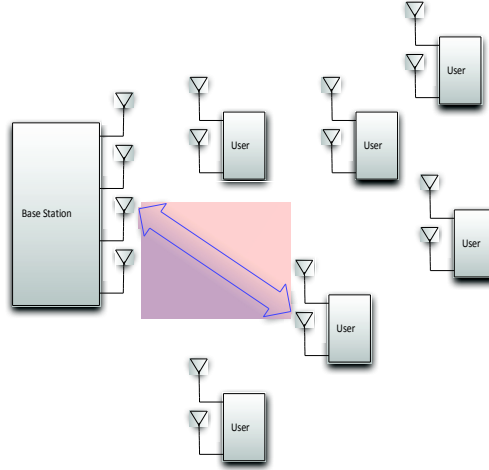
We consider the uplink traffic in a single-cell multiple-access model, comprising a single base station (the receiver) which is equipped with $r$ receive antennas, and $K$ users (the transmitters), each equipped with $t$ transmit antennas (**Figure 5**). We assume that time is slotted, whereas throughout most of this paper we assume that only one user can utilize the channel in each time slot,

i.e., single user MIMO (SU-MIMO). As previously mentioned, the focus of this study is on selecting this user distributively, and computing the resulting capacity. We extend our results to cases where more than one user can utilize a single slot.



**Figure 12 MIMO system with K users**

Denote by $H_i^s \in \mathbb{C}^{r \times t}$ the channel gain matrix between user $i$ and the base station at time-slot $s$. We assume $H_i \in \mathbb{C}^{r \times t}$ is a complex random Gaussian channel matrix. Note that when all users are identically and independently distributed, it is common to assume the entries of $H_i$ have mean zero and variance $1/2$ imaginary and real parts for all users. However, in the major part of this paper, the focus will be on cases where the channel matrices $\{H_i\}_{i=1}^K$ are **not identically distributed**. Specifically, we will assume different parameter values for the channel distribution of each user. We assume that the channel is memoryless, that is, for each channel use (slot), independent realizations of $\{H_i\}$ are drawn. Accordingly, for ease of notation, throughout this study we discard the dependence on $s$, as well as the dependence on $i$ when clear from the context.

The signal received by the base station when only user $i$ is active, denoted by $y_i \in \mathbb{C}^r$, is given by:

$$y_i = H_i x_i + n$$

where $x_i \in \mathbb{C}^t$ is the transmitted vector. $x_i$ is constrained in its total power to $P$, i.e., $E[x^\dagger x] \leq P$. $n \in \mathbb{C}^r$ denotes the uncorrelated Gaussian noise.

In this study, we are interested in the capacity under a scheduling scheme. Obviously, the capacity of such a system is highly dependent on the scheduling scheme, i.e., the choice of which user transmits at each time slot. For example, the base station can schedule the users according to some equal

time share schedule (e.g., Round-Robin). The expected capacity of such a scheduling scheme when all users have the same channel statistics is the same as if there is only a single user, that is $E\left[\log\det\left(I_r + \frac{P}{t}HH^\dagger\right)\right]$, [9]. In several cases, it is beneficial to represent the capacity in terms of the eigenvalues $\lambda_1, \dots, \lambda_{min\,(r,t)}$ of the Wishart matrix $H^\dagger H$, i.e.,

$$C = E\left[\sum_{i=1}^{\min(r,t)} \log\left(1 + \frac{P}{t}\lambda_i\right)\right]$$

When users are not homogeneous and different users have different channel statistics, the expected capacity is amended according to the fraction of time each user receives and its associated expected capacity.

If the goal is to maximize the capacity, it is beneficial to schedule the **strongest user** in each time-slot (**Figure 6**). The expected capacity will then be ( [9]):

$$C_{max} = E\left[\max_{i=1,\dots,K} \log\det\left(I_r + \frac{P}{t}\,H_i H_i^\dagger\right)\right] \qquad \textbf{\textit{2.1}}$$

This is the essence of multi-user diversity.



**Figure 13 Achieving the upper envelop by scheduling the strongest among four users. The four graphs depict the capacities of the four users at each time-slot, while the circled points depict the chosen user in each slot.**

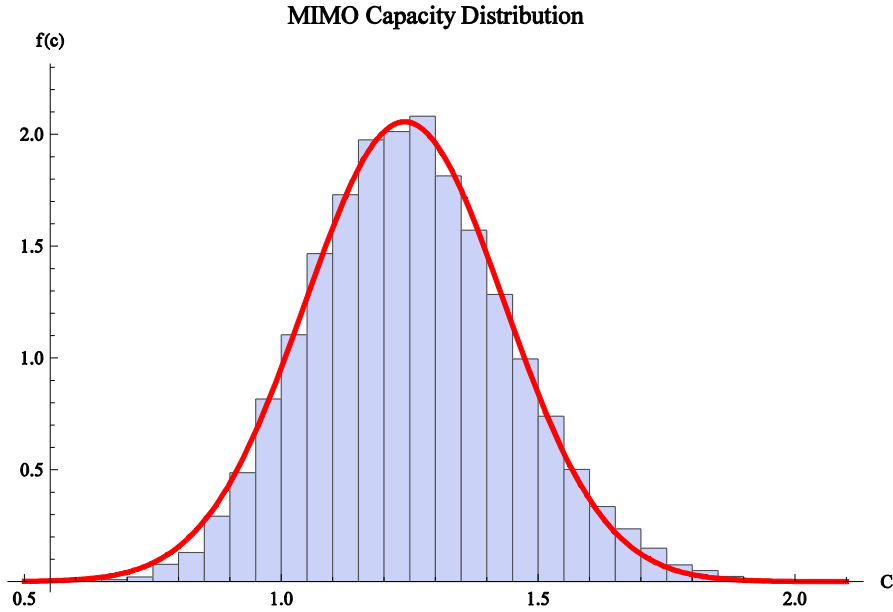The majority of this work is devoted to calculating the capacity for a scheduling scheme which allocates each time slot to the user (users) with the strongest received signal at the base station, i.e., the expected capacity as presented in 2.1, first in the case that user channel statistics is i.i.d. and more importantly in the more complicated cases where channel matrices are not identically distributed, or for scheduling schemes in which the time is allocated to users according to various quality of service and fairness constraints. Moreover, our focus will be on gaining the multi-user diversity and achieving 2.1 in a distributed manner, that is, without central scheduling and without communication among the users. Throughout this study, we assume that the exact values of the channel matrix $H$ are not available to the transmitter. Yet, the transmitter can approximate its channel capacity by estimating the SNR of a pilot signal sent from the receiver.

### 2.4.2.2 Multi-User Diversity In the i.i.d Case

In [10, 11, 12], it was shown that when the elements of the channel gain matrix, $H$, are i.i.d. with zero mean and finite moments up to order $4 + \delta$, for some $\delta > 0$, the distribution of the single-user-MIMO capacity follows the Gaussian distribution with mean that grows linearly with $\min(r, t)$, and variance which is mainly influenced by the power constraint $P$. An illustrative example of this result is presented in Figure 7.

***Figure 14: MIMO capacity distribution.*** $30,000$ ***channel matrices were drawn, with*** $min(r,t) = 4$ ***transmitting antennas. The solid line depicts the analytical results of*** ***[11], with*** $\mu = 1.25$ ***and*** $\backslash sigm = 0.19$***, while the bars represent the simulation. A very good fit is clear even in this non-asymptotic case, with a reasonable number of antennas.***

A common tool used to approximate the distribution of the maximum in 2.1, for large number of i.i.d. users, is **Extreme Value Theory (EVT)** (e.g., [13, 14, 15]**)**. We first briefly review key EVT results that are used in the asymptotic analysis of the capacity gain.

As mentioned, EVT is at the heart of the capacity gain analysis. Let $x_1, x_2, \dots, x_n$ be a sequence of random variables and let $M_n = \max(x_1, x_2, \dots, x_n)$. In our case, the variables $\{x_i\}$ represent the capacities experienced by the users, that is, $\log \det \left( I_r + \frac{P}{t} HH^\dagger \right)$, while $M_n$ represents the capacity experienced by the maximal user (with $n = K$). If there exists a sequence of normalizing constants $a_n > 0$ and $b_n$ such that as $n \to \infty$, $\Pr(M_n \le a_n x + b_n) \xrightarrow{i.d.} G(x)$ for some non-degenerate distribution $G$, then $G$ is of the **Generalized Extreme Value** (**GEV**) distribution,

$G(x) = e^{-(1+\xi x)^{-\frac{1}{\xi}}}$ where $\xi$ is a **shape parameter**. In particular, if the random variables are normal $N(\mu, \sigma^2)$ and *i.i.d.*, the asymptotic distribution of $M_n$ is a Gumbel distribution. Formally,

**Theorem 1** $[14]$ Theorem 1.5.3

If $\{x_n\}$ is an *i.i.d.* standard normal sequence of random variables, then the asymptotic distribution of $M_n = \max(x_1, x_2, \dots, x_n)$ is a Gumbel distribution. Specifically,

$$\Pr(M_n \leq a_n x + b\_n) \to e^{-e^{-x}}$$

*2.2*

where

$$a_n = (2\log n)^{-\frac{1}{2}}$$

*2.3*

and

$$b_n = (2\log n)^{\frac{1}{2}} - \frac{1}{2}(2\log n)^{-\frac{1}{2}}(\log\log n + \log 4\pi)$$

*2.4*

As a result, if $\{x_n\}$ follows the Gaussian distribution with mean $\mu$ and variance $\sigma^2$, then the above normalizing constants are:

$$a_n = \sigma(2\log n)^{-\frac{1}{2}}$$

*2.5*

and

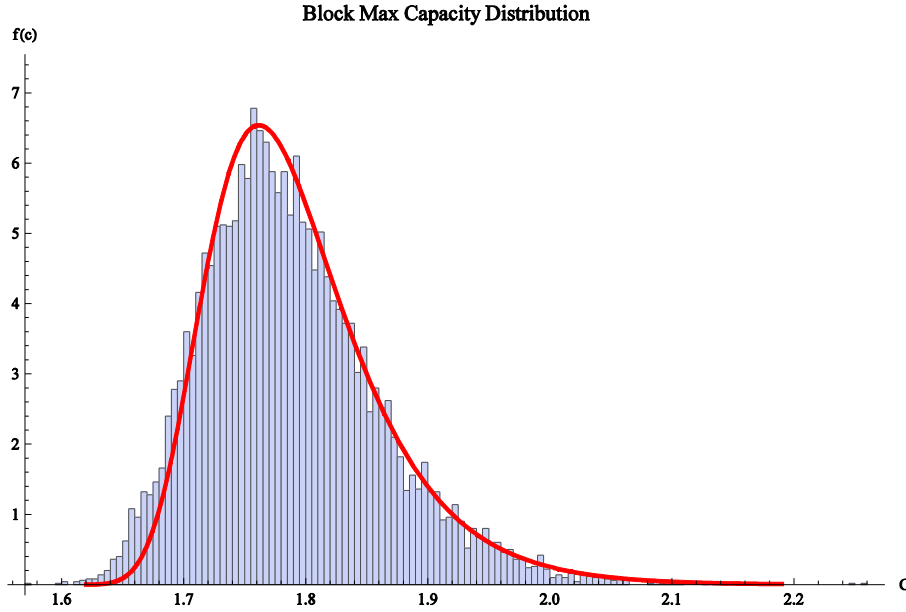$$b_n = \sigma\left[(2\log n)^{\frac{1}{2}} - \frac{1}{2}(2\log n)^{-\frac{1}{2}}[\log\log n + \log 4\pi]\right] + \mu$$

*2.6*

It follows that for a Gaussian distribution,

$$a_n = \sigma(2\log n)^{-\frac{1}{2}} \to 0$$

*2.7*

which implies that

$$M_n \sim b_n \sim \sigma(2\log n)^{\frac{1}{2}} + \mu$$

*2.8*

Some simulation results which illustrate 2.8 are depicted in Figure 8.

***Figure 15: Maximal MIMO capacity among*** $300$ ***users, each following the Gaussian distribution with*** $\mu = 1.25$ ***and*** $\sigma = 0.19$***. The solid line represents the analytical results, corresponding to the Gumbel density function in the range*** $[\mu + 2\sigma, \mu + 5]$ ***while the bars represent simulation results***

From the above discussion, it follows that assuming a MIMO uplink model with $K$ users, with perfect CSI at the receiver, the asymptotic expected capacity (at the limit where the Gaussian approximation is applicable and for a large number of users) achieved by scheduling the user with the maximum capacity in each time slot, follows the expected value of a Gumbel distribution with parameters $a_K$ and $b_K$ (e.g., [16]). That is

$$\mathrm{E}[\boldsymbol{M_K}] \quad =^{(a)} \sigma(\mathrm{b_K} + \mathrm{a_K}\gamma) + \mu$$

$$=^{(b)} \sigma\left[(2\log K)^{\frac{1}{2}} - \frac{1}{2}(2\log K)^{-\frac{1}{2}}[\log\log K + \log 4\pi]\right. \qquad \textbf{\textit{2.9}}$$
$$\left. + \gamma(2\log K)^{-\frac{1}{2}}\right] + \mu$$

where $\gamma \approx 0.57721$ is Euler-Mascheroni constant, $(a)$ follows from the expectation of the Gumbel distribution and $(b)$ follows from 2.3 and 2.4. Hence, for large enough $K$,

$$\mathrm{E}[\mathbf{M}_K] = \sigma(2\log K)^{\frac{1}{2}} + \mu + \mathrm{o}\left(\frac{1}{\sqrt{\log K}}\right)$$

This gives rise to the $\sqrt{2\log K}$ factor of the expected capacity in a scheduled MU-MIMO system with $K$ users and perfect CSI at the base station. However, the analysis above applies only to perfect CSI at the base station and *i.i.d.* users. The goal of this work is to extend it to the cases where the scheduling is distributed, that is, there is no central decision making at the base station, and where the users are not necessarily identically distributed.

### 2.4.2.3 Threshold Based Algorithms and Point Process Approximations

We are now ready to present the tools, namely, Point Process Approximations (PPA), that will allow us to suggest distributed algorithms and analyze their performance, both in the *i.i.d.* model which was solved using the EVT discussed above, and in the new models we suggest in this study, which are currently unsolved.

In order to describe the PPA let us examine the following algorithm which was utilized by some of the opportunistic schemes discussed above.

**Algorithm 1:** *Set a capacity threshold such that a fraction of the users will exceed it. If the capacity seen by a user is greater than the capacity threshold, it transmits. Otherwise, it remains silent.*

At first, we assume the base station can successfully receive the transmission if no collision occurs, i.e., exactly one user exceeds the threshold. Later on, we will relax this assumption.

Let $C_{av}(u_k)$ denote the expected capacity, given a threshold $u_k$ such that $k \ll K$ users exceed it on average. The threshold $u_k$ influence on the average capacity is twofold. First, it determines the average number of users that will attempt transmission. Second, it impacts the expected capacity seen by a user, **given the user exceeded it**. Accordingly, we classify the slots into three types: *(i)* **Idle slots** in which no user attempts to transmit; *(ii)* **Collision slots** in which more than one user attempts to transmit; *(iii)* **Utilized slots** in which exactly one user transmits.

Since only utilized slots contribute to the capacity, the overall expected capacity $C_{av}(u_k)$ has the form:

$$C_{av}(u_k) = \Pr(\text{utilized slot}) \cdot \mathrm{E}[C|C > u_k]$$

*2.10*

The power of the PPA we now discuss, is in the ability to analyze 2.10 for several complicated scenarios, such as dependent and non-identical user distributions. This analysis was considered intractable with previously used

methods. Detailed explanation of PPA can be found in many textbooks such as [17, 18, 19].

### 2.4.3  Summary of the results

In the sequel we present a brief summary of our results.

#### 2.4.3.1 Establishing The Expected Capacity in the *i.i.d.* Case

In the *i.i.d.* case, the PPA can be easily harnessed to reproduce the results previously established using the extreme value distribution directly. Assume a threshold $u_k$ is set such that the expected number of users who pass it is $k$. An evaluation of the expected capacity 2.10 should be done in three steps. First, the probability of a utilized slot should be evaluated using the PPA. Then, the conditional distribution of its capacity should be computed. Finally, the result can be optimized over $k$.

Accordingly, we used the PPA to analyze a distributed, threshold based algorithm to schedule ***i.i.d. users in a MIMO uplink system***. Let $C_{av}(u_k)$ denote the expected capacity, given a threshold $u_k$ such that $k \ll K$ *i.i.d.* users exceed it on average. For sufficiently large number of users, $K$, we showed that the following proposition holds.

**Proposition 1:**
The expected capacity when working with a single user in each slot is:
$$C_{av}(u_k) = \mathrm{k} \cdot e^{-k}(\mathrm{u}_k + \sigma \cdot \mathrm{a}_K) + \mathrm{o}(\mathrm{a}_K) \qquad \textbf{\textit{2.11}}$$
where $a_K$ is the normalizing constant of the extreme value distribution.

For example, ***Figure 9*** depicts the capacity distribution of users which exceeded the threshold $u_k$. ***Figure 10*** implies that the system will be idle $e^{-1}$ of the time when setting the optimal threshold, such that a single user exceeds on average ( [20] Proposition 4).

Capacity Distribution given C> threshold



***Figure 16: Tail distribution of the capacity for a high threshold. Bars represent simulation results while solid lines are from direct tail analysis and EVT-based analysis. The statistics are for 11,722 observations out of 50,000,000 that exceeded a threshold of*** $3.5$***, which is*** $\approx 1 - \Phi(3.5)$ ***of the observations.***

Expected Capacity when k users exceed on average



***Figure 17: Threshold estimation. Expected capacity gain for*** $K = 1000$ ***users, when the threshold is such that*** $k$ ***users exceed on average (dashed line), compared to the expected capacity of the optimal multi-user diversity centralized scheme (dot-dashed line)***

### 2.4.3.2 Non-Homogeneous Users

As mentioned, previous work focused on algorithms and capacity evaluation when users are assumed *i.i.d.* In this study we relax the homogeneity

assumption. That is, we assume that users, due to, for example, different distances from the base station, experience different channel statistics.

To do this, we will aim to obtain an approximating Poisson process for non identically distributed random variables. Such an approximation can be achieved by considering the threshold arrival rate of **a single user, with respect to $M$ consecutive time slots**. Since, herein, we still assume independence and identical distribution **in time**, this arrival rate can be approximated by a Poisson process, but this time matched to the capacity distribution of the specific user. E.g., we compute an arrival rate $\Lambda_i$ for user $i$, and analyze the capacity of the distributed algorithm based on multiple, independent Poisson processes.

Specifically, we considered the case where the *i*-th user capacity follows a Gaussian distribution with mean $\mu_i$ and variance $\sigma_i^2$. Again, we used the PPA to prove the following.

### Theorem 1

The expected capacity for non-homogeneous users, where the $i$-th user capacity is approximated with mean $\mu_i$ and variance $\sigma_i^2$, is:

$$C_{av}^{nu}(u) = \frac{1}{K}\Lambda_T e^{-\frac{1}{K}\Lambda_T} \sum_{i=1}^{K} \frac{\Lambda_i}{\Lambda_T}\big(u + \sigma_i a_K + o(a_K)\big) \qquad \textbf{2.12}$$

where $\Lambda_i = e^{-\frac{u-(\sigma_i b_K + \mu_i)}{\sigma_i a_K}}$ is the average threshold exceedance rate of the *i*-th user, $\Lambda_T = \sum_{i=1}^{K} \Lambda_i$.

$u$ is a threshold greater than zero that is set for all users, and can be optimized to maximize $C_{av}^{nu}(u)$.

*Figure 11* compares the analytical results of Theorem 1 to the simulations. Indeed, the PPA results in very accurate approximation.

*Figure 18: Expected capacity in a non-homogeneous environment. Bars represent simulation results, while solid lines represent our analytic results. The blue (middle) lobe depicts the expected capacity for $K = 1000$ users in a non-uniform environment, where the channel capacity of each user follows the Gaussian distribution with $\sigma_i \sim U[0.03, 3]$ and $\mu_i \sim U[\sqrt{2}-1, \sqrt{2}+1]$. The red lobe (right) depicts the capacity when all users have the same capacity as the strongest user. The black lobe (left) depicts the capacity when all users have the same capacity as the mean user.*

### 2.4.3.3 Capture effect

In our previous models we have assumed that simultaneous transmissions of two or more users result in collision and unutilized slot (packet loss). Of course, such collisions can be avoided when coordinating the users, or, in the more general scenario, when using multiple access channel coding and decoding. As previously mentioned coordinated system do not scale to large number of users, and multiple access coding techniques are often neglected due to the required complexity or coordination as well.

Nonetheless, previous studies have shown that based on what is known in literature as the *physical capture effect*, in many cases a receiver can correctly receive the strongest of multiple transmitted signals ( [21, 22, 23]). Obviously, in such cases capacity is expected to grow as some of the unutilized slots will become utilized. This is done without the overhead of coordination or complex coding strategies.

Finally we analyzed the capacity of a system with non-homogeneous users in the presence of the capture effect. We assume throughout that whenever two transmissions are transmitted simultaneously, the strong one is received

correctly and whenever three or more transmissions are transmitted simultaneously all of them fail, regardless of the strength of the strong signal or the strength of the weaker signals. Once again, the PPA facilitates the analysis, and the resulting capacities are significantly higher than the single user case. Moreover, from **Figure 12** it is clear that the optimal threshold is lower, resulting in more users passing it.



*Figure 19: Capture effect capacity gain for 1000 i.i.d. users. Solid line represents the expected capacity when setting a threshold such that $k$ users exceed the threshold on average. Dashed line represents the expected capacity when $k$ users that are subject to the capture effect, exceed the threshold on average. The dot-dashed line represents the expected capacity of the optimal multi-user diversity centralized scheme. Note the difference in the maximizing values of $k$*

### 2.4.4  FLAVIA support

The discussed above distributed opportunistic scheduling scheme takes FLAVIA a step forward in terms of flexibility and convergence between scheduled based and random access systems.

Specifically, with FLAVIA's services and interfaces, we can implement a system, where the downlink side is fully controlled and scheduled by the base station, and, yet, the uplink side supports (in part of the band or over the whole band) distributed opportunistic scheduling where the mobiles decides whether to access the channel. Such an unusual system is feasible by modification of merely the Scheduling Strategy and the MAC scheduler services at both the base station and mobile station sides. At the mobile station side, it further requires some modification in the Link Adaptation, which provides the input to the MS scheduler whether to access the channel or not (based on the channel conditions).

## 2.5 Joint scheduling and power control under noise rise constraints (BGU)

### 2.5.1 Problem statement and motivation

The desire to provide integrated broadband services while maintaining Quality of Service (QoS) guarantees bestows growing interest in scheduled access techniques used in multiple-access protocols for future broadband radio systems. Such schedule-based techniques are utilized to ensure that a transmission, whenever made, is not hindered by any other transmission and is therefore successful. Accordingly Orthogonal Frequency-Division Multiple Access (OFDMA) has been widely adopted as the core technology for various broadband wireless data systems, including the next generation cellular systems, 3GPP Long Term Evolution (LTE), \cite{sesia2011lte}, and IEEE 802.16e/m (WiMAX), \cite{IEEE_802.16e,IEEE_802.16m}. In these systems, the BS allocates (schedules) distinct frequency-time chunks among the active MS within its cell, both for their downstream (BS to MS) and for their uplink (MS to BS) traffic. In addition to the frequency-time allocation, the BS also determines the uplink transmission power of the preselected (scheduled) MS, a.k.a. uplink power-control.

Common power-control approaches to the uplink resource allocation problem are either to assign transmission power to the MSs such that all are received at the BS with the same Signal to Interference-plus-Noise Ratio (SINR), or to allow MSs to transmit at their maximal available power. Both of these techniques optimize the MS cell throughput (intra-cell throughput), neglecting the interference injected to neighboring cells (inter-cell interference). However, since OFDMA systems are sensitive to inter-cell interference, the interference from neighboring cells can dramatically decrease the SINR received at the BS, hence reduce the MS throughput. Moreover, without knowing in advance the interference a BS is expected to experience in a transmission, an MS is unable to fine-tune its modulation and coding scheme to the expected SINR at the receiving BS. Accordingly, power control plays a decisive role in providing the desired SINR, not only by controlling the MS received signal strength at its intended BS, but also by controlling the interference caused to neighboring cells. This double role is challenging, as on the one hand as far as intra-cell throughput is concerned, an MS in the proximity of the BS is expected to have high quality link, hence high throughput, even when transmitting in low power, while a distant MS needs to transmit at much higher power to attain the same throughput, and on the other hand as far as inter-cell interference is concerned, MSs near a BS can transmit at high power since they are not in the proximity of other cells, while distant MSs which can be in the proximity of

other cells should not transmit at high power as they can interfere with other (neighboring) BSs.

### 2.5.2  Proposed solution

We suggest a new approach for joint uplink power control and scheduling, which not only ensures high per-cell throughput but also guarantees proportional fairness while alleviating the inter-cell interference. The new approach controls the inter-cell interference, yet does not require any cross-deployment communication or coordination. Accordingly, the aggregate uplink inter-cell interference that all MSs in a cell are allowed to create is bounded. This limited egress interference budget, termed Noise Rise, is treated as an additional limited resource which is allocated to MSs by the BS in conjunction with the ordinary resources (time and frequency), according to some fairness criterion and channel condition. We show that controlling the interference generated by each cell also controls the average interference level sensed by each BS and provides a more predictable uplink SINR, which allows lower interference margins and more efficient rate selection. Hence, it obtains higher capacity and better coverage.

In more details, we based our approach on two main observations. First, assuming a fully homogeneous deployment (the distribution of MSs in the cell is identical at all cells), we have that on the average the ingress interference that neighboring cells inject to a BS equals the egress interference caused by the BS to the neighboring cells. Second, we show that a BS can estimate the normalized interference an MS ejects to surrounding cells from the MS' downlink channel state reports (downlink Signal-to-Interference-Ratio) without the need for inter-cell communication or coordination.

With that at hand, the scheduling problem is formulized under the noise rise constraint as a convex constrained optimization problem, and we provide an efficient iterative algorithm that is proved to solve it optimally. Moreover, we suggest a second setting, in which instead of bounding the average noise rise over all channels, allowing some sub-channels to contribute more noise rise at the expense of further limiting the noise rise on others, we bound the noise rise on each sub-channel to the exact same value. The latter setting allows the decoupling of the scheduling algorithm from the power control and thus facilitates an even simpler algorithm.

### 2.5.2.1 FLVIA support

Both scheduling and power control schemes requires some MAC flexibilities, mainly at the Scheduling Strategy and the Link Adaptation schemes. Specifically, both schemes requires an additional data to perform uplink power
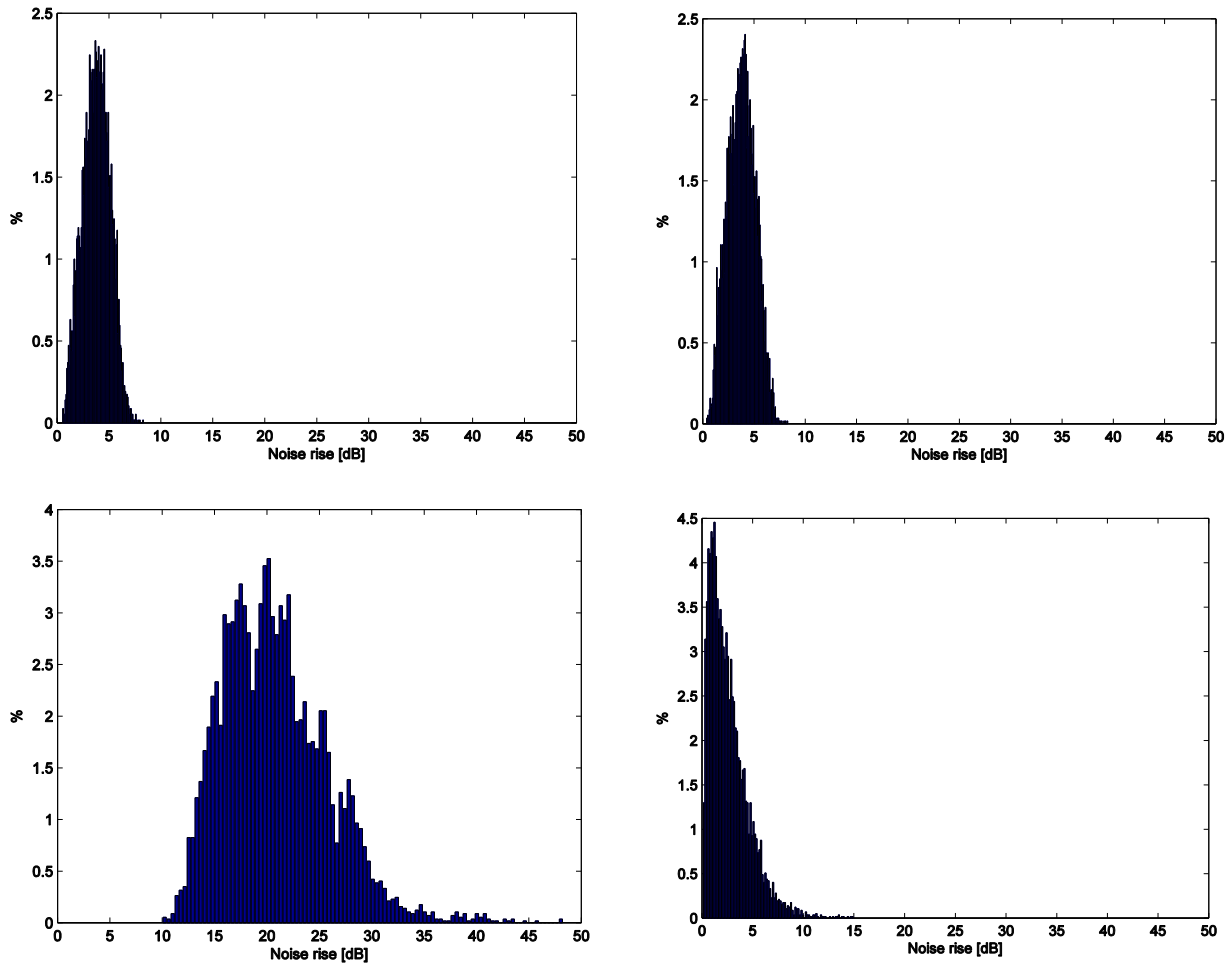
control, namely, estimation of downlink Signal-to-Interference-Ratio. This estimation (performed at the mobile station) can be made available at the base station Link Adaptation and Scheduling Strategy by FLAVIA's the inter-entity interface (IAP2) between the Link Adaptation services (at the MS and BS), and then by the MAC Services function interface (IAP3) between the Proportional Fairness function (of the scheduling Strategy service) and the Power control function (of the Link Adaptation service).

The implementation of the first scheme is more challenging as it requires a closer interaction between the Scheduling Strategy service and the Link Adaptation service, where part of the power control function is performed at the Scheduling Strategy and feed-back (via IAP3) to the Rate Adaptation function of the Link Adaptation service.

### 2.5.3 Some insights on resulting advantage and benefits

The noise rise concept has been thoroughly evaluated via an extensive set of simulations, using both an all-inclusive simulator as defined by IMT-Advanced and numerical results for the exact expressions we analyze utilizing the Shannon-capacity based approach. Our numerical results clearly depict that the suggested approach dramatically increases the overall throughput achieved in each cell compared to the traditional approach, while maintaining fairness. The results obtained by the IMT-Advanced simulator include a more realistic setup which takes into account modulation, coding and several other practical aspects, and show that even though MSs in the proximity of the BS (hence can take advantage of transmitting in high power and high modulation rates), lose throughput due to the noise rise constraint, MSs further away from the BS, and in particular those closer to the cell edge, gain dramatically due to the noise rise constraint.

First, in Figure 13, we compare the ingress noise rise level at the BSs of the two suggested schemes with two traditional power control schemes, namely, (i) maximal transmission power, and (ii) target received SINR. One can see that both the constrained noise rise scheme (termed N.R.) and the relaxed constrained noise rise density scheme based (termed N.R. density) obtain a relatively narrow histogram around the target noise rise (of 4 [db]). Alternatively, the maximal transmission power scheme (termed MaxP) and the target received SINR based scheme (termed SINR) result in a much wider histogram (especially the max power scheme). Such wide histograms corresponds to an unpredictably highly variant uplink interference.

*Figure 20: Histograms of the ingress Noise Rise in [dB] with the various power control schemes: (a) N.R., (b) N.R. density,(b) MaxP, (c) SINR*

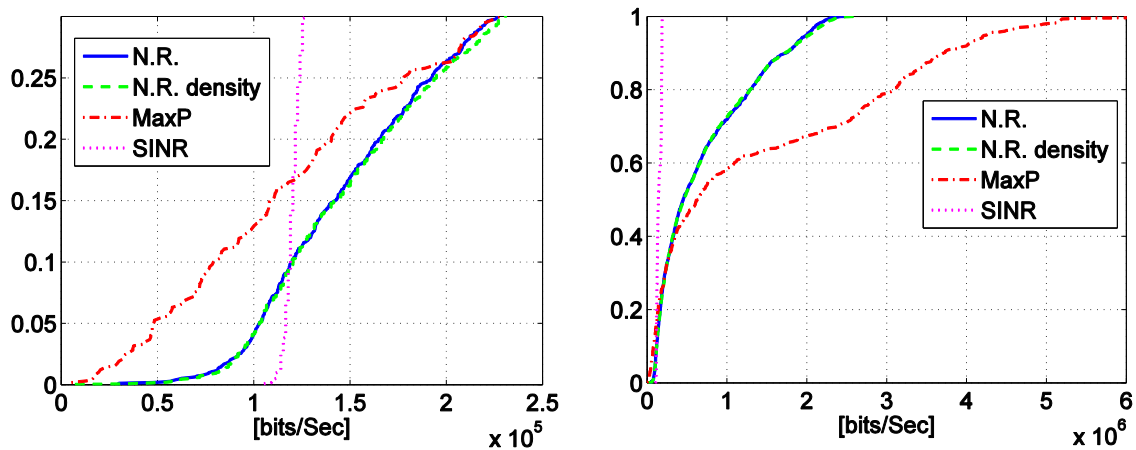To compare the performance of the four schemes we adopt two commonly used indicators, namely, the cell throughput and the cell edge MS throughput (both in [bits/Sec]). The cell edge MS throughput is defined as the 5th percentile point of the *CDF* (*Cumulative Distribution Function*) of MS



throughput.

***Figure 21: Uplink cell throughput and cell edge MS throughput for various values of the proportional fair decay factorβ: (a) cell throughput, (b) cell edge throughput***

Figure 14 depicts the uplink cell throughput as well as the cell edge MS throughput for various values of the proportional fair decay factor $\beta$. Indeed, the greedy approach of transmitting at maximal power obtains the highest throughput. However, it is at the expense of starvation of the cell edge MSs. Alternatively, our noise rise schemes provide a better trade-off between capacity and fairness. For example, the noise rise approach (N.R. in the figure) obtains cell throughput lower by 40% from the cell throughput with the maximal power approach, yet it is 440% higher than the throughput with the SINR based approach. Additionally, cell edge MSs with the noise rise approach gain 71% more throughput than with the maximal power scheme and only 35% less than the more fair SINR based approach. The schemes' fairness is further illustrated in Figure 15 that depicts the CDF of the MSs' throughput and a zoom in on cell edge users. Here, it is clear that the SINR based approach provides the best fairness, where all MSs get similar throughput. Clearly, fairness comes at the expense of total system throughput. On the other hand, the maximal power approach sacrifices about 25% of the MSs by allocating them unacceptably low throughput (which in practice would result in high blocking probability). Again, the noise rise approach provides a good trade-off between cell throughput and fairness.

***Figure 22: The CDF of the MSs' throughput, comparing the fairness of the various power control schemes (a) CDF, (b) zoom-in***

It is interesting to see from Figure 14 that at times the constrained noise rise density algorithm obtains better throughput than the noise rise scheme (e.g., for $\beta = 0.6$). This is due to the difference between the simulated IMT advanced EESM channel model and the Gaussian Channel model, which is fundamental to the noise rise approach. Alternatively, the N.R. density approach decouples the scheduling and power control schemes, allowing a link adaptation that does not assume the Gaussian Channel model.

More details are available in [24].

## 2.6 Inter-cell interference management via base station scheduling (IMDEA)

### 2.6.1 Problem statement and motivation

The continuously increasing demand for higher data rates results in increasing network density, so that inter-cell interference is becoming the most serious obstacle towards spectral efficiency. Considering that radio resources are limited and expensive, new techniques are required for the next generation of cellular networks, to enable a more efficient way to allocate and use radio resources. In this contest, and with the support of FLAVIA's architecture for scheduled systems, we have targeted the design of a frequency reuse 1 scheme, which exploits the coordination between base stations as a tool to mitigate inter-cell interference by separating the scheduling from resource allocation. While common approaches proposed in the literature focus on the optimal user scheduling, we tackle the problem from a different angle. In particular, we formulate a base station scheduling problem to decide whether a base station is allowed to transmit to any of its users in a given sub-frame, without causing excessive interference to any of the users of other scheduled base stations. Finding the optimal base station scheduling is NP-hard, so we formulate a heuristic algorithm to approximate the optimal solution with acceptable complexity cost. By means of numerical and packet-level simulations, we prove the effectiveness and reliability of the proposed solution as compared to the state of the art of inter-cell interference mitigation schemes.

We consider a multicellular LTE-like environment with $N$ base stations and $U$ mobile users. We address only downlink transmissions, for which no power control is adopted, as in the majority of state of the art proposals. Each base station schedules its users across sub-frames, as specified in LTE systems [13], each sub-frame lasting $1ms$. Users associate to the base station from which they receive the strongest signal, and transmission rates are selected, in each sub-frame, according to the Signal-plus- Noise Interference Ratio (SINR). In this system system, we focus on mitigating interference by deciding whether a base station can transmit during a given sub-frame. We refer to this decision as *base station scheduling*, to be distinguished from legacy *user scheduling* which occurs at each base station when it is allowed to transmit.

We aim at guaranteeing a minimum SINR to every user in the system by allocating in each sub-frame a subset of the available base stations while minimizing the number of sub-frames needed in order to schedule the complete set of base stations. To achieve this goal, we propose to schedule

base station transmissions in each sub-frame so that the SINR is greater than a threshold Th for every user $u$ in the network. The problem of minimizing the total number of sub-frames used to schedule once all base stations in the system, for a given minimum SINR (or threshold Th), is formulated using the following notation:

$$SINR_u = \frac{I_b^u}{N_0 + \sum_{j \neq b} I_j^u} , \qquad u = 1..U, \qquad u \text{ is connected to base station } b;$$

$$SINR_u \geq \text{Th} \Rightarrow \sum_{j=1}^{N} I_j^u \leq I_b^u \frac{1 + \text{Th}}{\text{Th}} - N_0 \triangleq \text{Th}^u , \qquad u = 1..U,$$

where $I_j^u$ is the signal received by user $u$ from base station $j$, with $b$ being the base station of $u$. With this notation, the problem can be stated as follows:

$$\begin{cases} \text{minimize } Z = \text{number of sub-frames needed to allocate all base stations once,} \\[2mm] \qquad\qquad \sum_{j=1}^{N} I_j^u \, x_{i,j} \leq \text{Th}^u, \qquad u = 1..U, \\[2mm] s.t. \qquad\qquad \sum_{i=1}^{N} x_{i,j} = 1, \qquad j = 1..N, \\[2mm] \qquad x_{i,j} \in \{0,1\}, \qquad i = 1..N, \qquad j = 1..N, \end{cases}$$

where $x_{i,j}$ equals 1 if base station $j$ is scheduled in sub-frame $i$, and 0 otherwise.

### 2.6.2 Proposed solution

The problem described so far is NP-hard since it is equivalent to a k-dimensional vector bin-packing problem with $U$ dimensions, and the lower bound for $Z$ is then easily computed as follows:

$$Z \leq L = \max_{u=1..U} \left( \left\lceil \frac{\sum_{j=1}^{N} I_j^u}{\text{Th}^u} \right\rceil \right).$$

To solve the problem of minimizing $Z$, we propose a heuristic consisting in a greedy algorithm for the mapping of base stations to sub-frames. We name the algorithm BASICS, which stands for BAse Station Inter-Cell Scheduling. The algorithm is designed to dynamically mitigate inter-cell interference caused to any possible user in the system under any possible user scheduling decision taken by the base stations. As a result, our algorithm is _user-scheduling-agnostic_ and does not require coordinated scheduling among base stations.

We propose a new heuristic rather than using existing heuristics for two main reasons. First, existing heuristics for multidimensional vector bin-packing problems are simple extensions of solutions designed for the one-dimension problem. Second, existing heuristics do not take into account the nature of the dimensions that describe the items to be allocated. In particular, they assume

that the size of an object is the same in any of the possible combinations of items in a bin. In contrast, in our case, the size of an object is the interference caused to mobile users *belonging to the scheduled base stations only*. Therefore, the weight associated to a base station (i.e., its *size*) changes any time a base station is removed from the list of candidate transmitters (e.g., since it is allotted to a sub-frame).

Our heuristic is based on a sum-based algorithm that solves k-dimensional vector bin-packing problems by collapsing all problem dimensions (i.e., the interferences to different users) into one unique value. This value is computed for each base station, and consists in the total interference caused by the base station to users belonging to other scheduled base stations. However, differently from conventional bin-packing problems, here the size of the items (base stations) to be allocated into bins (sub-frames) *changes* at any iteration of the algorithm. In particular, our heuristic represents a modification of the FFDSum algorithm [25] in which (*i*) the size of each item to be accommodated changes at each iteration, and (*ii*) items are accommodated into bins in order, beginning with the smallest one. Our algorithm allocates a new bin only when there is no more room left in the old bins to accommodate the remaining items. Note that existing algorithms for bin-packing would rather sort items from the largest to the smallest. The rationale behind our approach is as follows. First, when we start allocating base stations from the least interfering one, we have a chance to schedule together the highest number of not-previously-allocated base stations in the same sub-frame. This eliminates the highest number of base station candidates for the next sub-frame allocation. In turn, considering a uniform distribution of users, this procedure eliminates the highest number of users from the set of interfered users in the next iteration of the algorithm. As a result, the cumulative interference over the remaining users, due to the remaining candidate base stations in the next iteration, is likely to be much lower than in the previous iteration. In contrast, if we removed from the candidate set a base station generating less interference, we would have a high probability that that base station interfered fewer users. Thus, removing the least interfering base station would not only bring less benefit to the current sub-frame, but also would not reduce much the impact of that base station in the next sub-frame allocation (since the set of potential interfered users did not change much).

### 2.6.2.1 FLAVIA support

Base station scheduling uses signal quality measurements that are currently available at UEs and base stations. In addition, it involves the deployment of novel inter-base station cooperation features that rely not only on inter-cell coordination and cooperation features, but also on SON features. E.g., SON

operations are responsible for the access to the LTE Evolved Packet Core (EPC) to compute and distribute the scheduling of base stations periodically, and to force the *sub-frame blanking* of base stations when they are not scheduled.
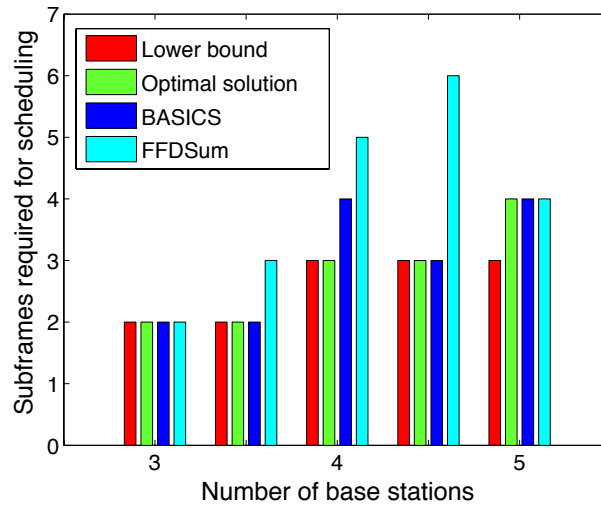
In the frame of FLAVIA, the BASICS algorithm can be implemented by acting on four services defined in WP2 and WP3: MEAS (the measurement service), ICIC, SSON, and SCHE (the MAC scheduler service). First, while BASICS runs in the EPC, it has to be fed with signal levels which are commonly measured by UEs and base stations and are collected at the base station side (through the MEAS service). These measurements shall then be used by the ICIC service to prepare the data to be communicated to the EPC via the inter-entity interface defined in the FLAVIA architecture, namely IAP2. ICIC needs also to keep the SSON service updated with the list of neighbors and potential interferers which are involved in self-organized networking operations. To this aim, the intra-MAC interface can be used, i.e., IAP5. SSON is responsible for managing the communication with the EPC and to pass to the SCHE service the base station decisions taken by BASICS (i.e., to fix the sub-frame blanking pattern to be used by the base station scheduler).

Therefore, the architecture of FLAVIA is ready to support frequency reuse 1 schemes with advanced ICIC operations exploiting SON and EPC features not yet present in the official standards.

### 2.6.3  *Some insights on resulting advantage and benefits*

The number of operations to be performed by BASICS to allocate base stations to sub-frames is at most $U \cdot \sum_{j=1}^{N} \frac{j \cdot (N-j+1)^2}{N}$. Therefore, the complexity of the algorithm is $O(U \cdot N^3)$, and the algorithm scales with the number of users in the system.

To compare performances achieved with our proposed heuristic algorithm, the performance of original FFDSum, and the one of an optimal (brute force) algorithm, we simulate a network with a variable number of base stations and users, using the OPNET simulator and Matlab.

***Figure 23: Number of sub-frames used with our proposal (BASICS) and with other base station scheduling approaches.***



***Figure 24: Throughput achieved with different base station scheduling algorithms.***

Figure 16 shows that BASICS finds the same number of sub-frames as the optimal solution, except for the case of 5 base stations in which it uses only one extra sub-frame. Furthermore, BASICS uses at most one sub-frame more than the theoretical lower bound. In contrast, FFDSum achieves significantly worse results (as expected from the discussion above). Figure 17 further illustrates the throughput performance obtained from these algorithms,

and reveals that FFDSum not only uses more sub-frames, but also provides worse throughput. In contrast, BASICS achieves near-optimal results.



*Figure 25: Throughput and fairness performance comparison of base station and user schedulers with 150 users.*



*Figure 26: Performance comparison of base station and user schedulers with fixed number of users per base station (8 users per base station).*

We next evaluate the performance of BASICS in terms of throughput and

fairness by using MATLAB, which allows to explore the impact of a large number of base stations and users under various network configurations. In order to assess its performance, we compare BASICS against the following two approaches: (*i*) normal network operation, in which all base stations are allowed to transmit in any sub-frame (referred to as "Legacy" in this deliverable), and (*ii*) a frequency reuse 3 scheme that partitions the network into three parts. For all cases, two intra base station schedulers are considered: round robin and proportional fair scheduling (for clarity of presentation, for BASICS and frequency reuse 3 we only show results achieved with the proportional fair scheduler, which are slightly better than with the round robin scheduler). We measure network performance in terms of the sum of the logarithms of the throughputs, as this is a well-accepted metric to compare different scheduling mechanisms in terms of efficiency as well as fairness. Note that for the case of frequency reuse 3, we normalize the throughput to the number of carriers utilized, i.e., 3. Figure 18 shows the performance of the above approaches when we fix the total number of users in the system to 150, and vary the number of base stations from 3 to 9. We observe from the figure that BASICS significantly improves network performance with respect to both legacy approaches and frequency reuse 3. Results achieved with frequency reuse 3 are similar to the ones achieved with BASICS only for scenarios with very few base stations. Next, we evaluate the impact of the number of users. To this end, we consider a network in which the number of users is proportional to the number of base stations. Specifically, we simulate 3 to 10 base stations with 8 users each. Figure 19 depicts the sum of logarithmic throughputs as a function of the number of base stations. Also for this case, BASICS exhibits the best performance over all the other approaches. On the one hand, the gain of BASICS over the legacy schemes is of several logarithmic units (and hence substantial in a linear scale). On the other hand, the gain over frequency reuse 3 is lower until the number of base stations reaches 10 (which is explained by the fact that frequency reuse becomes less effective as the network density grows). Taking into account that frequency reuse requires multiple carriers to achieve worse results, we conclude from these results that BASICS provides substantial improvements in performance also for this case.

In summary, the key contributions of our proposal are as follows: (*i*) we formulated a novel base station scheduling problem, which is NP-hard; (*ii*) we designed an algorithm, that runs in polynomial time and scales with the number of users; (*iii*) we showed that our algorithm not only achieves better throughput performance with respect to state of the art schedulers, but also significantly improves fairness among users.

# 3 Novel cellular architectures and scenarios

## 3.1 Opportunistic scheduling with clusters (IMDEA)

### 3.1.1 Problem statement and motivation

Opportunistic scheduling was initially proposed to exploit user channel diversity for network capacity enhancement. However, the achievable gain of opportunistic schedulers is generally restrained due to fairness considerations that impose a tradeoff between fairness and throughput. In FLAVIA, we have shown via analysis and simulation that opportunistic scheduling not only can be used to increase network throughput dramatically, but also can be fair to the users when they cooperate, in particular by forming clusters.

Opportunistic schedulers have become a promising solution to cope with the mobile consumer traffic boom in cellular networks. This class of schedulers exploits multiuser diversity to reorder transmissions so that each user is served, with high probability, when it is in a good channel state. However, opportunistic schedulers are often not practical and not yet widely adopted since they achieve poor fairness levels with the scarce memory and computational resources available at the base station. In particular, schedulers that could potentially achieve a good tradeoff between throughput and fairness, e.g., the renowned Proportional Fair scheduler (PF) are too complex for the centralized architecture of the cellular networks and do not scale with the number of users in the network.
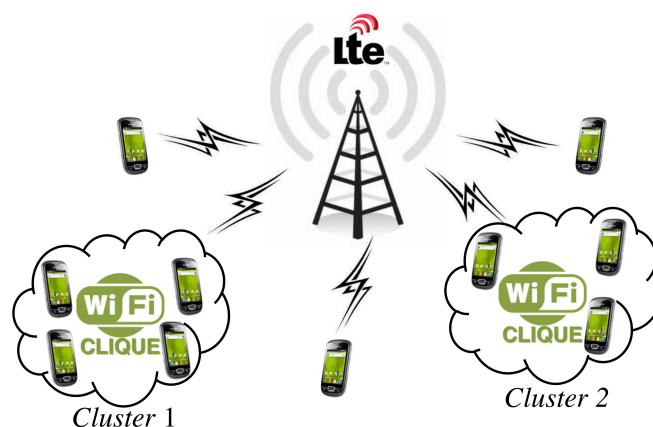
Aiming at increasing the capacity of the cellular network while providing user-level fairness and scheduling scalability, the main contributions of this FLAVIA research activity are as follows: (*i*) we have designed a novel network architecture based on cooperative communications with clusters of mobile users and opportunistic scheduling in cellular networks with dual-radio mobile users; (*ii*) we have designed novel opportunistic and cluster-based scheduling algorithms; (*iii*) via extensive numerical simulations, we have evaluated the performance of the proposed scheduling solutions and shown that they enable dramatic throughput gain and extremely fair throughput distributions among mobile users.

### 3.1.2 Proposed solution

*We propose to leverage smartphone's dual-radio interface capabilities to form clusters among mobile users, and we design simple and scalable cluster-based opportunistic scheduling strategies that would incentivize mobile users to form clusters.*

Interestingly, today's mobile devices are equipped with high processing power, large memory, and multiple radio interfaces. Therefore, mobile users may connect to each other using 802.11 interfaces, e.g., endowed with WiFi-Direct capabilities [26]. Particularly, multiple radios could be exploited for establishing cooperative communications, e.g., to form clusters, as illustrated in Figure 20.



*Figure 27: Cellular network with clusters of dual-radio mobiles.*

We propose to use the multi-radio capabilities of newly designed mobile devices to form cooperative clusters. The presence of clusters simplifies the scheduling operation and enables efficient radio resource utilization. *We design a multi-layer cluster-based scheduling mechanism in which the base station schedules clusters instead of users, while intra-cluster resource distribution is left to cluster members*. After cluster formation, the base station is notified of clustering decisions. Thus, whenever a packet is destined to a cluster member, the base station simply sends it to the cluster member that maximizes the throughput at that epoch, namely the *cluster head*. Therefore, *unlike existing clustering approaches, we propose to select the cluster head opportunistically in each frame*. The resulting scheduling mechanism consists of two elements: an algorithm to schedule clusters, and an algorithm to select cluster heads opportunistically within each cluster. As for cluster scheduling we propose simple and scalable algorithms such as round robin (RR), weighted round robin (WRR) and MaxRate (MR). We name these cluster-based algorithms CL(RR), CL(WRR) and CL(MR), respectively. In particular, in CL(WRR) we assume that airtime resources in the cell are allotted to cluster proportionally to the number of cluster members. Eventually, in order to allow non-head cluster member to communicate with the base station, cluster members use WiFi-Direct to exchange packets with the cluster head, i.e., all intra-cluster communications take place over a local wireless network.

### 3.1.2.1 FLAVIA support

The proposed two-tier scheduling strategy requires changes in the MAC layer of LTE and/or 802.16 or similar scheduled systems. To implement clustering and opportunistic scheduling with clustering, some novel functionality is required, which is not readily available in off-the-shelf platforms. First we need to change the scheduling algorithm in order to account for *cluster* as scheduling candidates, i.e., the scheduler has to merge the CQI info available on cluster members, select the cluster head, and rank cluster heads as scheduling candidates. This is possible with FLAVIA, and can be achieved by modifying the following services: Scheduling Strategy, QoS Strategy and Link Adaptation Service, which are defined in the architectural work of WP2 and WP3. Second, we need to modify the behavior of the data transport in both the base station and the UE, so that packets can be routed at layer 2 among cluster members. This can be achieved in the frame of the FLAVIA architecture by enhancing the standard data transport functions, and in particular by modifying the Data Transport Service defined in WP2 and WP3.

Therefore, thanks to the architecture proposed in FLAVIA, it would be possible to leverage the diffusion of currently available multi-radio devices to benefit from cluster-based scheduling.

### 3.1.3 Some insights on resulting advantage and benefits

We assume that user channels are independent and characterized by stationary Rayleigh fading, so that the average SNR experienced by each user does not change over time. However, each user can experience a different average SNR (i.e., network conditions can be heterogeneous). For the sake of tractability, we assume that mobile user's SNR can take one of three predefined average SNR values, which correspond to *poor*, *average*, and *good* users. The three designated average SNR values are chosen in a manner that the mean achievable rates—computed according to the MSC reported for an LTE system in Table 1 [27]- for *poor*, *average*, and *good* users are 20%, 50%, and 80% of the maximum transmission rate achievable in the system, respectively. With the thresholds and MCS values reported in Table 1, the designated SNR values are 7dB, 16dB and 23dB, respectively for *poor, average*, and *good* users.

***Table 3 Modulation and coding schemes, and respective thresholds***

| Modulation | Coding Rate | SNR [dB] | Implementation Margin (IM) [dB] | SNR+IM [dB] | Bits per symbol |
|------------|-------------|----------|--------------------------------|-------------|-----------------|

| | | | | | |
|---|---|---|---|---|---|
| QPSK | 1/8 | -5.1 | | -2.6 | 0.25 |
| | 1/5 | -2.9 | | -0.4 | 0.4 |
| | 1/4 | -1.7 | | 0.8 | 0.5 |
| | 1/3 | -1 | 2.5 | 1.5 | 0.67 |
| | 1/2 | 2 | | 4.5 | 1.0 |
| | 2/3 | 4.3 | | 6.8 | 1.3 |
| | 3/4 | 5.5 | | 8.0 | 1.5 |
| | 4/5 | 6.2 | | 8.7 | 1.6 |
| 16QAM | 1/2 | 7.9 | | 10.9 | 2.0 |
| | 2/3 | 11.3 | 3.0 | 14.3 | 2.66 |
| | 3/4 | 12.2 | | 15.2 | 3.0 |
| | 4/5 | 12.8 | | 15.8 | 3.2 |
| 64QAM | 2/3 | 15.3 | | 19.3 | 4.0 |
| | 3/4 | 17.5 | 4.0 | 21.5 | 4.5 |
| | 4/5 | 18.6 | | 22.6 | 4.8 |

Note that our clustering proposal does not conflict with the scheduler deployed at the base station. In fact, from the base station point of view, clusters are seen as users (characterized by higher demand than regular users). We assume that the user selected for transmission represents the cluster head; therefore, it can receive traffic for any other user in the cluster. As a consequence, we focus on the aggregate per-cluster throughout rather than on the per-user throughput. Note that per-user throughput can be derived from the aggregate per-cluster throughput if the cluster resources are divided equally among users. Indeed, through this paper, we assume that cluster resources are shared equally among users unless otherwise specified.

The incentive behind clustering can be easily observed by comparing the user's channel state probabilities. Figure 21 shows the impact of clustering on user's with *good*, *average*, and *bad* channel quality. The clustering impact is depicted in terms of MCS probabilities. As shown in the figure, cluster formation highly boosts the transmission rate of *poor* and *average* users, but it may be not that helpful for *good* users. Nonetheless, we will show that with our proposal the throughput of *good* users improves as well by a non-negligible quantity. Therefore, *good* users are incentivized to help users with lower channel qualities. In practice, *good* users can be encouraged to participate in clustering by receiving an extra quota for the portion of traffic that they forward for others users.

*Figure 28: Impact of clustering on channel state (MCS) pdf.*



*Figure 29: Clustering with users belonging to the same cell.*

Let us evaluate throughput and fairness in the network by numerically simulating a simple scenario including 3 clusters with fixed number of members, all of them belonging to a single cell, see Figure 22. In this scenario, clusters C1, C2, and C3 have 5, 10 and 15 users, respectively. We use a uniform distribution to select the SNR classes of users in each cluster. The experiment results include the mean, 5th and 95th percentiles over 2000 simulation runs. Note that clusters have no effect on user-based scheduler results as users are scheduled individually under RR and PF schedulers. Since it is not practical to show the throughputs of all users in one graph, we computed the average throughput for every cluster. When we show per-user throughputs, we report per-cluster throughputs normalized to the number of cluster members. Therefore, the per-user throughput can be interpreted as the average throughput received by a cluster member or, equivalently, as the

amount of resources allotted to the cluster for each member forming the cluster.
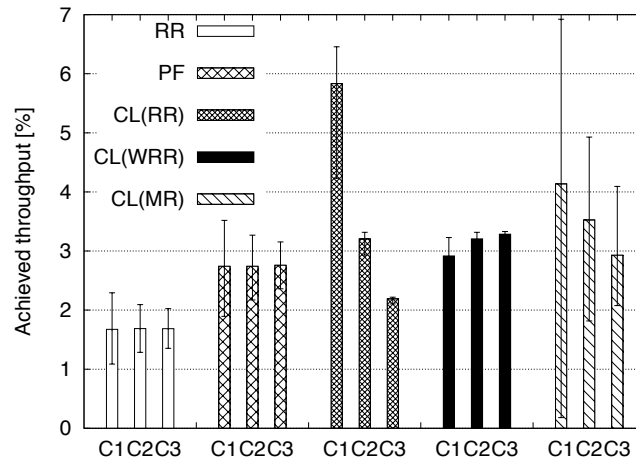
In Figure 23, we observe that users receive the least throughput under RR. The performance improvement due to PF is remarkable, but it is irrespective of cluster sizes. Cluster-based schedulers have different performances. First, CL(RR) exhibits an unfair distributions of per-user throughputs, due to the fact that it allocates equal airtime to all clusters regardless of their size. This behavior explains the reason why cluster C1, with the smallest size, has the highest per-user throughput. On the contrary, CL(WRR) allots airtime to clusters according to their size, which explains why users in all three clusters have almost the same throughput. Small differences in throughput distribution are due to the fact that clustering gain grows with cluster size. Eventually, CL(MR) operates based on the cluster's equivalent channel state, i.e., it schedule the cluster whose leader has the best channel state in the network, and does not take into account the size of clusters in scheduling decision, like in CL(RR). The maximum throughput is achieved under CL(MR), since it is a pure opportunistic scheduler. However, throughput distribution across users is unfair. Moreover, the high throughput variation observed for cluster C1 is due to the random user channel distribution within a cluster, which has higher impact on clusters with small number of users.

Figure 24 illustrates the aggregate throughput achieved in clusters, for the same scenario of Figure 22. As it can be noticed, under CL(RR) all clusters receive similar throughput regardless of their sizes, whereas with CL(WRR) the clusters with bigger size obtains more resources. CL(MR) achieves the best throughput, but also the highest variability.
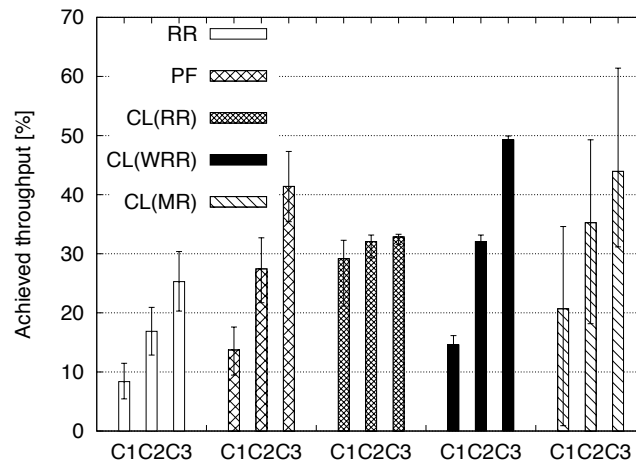
The total cell throughput is shown in Figure 25, from which it is clear how RR and PF are both outperformed by cluster-based schedulers. As expected, CL(MR), that always serves the cluster in the best channel, has the highest throughput. However, CL(RR) and CL(WRR) achieve comparable performance to CL(MR). Note that the difference between RR and CL(RR) consists in the clustering gain which in this scenario amounts to 45% of the total cell capacity. Overall, while CL(RR) is fair with respect to clusters as a whole, CL(WRR) results in fairer throughput distributions among users. CL(MR) achieves the best throughput, but it is neither cluster-fair nor user-fair. Figure 26 gives insights into the levels of fairness achieved by the different schedulers. In the figure, we report the Jain's fairness index among per-user throughputs. It is interesting to observe that our clustering proposal not only increases the throughput, but also it increases the fairness level. In particular, CL(WRR) advantages are three folds: (*i*) it provides nearly perfect fairness among users; (*ii*) it offers the possibility to gain a high throughput with respect to legacy RR
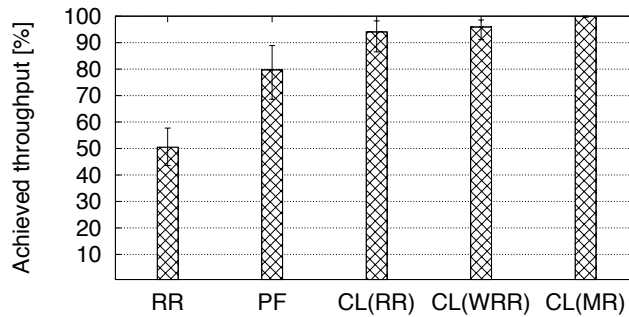
and PF schedulers; (*iii*) it allows each cluster to exploit the clustering gain proportionally to its size.



*Figure 30: Per-user per-Cluster throughput under different scheduling mechanisms.*



*Figure 31: Aggregate per-Cluster throughput under different scheduling mechanisms.*

*Figure 32. Aggregate cell throughput under different scheduling mechanisms.*



*Figure 33: Fairness under different scheduling mechanisms.*

In conclusion, the novel cluster-based scheduling schemes developed in FLAVIA avoid the need of trading off fairness and throughput by integrating the concepts of opportunistic scheduling and cooperative communications. We have shown that cluster-based scheduling substantially ameliorates the throughput (up to 50%) while maintaining high fairness among users. In particular, CL(WRR)—which assigns resources to the clusters in weighted round robin and selects clusters heads opportunistically—achieves throughputs close to the ones achieved by the MaxRate scheduler, but with nearly perfect fairness among users.

## 3.2 Dynamic assignment of UL/DL subframes in TDD systems (NEC)

In deliverable 5.2, we introduced a novel method of assigning uplink and downlink asymmetrically at adjacent base stations. In this deliverable, we briefly revisit the two-way interference channel and present a more specific approach as well as performance figures for it. The described approach as well as performance results are further detailed in [28].

### 3.2.1 Two-Way Interference Channel



*(a) Symmetric Assignment*       *(b) Asymmetric assignment*

**Figure 34 The two-way interference channel**

The two-way interference channel is illustrated in Figure 27. It considers two base stations (1 and 2), two user terminals (3 and 4), and interference between both paths. On the left hand-side, we see the symmetric assignment of uplink and downlink, i.e. we have the same interference channel in both forward and backward channel. On the right hand-side, we see that both paths operate asymmetrically, i.e. while one path is in uplink, the other one is in downlink. In D5.2, we already showed that the asymmetric assignment provides potential improvements of achievable rates particularly towards the cell-edge. We further detail this analysis by describing a specific transmission model which takes imperfect channel knowledge into account and by describing a particular cooperation scheme which exploits the potential benefits.

### 3.2.2 Transmission Model

In contrast to D5.2, we consider in this deliverable a transmission model which takes imperfect channel knowledge at the transmitter and receiver into account. We use again the same system setup as described in [D5.2, Section 2.2], however, the transmission model is defined by two phases with different
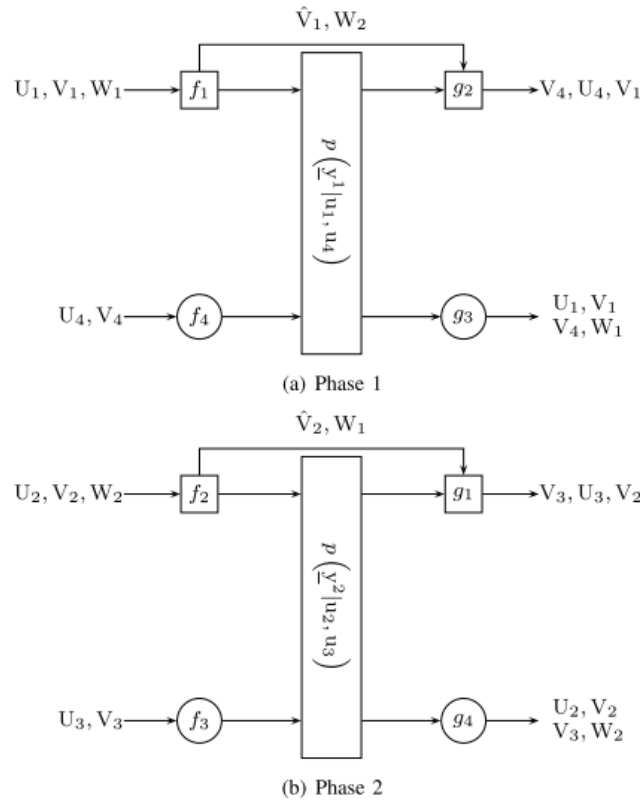
channel estimation error depending on the effective SNR in the uplink and downlink direction. Based on the estimated channel as the sum of actual channel and a Gaussian error term, we are able to express the transmission equation as the sum of an effective channel, an error term depending on the quality of the channel estimation, and the Gaussian noise. Details for this model are given in [29].

### 3.2.3 Slepian-Wolf Coding based Cooperation

In the following, we present a protocol for asymmetrically operating cells which is based on superposition coding and Slepian-Wolf (SW) coding. SW coding refers to lossless source-coding where the decoder has access to side-information while the encoder has access to the joint statistics of source-signal and side-information. Assume that the encoder wants to communicate message $W_1$ while the decoder has access to message $W_2$. Then, the encoder needs rate $R_1 \geq H(W_1 \mid W_2)$ to reliably communicate message $W_1$ while only having access to the joint pdf $p_{W1,W2}(\bullet,\bullet)$.

In our protocol, SW coding is applied to detect, decode, and subtract the interference caused by the transmitting base-station to the receiving base-station. Specifically, the receiving base-station uses its own channel output as side-information to decode the support message sent on the backhaul. After cancelling the interference, the receiving base-station detects and decodes the message from its assigned user terminal. Similarly, the receiving user terminal applies SW coding to the interference from the transmitting user terminal. Furthermore, we apply superposition coding [30], where the channel input of a node depends upon another message. This allows to overlay information for different receivers as in the broadcast channel [30]. This has also been applied by Han and Kobayashi [31] for the interference channel where a private and common message are overlaid. The private message is only decoded by the assigned communication partner while the common message is decoded by both receivers.

**Figure 35 Message exchange of the cooperation approach for the asymmetric assignment of uplink and downlink**

Figure 28 shows the messages which are exchanged by the described protocol in phases 1 and 2. In the following, we explain in detail the messages which are exchanged in phase 1 by the transmitting base-station 1 and user-terminal 4, and receiving base-station 2 and user-terminal 3. User-terminal 4 divides its message into two parts $U_4$ and $V_4$, which are both decoded by base-station 2 and user-terminal 3 decodes only the common message $V_4$ to reduce the interference signal power for the useful signals from base-station 1. Base-station 1 divides its messages into three parts, $U_1$, $V_1$, and $W_1$. The first part, $U_1$, is a private message which is only decoded by user-terminal 3, the second part, $V_1$, is a common message decoded by base-station 2 and user-terminal 3. Finally, the third part, $W_1$, in block b depends on message $V_4$ in block b-1 and supports user-terminal 3 during the decoding process of $V_4$. In addition, base-station 1 forwards support messages $\hat{V}_1$ and $W_2$ over an errorless backhaul to base-station 2, which are used to support the decoding process at base-station 2 and user-terminal 4.

In the first phase, base-station 1 transmits the superimposed messages $U_1$, $V_1$, and $W_1$ while user-terminal 4 transmits messages $U_4$ and $V_4$. Base-station 2
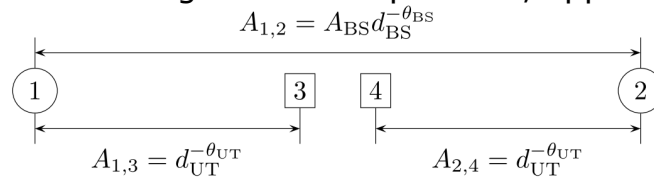
decodes both messages of user-terminal 4 and message $V_1$. Base-station 1 further supports the decoding at base-station 2 by providing additional redundant information over the backhaul ($\widehat{V}_1$) which is used together with side-information provided by the channel-output at base-station 2. The decoding process at user-terminal 3 follows along similar lines. At first, it decodes message $W_1$ which support the decoding process of $V_4$ (sent in the previous block), which itself is used to reduce the interference level for messages $U_1$ and $V_1$. Base-station 1 further provides additional information with $W_1$, which was forwarded by base-station 2 in the previous block. Based on this description, [32] derives detailed rate expressions for the deterministic Markov channel as well as the Gaussian channel. However, these rates are not repeated here but we rather focus on the results.
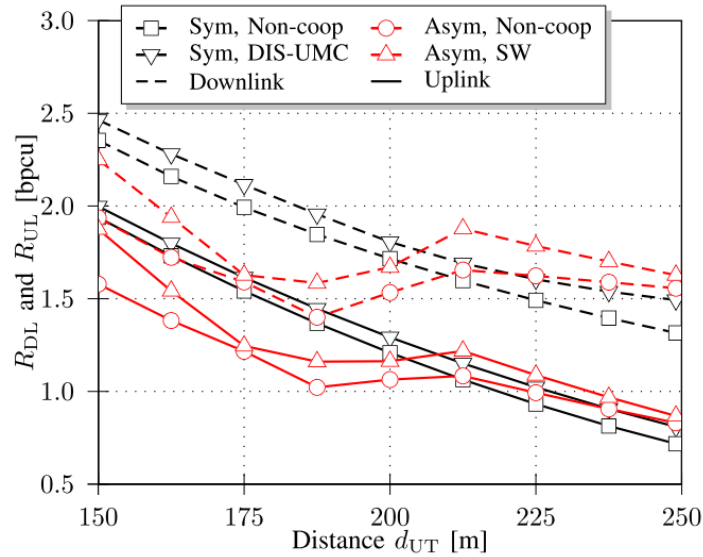
### 3.2.4  Results and Comparison

We compare the introduced protocol with with two cooperative multi-point (CoMP) approaches which were introduced in [33] [34]. Both approaches are abbreviated with DIS (distributed interference subtraction, applied in uplink) and UMC (unquantized message based cooperation, applied in downlink).



$$A_{1,2} = A_{BS}d_{BS}^{-\theta_{BS}}$$

$$A_{1,3} = d_{UT}^{-\theta_{UT}} \qquad A_{2,4} = d_{UT}^{-\theta_{UT}}$$

*Figure 36 Evaluation setup*

We compare these approaches based on the setup illustrated in Figure 29 where we assume that both base-stations are equipped with two antennas each and both user-terminals are equipped with a single antenna each. We apply an exponential path-loss model as shown in Figure 29 with $\theta_{BS}$=2.5 and $\theta_{UT}$=3.5. Furthermore, $A_{BS}$=10dB is applied in order to take the directivity of the base-station antennas into account. The compound channel matrix is chosen similarly to [34] with constant envelop but random phases. We evaluate the maximum common rate in downlink and uplink $R_{DL}$ and $R_{UL}$, respectively. Furthermore, we assume $N_p$=4 pilots, a cell-edge SNR in downlink of $\gamma_{DL}$=10dB and $\gamma_{UL}$=3dB in the uplink, which both result from a typical urban micro scenario as defined in IMT-Advanced [35].

***Figure 37 Performance results depending on the distance between user terminals and base stations under backhaul constraint***

Figure 30 shows the performance results depending on the user distance and under a backhaul constraint of ß=4. We can see from these results that the symmetric assignment outperforms the asymmetric assignment towards the cell-center by about 20% while towards the cell-edge the asymmetric assignment of uplink and downlink outperforms the symmetric assignment by up to 25%. When both user-terminals are closer to the cell-edge, they resemble a strong interference-channel where it is possible to detect and cancel the interference. However, towards the cell-center both user-terminals cause a weak interference channel where interference can neither be ignored nor reliably canceled, which implies the performance drop at about d=175m. Interestingly, the asymmetric assignment outperforms CoMP even though it implies much less complexity.

***Figure 38 Performance results depending on the backhaul constraint at distance d=225m***

Figure 31 shows the performance results depending on the available backhaul resources and at distance d=225m (cell-edge). We can see that the asymmetric assignment achieves its maximum performance with significantly less backhaul resources compared to CoMP, i.e. while the maximum for the asymmetric assignment is achieved at ß=2, CoMP requires about ß=10 to achieve its maximum performance. This confirms that the asymmetric assignment is more backhaul-efficient.

## 3.3 Coded unicast downstream traffic in a wireless network

### 3.3.1  Problem statement and motivation

The inherent broadcast nature of the wireless medium, which allows each transmission to be heard by all users simultaneously, makes network coding techniques pertinent. In such techniques, nodes do not necessarily forward incoming packets. Rather, they can transmit a manipulation (usually a linear combination) of their incoming data.  However, in order for such a combination to be valuable to multiple users, each such user needs to possess different piece of the information encoded into the combined packet. Accordingly, one of the key challenges in network coding techniques is to decide which packets to manipulate in each transmission. While efficient algorithms answer this challenge in the multicast setting, the problem of multiple unicast remains open.

On the down side, the wireless medium characteristics make wireless transmissions susceptible to losses due to noise and interference (i.e., low SNR and SINR). In order to cope with packet loss in MAC layer, conventional wireless protocols rely on retransmissions. In such protocols, each packet has to be acknowledged by the intended receiver. Packets which are not acknowledged are retransmitted over and over again until they are received successfully by the receiver, or until dropped by the sender.

Typical last mile wireless Internet access architecture comprises a gateway, e.g., an Access Point (AP) or a Base Station (BS), to which all clients are wirelessly connected (e.g., WiFi, WiMAX, LTE). In such architecture, all traffic to and from the wired Internet must pass through the gateway via the wireless medium. Accordingly, all transmissions by the gateway are potentially heard by all clients associated with this gateway. In this paper, we utilize these aforementioned wireless properties of channel, protocol and last mile architecture and suggest coded wireless retransmissions for downstream traffic. In particular, we suggest a novel scheme which is based on *Markov Decision Process* (*MDP*), that combines multiple MAC layer retransmissions which are intended to different receivers, into a single packet transmission.

At the basis of our work stands the already well understood concepts of *network coding* [36]. In this pioneering work, intermediate nodes in the network perform coding operation on the data in order to achieve certain rate goals. Indeed, it was shown that network coding can improve the network throughput significantly, and achieve the optimal performance in the multicast scenario. The theory of network coding includes linear as well as non-linear coding techniques. In this study, we focus on coherent linear network coding.

Following several important works discussed various practical issues in network coding. [37] introduced the idea of generations, and suggested coding only over packets of the same generation. As generations advance, old generations are flushed. In a sense, the concept is useful in this work as well, when we suggest that if a user acknowledges receiving a packet intended only to him, neighboring users who overheard it in a previous transmission and buffered it, discard it. Practical aspects of network coding also include several key works on opportunistic coding. That is, protocols, algorithms and analysis aimed at understanding which packets to send coded, and which coding coefficient to use, given the senders (maybe limited) knowledge on the data available at the receivers. Coding using only local information and opportunistic network coding was first introduced in [38] [39] as COPE. While decentralized, this ground-breaking work was not tailored to the multiple unicast with one sender scenario we consider in this paper. A polynomial time centralized algorithm, yet with guarantees only for the multicast scenario, was given in [40].

On the practical side, several works implemented network coding concepts on real networks. We focus here only on works whose implementation is below the application layer (i.e., excluding network-coded content distribution and related works). A pioneering work in this context is the already discussed [39]. The COPE header in this work resides between the routing and the MAC headers. The implementation runs on a 802.11a network as a user space daemon, that is, it sends and receives raw 802.11 frames. Random linear network coding on the iPhone was studied in [41], though, again, the implementation is on top of the WiFi driver, and not within it. Using Nokia mobile devices was suggested in [42]. In [43], the authors considered a chain topology, and gave numerical evaluation of the suggested iCORE scheme. Still, iCORE is a user space daeamon, which uses WiFi, but does not alter it. In this study, the implementation was within the WiFi driver, rendering the coding procedure transparent to all above layers.

### 3.3.2 Proposed solution

We suggest an ongoing process in which the BS alternates between transmitting uncoded and coded packets. Receivers acknowledge packets they have received successfully and in addition provide feedbacks to the BS regarding packets they overheard which were not meant for them. Based on these feedbacks the BS chooses which retransmitted packets (if any) should be coded in each retransmission. Our model is inherently not multicast each receiver has a different stream as its demand. Moreover, we assume an infinite horizon model, where there is no point in time in which all demands are met and the system reaches a terminating step. Packets arrive at the BS continuously, and only the current packets for each user are available for

coding. Based on this model, we are able to analytically solve throughput problems. We believe that these two aspects of our model are of key importance, since this is the typical use of most wireless Internet access networks.

Second, we show that the aforementioned continuous transmission process can be modeled as a discrete time stochastic process, in which at each state the next state is determined solely based on the BS decision which packet to transmit next (i.e., which coded or un-coded packets should comprise the next transmission) and based on the channel state of each and every receiver which determines which nodes receive the next transmission. We suggest a BS policy which is based on Markov Decision Process theory, in which the reward attained in each iteration corresponds to the number of successful packets received in each transmission.

Third, we leverage this continuous, infinite time stochastic model, to compute stationary behavior, which in turn allows us to define convergence, calculate the resulting asymptotic performance efficiently (using only a set of linear equations), and assess the benefit in coding directly and analytically. Specifically, we give the matrix equation that computes the cumulative expected reward (equivalent to the system throughput when a unit reward is given to decoding of one packet) for any state in the system given the transition probabilities and the reward vector. This enables us to directly compute the performance of any coded or un-coded strategy. For the two user case, we indeed give a few possible strategies and compute the resulting performance.

Fourth, we show that in order to reach an optimal decision, the BS needs to consider all possible future states of the system, channel states of all users and all possible actions and outcomes. This procedure certainly cannot scale to large number of users. Accordingly, we suggest a greedy approach in which at each transmission the BS tries to maximize the instantaneous reward received for each transmission (as opposed to maximizing an expected or discounted reward, which takes into account the expected rewards at future states). We further suggest an enhancement to the greedy approach, termed semi-greedy approach, which takes some concern into the future, without adding significant complexity to the greedy approach. In the semi-greedy approach, we also suggest a direct analysis for the simple case of two receivers, which besides the analysis of this simple case, also provides some insight into much larger scenarios. We evaluate both schemes via an extensive set of simulations, which show that our approach attains high gain over the traditional un-coded transmissions while maintaining long time fairness. Moreover, we show that the semi-greedy approach exploits the multi-user diversity in the system,

putting more emphasis on serving the users with the best channels conditions at any given time.

### 3.3.2.1 FLAVIA support

Network coding is a promising MAC enhancement that is not supported by current wireless architecture. However, such an enhancement is feasible with FLAVIA. Specifically, the data transport service and its IAP2 (Inter-entity interface) and IAP3 (MAC Services Functions interface) are ideal for such a scheme.
Since, the data transport service was not implemented in the scheduled based prototype, we chose a WiFi platform to implement and experiment our network coding scheme, as described next.

### 3.3.3 *Some insights on resulting advantage and benefits*

We implemented our scheme on a WLAN topology in which a single AP transmits unicast traffic to two receivers. We show that the suggested scheme can be easily implemented over a typical 802.11 card, with some modifications to the wireless driver. To the best of our knowledge, this is the first implementation of these concepts within the WiFi driver, and transparently from the upper layers. We further show that at least for this simple case, the experimental gain agrees with the one predicted by our analytical model.
We shall describe the simple WiFi implementation of the suggested scheme using off the shelf 802.11 devices. Specifically, we implement the scheme on cards with an Atheros chipset (ATHEROS AR5007G chipset), operating the open source ath5k drivers (2.6.32 version) for a Linux environment (kernel 2.6.32). We realize the suggested scheme in a simple network comprising two stations and an AP. Due to some hardware limitations we implement a slightly modified version of the suggested scheme which we describe below. We distinguish between the enhancements required by the AP (transmitter), required by the stations (receivers) and those necessitated by the standard.
**Frame format -** In order for users to keep track of the packets received intended for other users, the header of each uncoded frame includes a two byte sequential frame index. We utilize a special multicast address to mark all coded frames. In addition, in all coded frames, besides the two byte sequential frame index which is included in the header, an additional two times two bytes are included in the header, indicating the sequence numbers of the two frames that are coded.
**Station -** In contrast to 802.11, and in order to support the NC procedure, a station needs to receive packets not addressed to it. Accordingly, we set each station Network Interface Card (NIC) to work in promiscuous mode, which means that each station captures all frames sent by the AP, even if it is not its intended addressee. If a regular frame is received successfully at the

addressee station, the station sends immediately an ACK, in accordance to the 802.11 standard. On the other hand if the station receives a frame which is not addressed to it, it stores it locally in a hashed buffer, as illustrated in Figure 32(a). Once a coded frame is received, (recognized according to the designated multicast address) the relevant two byte header with the frame sequence number is used to locate the hashed frame. Then, the station retrieves the missing frame by XORing the XORed received frame with the hashed frame (if available). If a frame cannot be decoded from the NC retransmission (e.g., due to unavailability of both frames), no action is taken. Note that in such an event the packet is not going to be retransmitted, i.e., it is going to be dropped. Furthermore, it is important to note that even though we ran our experiment only on two receivers hence when receiving a coded packet the receiver knows that one of the XORed packets is intended for it, in the first part of the payload we also  XORed the two MAC addresses of the intended receivers, such that our implementation also applies to more than two receivers.



*Figure 39: (a) System topology and settings. (b) Access point buffer. (c) Access point retransmission scheme.*

**Access Point -** The first modification in the AP driver to allow coded retransmissions, is to stop its automatic retransmissions. Accordingly, we set the 802.11 retry limit to 0, i.e., no retransmission attempts. Nonetheless, in contrast to 802.11 where a frame is dropped when it reaches the retry limit, in our implementation if the AP does not receive an ACK message for a certain packet, it stores it in a pre-allocated buffer. A separate buffer is allocated to each user, Figure 32(a). As soon as two "failed" frames (waiting for retransmission) addressed to each of the two stations are available, the AP codes the two frames into a single retransmission frame. Coding is obtained by using a XOR operation. Then, the XORed frame is transmitted to a predefined

multicast address which differentiates coded frames from regular uncoded frames, Figure 32(c).

In contrast to the theoretical algorithm the AP does not work in a Stop and Wait manner, in which it stops all transmissions to a certain station upon frame failure until the frame is either received correctly or dropped. Rather, in our implementation it stores the un-acknowledged frame in the aforementioned buffer and continues sending subsequent frames to this user (i.e., selective repeat manner). Note that the unique sequential frame index included in each frame enables the support of such selective repeat mechanism and allows frame reordering. In order to avoid buffer overflows due to the selective repeat mechanism (which has an infinite window size), as soon as the buffer size crosses a threshold, all frames in the buffer are transmitted uncoded and the buffer is flushed. Additionally, for each un-acknowledged packet which is stored in the buffer, a time stamp is attached. A time-out mechanism is implemented such that the packet is retransmitted uncoded upon expiration of the time-out. It is important to note that in contrast to the algorithm presented earlier, no status packets are sent by the users to indicate which packets in their buffer are meant to the other user and were not acknowledged. In our implementation the AP assumes that each un-acknowledged packet is received by the other receiver, hence can be used for the coded packets. Consequently, the AP can send a coded packet that one or both receivers cannot really decode. Accordingly, as previously mentioned a retransmitted packet which cannot be decoded is lost. An enhancement in which a user periodically or upon request sends a status message which includes the unreceived packets can be easily implemented.

In our experiment, and in accordance to the algorithm, the AP generated packets to each one of the users at constant rate. Whenever two un-acknowledged packets were stored at the AP, one for each user, a single XORed packet was sent. In order to control the loss probability on each link, we artificially dropped packets at the receiver according to a fixed probability denoted by $p$. We varied the loss probability between zero (i.e., all packets are received successfully) and one (i.e., all packets are dropped). Note that coded packets were also subject to losses, according to the same probability $p$ as uncoded packets. For comparison we also show in some of the figures analytical results of three other schemes, (i) no retransmission - each packet is transmitted exactly once, accordingly an unsuccessful attempt on the first transmission results in packet loss (denoted by dashed line in the figures). (ii) Uncoded with single retransmission - each packet can be retransmitted at most once (i.e., retry limit is set to one) where the retransmitted packets are sent uncoded (denoted by dashed dotted line). Since the uncoded retransmissions scheme allows for the same number of sent packets more overall

transmissions than the coded scheme, we also examine a hybrid of the two first schemes, in which the number of transmissions (rather than the number of transmitted packets) is the same. (iii) Hybrid scheme - each un-acknowledged packet is retransmitted once more with probability half or dropped with probability half (denoted by dotted line).

We first examined the effect of network coding on air time utilization. We compare the number of successfully received packets by both receivers with the total number of packets sent, i.e., $\frac{total\ received\ packets}{total\ transmitted\ packets}$. Note that system utilization for all uncoded schemes is the same as it counts the number of successfully received packets, regardless of whether or not the received packet is a retransmission. Figure 33(a) depicts the results.
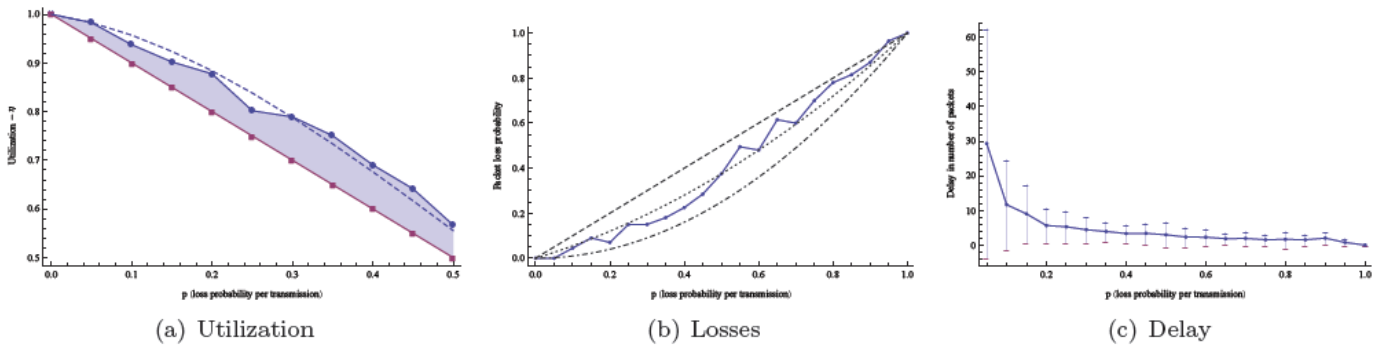
As can be seen in Figure 33(A), the coded schemes (denoted by circles in the figure) are always better than the uncoded schemes (squares) as far as airtime utilization is concerned. Note that per-packet overhead is not taken into account in the figure. Nonetheless such overhead is negligible, an extra 2 Bytes for uncoded packets and an extra 12 Bytes for the coded packets. The dashed line represents the expected analytical results for the coded scheme (expressed as the fraction of successfully received packets).

Next we evaluate the effect of not sending status packets. Recall that in our implementation the AP does not know which packets were received by each unintended user, and assumes that each unacknowledged packet was received by the other user. Furthermore, XORed packets are not acknowledged by the receivers. Accordingly, a XORed packet which is not received or cannot be decoded (i.e., its coupled packet was not received) by the receiver is lost. In Figure 33(b) we show the packet loss probability as a function of the link loss probability, $p$.

As expected the no retransmission scheme generates the highest packet drop for all values of $p$. Since in the uncoded single retransmission, the retransmissions are not coded and are dedicated to the intended receiver, the least drop packets are produced. Obviously this drop packet gain comes at the price of more packets being sent altogether. Interestingly the hybrid scheme is inferior to the coded scheme with respect to packet loss probability, i.e., more packets are dropped than at the coded scheme for $0 < p < 0.5$ and is superior as far as packet loss probability is concerned for $0.5 < p < 1$. The reason is that for the coded scheme to receive a coded packet successfully, relies not only on the acceptance of the coded packet itself but also on receiving the coupled packet successfully. Accordingly the mean packet loss probability is $p(p + (1 - p)p)$, where the first $p$ relates to the original transmission loss, and the terms in the parentheses refer to the retransmission of the coded packet which can be either lost or received successfully but its

coupled packet was lost. For the hybrid scheme, the mean packet loss probability is $\frac{1}{2}p + \frac{1}{2}p^2$, which indeed is greater than the coded loss probability for $p < \frac{1}{2}$ and less for $p > \frac{1}{2}$.



| (a) Utilization | (b) Losses | (c) Delay |

***Figure 40: {(a) Comparison of the total number of successfully received packets with the total number of packets sent. Circles - coded scheme, experimental; Squares - uncoded; Dashed - coded, theoretical. (b) Packet loss probability as a function of the link loss probability. Dashed - scheme (i); Dashed-dot - scheme (ii); Dotted - scheme (iii); Circles - coded, experimental. (c)Average number of packets transmitted between the first transmission attempt of an unsuccessful packet and its coded retransmission.***

Finally, we examine the retransmission delay due to the coding mechanism. Recall that the AP waits for un-acknowledged packets from both receivers before sending a retransmission. In Figure 33(c) we show the average number of packets which are transmitted between the first transmission attempt and its coded retransmission.

Obviously, and as can be seen in the figure, the greater the link loss probability $p$, the less the number of transmissions needed before the AP has two un-acknowledged packets, one for each receiver, hence can send a coded packet. On the other hand the less $p$ is, the longer a retransmitted packets needs to wait before it can be coupled with another lost packet to the other receiver.

More details can be found in [44].

# 4 Analysis Modelling and System simulations

## 4.1 Analysis of fundamental trade-offs in scheduling and resource allocation (NEC)

In this Section, we derive a framework to evaluate achievable rates in a multi-user OFDMA system.

The granularity of resource assignments in an OFDMA system determines the gains in spectral efficiency and requires to consider signalling overhead and finite channel coherence. Furthermore, practical systems suffer from partial CSI. These impacts are taken into account in the following framework in order to derive achievable net-rates which may significantly differ from gross-rates. Using this framework, we are able to optimize the MAC layer operation through the FLAVIA framework and based on the objective to optimize the net throughput.

### 4.1.1 Motivation

Multi-user systems offer the possibility to gain through multi-user diversity, i.e. if multiple users experience different channel conditions, the scheduler at the base station may choose the user with the best channel or the one with the highest utility function value. Hence, the benefits of multi-user scheduling depend upon the availability of channel knowledge at the base station as well as on the quality of this channel knowledge. Furthermore, the available resources are partitioned into resource blocks which are assigned to the individual users. Depending on the chosen resource block size, the system has to carry signalling overhead to identify the individual resources and the system has to cope with finite channel coherence, i.e. within one resource block the channel may not be constant and therefore the channel knowledge which is used to select a user may only be valid for part of a resource block. Hence, there is the need to find an optimal and flexible resource partitioning and assignment method that reduces the signalling load but takes the imperfect channel knowledge into account. Such an analytical framework is derived in the following, which allows to determine the achievable net-rates for opportunistic scheduling in a Rayleigh fading environment, taking into account exponential path-loss, partial CSI at transmitter and receiver, signalling overhead, and finite channel coherence in time and frequency.

### 4.1.2 System, Channel, and Scheduler Model

In order to analyse the mentioned system, we consider a TDD-OFDMA system with a central transmitter (base station) and K user terminals, each having a

single antenna. The system uses a bandwidth of B MHz divided into N subcarriers, which are divided into blocks of $N_f$ subcarriers and $N_t$ OFDM-symbols. Due to the consideration of a TDD system, we assume channel-reciprocity such that forward and backward channel are equivalent. The base-station estimates the channel based on pilots sent by the user terminal and uses this channel knowledge for its scheduling decisions.

All K users are distributed uniformly in a circular area with radius R=1 and distance $r_k$ of user k to the base station. The channel is described by

$$y_{f,t}^k = r^{-\eta/2} h_{f,t}^k \cdot x_{f,t}^k + n_{f,t}^k$$

where $h_{f,t}^k$ follows a Rayleigh distribution with an envelope correlation following the Jake-spectrum

$$\rho(\Delta f, \Delta t) = \sqrt{\frac{J_0^2(2\pi v/c\Delta t)}{1 + (2\pi\Delta f)^2 \sigma_\tau^2}}$$

Furthermore, user input signal $x_{f,t}^k$ and noise $n_{f,t}^k$ are i.i.d. Gaussian random processes. If for each resource block Np pilots are used, the capacity under imperfect channel knowledge can be bound by

$$C \geq R = \log\left(1 + \sigma_{eff}^2 \gamma_{f,t}^k\right), \qquad _{eff}^2 = \frac{N_p}{N_p + 1 + \frac{1}{\bar{\gamma}}} \bar{\gamma} < \bar{\gamma}$$

where $\bar{\gamma}$ is the expected user SNR (including path-loss) and under the assumption that a MMSE estimator and pilots are uniformly distributed, and $\gamma_{f,t}^k$ is the instantaneous SNR of user k on subcarrier f and OFDM symbol t.

In this framework, we assume a normalized-SNR based scheduler [45] which is asymptotically fair and achieves a similar performance as the proportional fair scheduler [46]. Hence, the base-station uses the estimated $|h_{ft}^k|^2$ and choses the users with the highest normalized channel gain.

### 4.1.3  Achievable Rates

In the following, we want to first solve the following problem

$$R_{cell}^* = \underset{N_f, N_t, N_p}{\arg\max} R_{cell}(N_f, N_t, N_p)$$

$$R_{cell}(N_f, N_t, N_p) = E_\gamma \left\{ \frac{1}{NN_t} \sum_{(f,t) \in N_t \times N} R_{N_f, N_t, N_p}(f,t) \right\}$$

$$\approx \frac{1}{NN_t} \sum_{(f,t) \in "central\ RB"} E_{\gamma_{f,t}} \left\{ R_{N_f, N_t, N_p}(f,t) \right\}$$

As detailed in  [32], we can show for users at equal distance R=1 that Rcell can be given by

$$R_{cell} = \frac{K}{\log(2)} \sum_{n=0}^{K-1} \left[ \binom{K-1}{n} \frac{(-1)^n}{1+n} \underbrace{exp\left(\frac{\sigma_{eff}^{-2}}{1 - \frac{n}{n+1}\rho_{f,t}^2}\right) \Gamma\left(\frac{\sigma_{eff}^{-2}}{1 - \frac{n}{n+1}\rho_{f,t}^2}\right)}_{R'(\bar{\gamma})} \right]$$

with the gamma-function $\Gamma(x) = \int_x^\infty \frac{e^{-t}}{t} dt$.

Since we are interested in uniformly distributed users, we need to further refine R':

$$\alpha_1 = \frac{N_p + 1}{N_p \sigma_x^2/\sigma_n^2} \qquad \alpha_2 = \frac{1}{N_p(\sigma_x^2/\sigma_n^2)^2}$$

$$\alpha_3 = \frac{k}{k+1} \frac{J_0^2\left(2\pi\frac{v}{c}\Delta t\right)}{1 + (2\pi\Delta f)^2\sigma_\tau^2} \qquad \alpha_4 = \frac{\sigma_n^2}{N_p\sigma_x^2}$$

$$R'(\bar{\gamma}) = \int_0^1 2r\, exp\left(\frac{\alpha_1 r^\eta + \alpha_2}{1 - \frac{\alpha_3}{1 + \alpha_4 r^\eta}}\right)^{2\eta} \Gamma\left(\frac{\alpha_1 r^\eta + \alpha_2 r^{2\eta}}{1 - \frac{\alpha_3}{1 + \alpha_4 r^\eta}}\right) dr$$

This integral cannot be solved in closed form, hence, we need an approximation to it. For path-loss exponent η=2 and under assumption of sufficiently large SNR, we can approximate R' by

$$R'(\eta = 2) \approx \int_0^1 2r\, exp(\beta r^2)\Gamma(\beta r^2)dr, \quad \beta = \frac{\alpha_1}{1 - \frac{\alpha_3}{1 + \alpha_4}}$$

$$= \frac{1}{\beta}\left(\gamma_e - \log\left(\frac{1}{\beta}\right) + e^\beta \Gamma(\beta)\right)$$

where $\gamma_e$ is Euler's constant.

### 4.1.4 Net-Rate Optimization

The previous part derived the achievable rates without considering net-rates. In order to derive the net-rates, we have to derive the relative signalling overhead 1-θ which gives the relative amount of resources that are required to communicate the resource map and for pilot signals. Given the resource width of $N_f$ subcarriers and K users, we need $N/N_f$ assignments of logK bits each. In a practical system, the resource map is communicated in a predefined region, hence, the expected data rate to each user can be approximated by

$$E_\gamma(R) = \int_0^1 2r \int_0^\infty \log(1 + \sigma_{eff}^2\gamma)e^{-\gamma}d\gamma dr$$

$$\approx \frac{1}{\alpha_1}\left(\gamma_e - \log\left(\frac{1}{\alpha_1}\right) + e^{\alpha_1}\Gamma(\alpha_1)\right)$$

Therefore, we need on average

$$N_\theta = \frac{Nlog(K)}{N_f} \cdot \frac{1}{E_\gamma(R)}$$

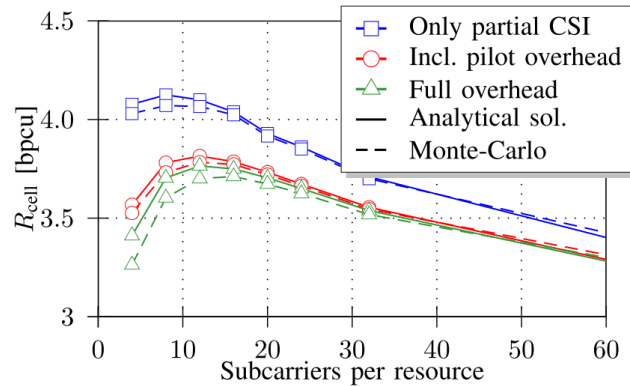bits to communicate the resource map. Now, taking also into account the overhead implied by pilot signals, we have

$$\Theta = \left(1 - \frac{\frac{N}{N_f}\frac{\log(K)}{E_\gamma(R)}}{NN_t} - \frac{N_p}{N_f N_t}\right)^+$$

$$= \left(1 - \frac{1}{N_f N_t}\left(\frac{\log(K)}{E_\gamma(R)} + N_p\right)\right)^+$$

Note that the multi-user scheduling gain is in the order of log log K while the overhead increases with log K, hence, eventually the overhead will outweigh the multi-user even though this only visible for very large numbers of users.

### 4.1.5  Results

The presented framework is evaluated using a typical OFDM setup with B=10MHz, N=1024 subcarrier, and delay-spread $\sigma_\tau$=2μs. We compare the pilot setup and resource block setup of LTE and WiMAX.
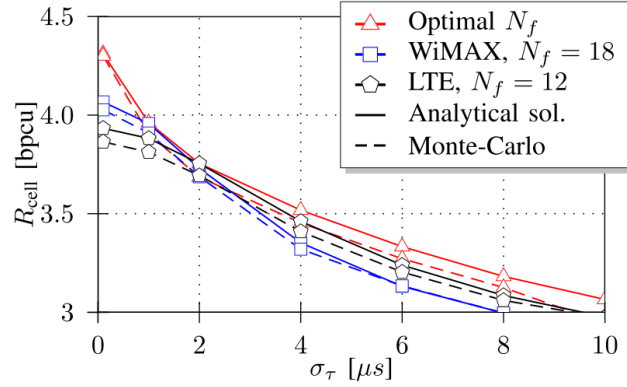
Figure 34 shows the achievable rates if imperfect channel knowledge is assumed, pilot overhead is included, and also the resource map overhead is considered. The results show both the analytical solution (solid lines) and based on a Monte Carlo simulation (dashed lines). The results are acquired for K=10 users, cell-edge SNR of 5dB, and $N_t$=6 OFDM-symbols.



*Figure 41 Achievbale rates as a function of resource block width $N_f$*

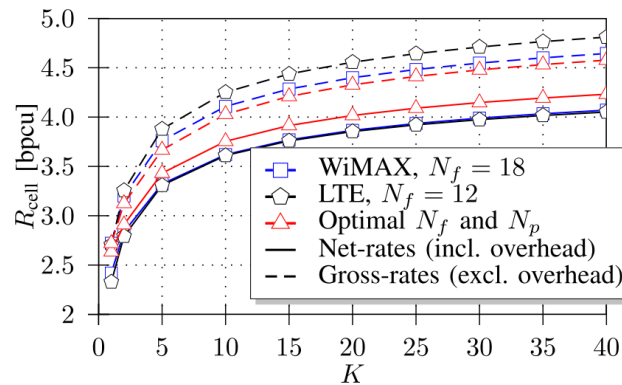We can see in Figure 34 that the loss due to overhead increases to about 10% for the highest net-rate value at $N_f$=12. The results further show that analytical and simulation results fit very well and the error does not exceed 4% and is for most parameter values less than 1%.

***Figure 42 Achievable net-rates as a function of the delay spread***

Figure 35 shows the achievable net-rates as a function of the delay-spread. We can see that at low delay-spread (higher coherence bandwidth), WiMAX provides higher net-rates than LTE due to the larger resource blocks. By contrast, at higher delay-spread (lower coherence bandwidth), LTE outperforms WiMAX due to its smaller resource blocks. However, at $\sigma_\tau$=2ms, which is considered as practically relevant, both WiMAX and LTE perform similarly and provide almost the same net-rates.



***Figure 43 Achievable rates depending on the number of users***

Finally, consider Figure 36 which shows the performance depending in the number of users. Firstly, we can see that the achievable rates increase with loglogK in the number of users. We can further see that the optimal choice of $N_f$ and $N_p$ provides lower gross-rate while providing higher net-rates. Similarly, WiMAX provides slightly higher gross-rates but provides the same net-rates as LTE.

# 4.2 An Extensible LTE Modelling Solution (CNIT)

### 4.2.1 Motivation

In order to develop new or improve existing system features in the FLAVIA architecture specified for 4G technologies, there needs to be a practical way to evaluate these mechanisms before implementing and deploying them in the actual system. Typically, such modelling is solved by the Discrete-Event Simulation (DES) technique.

However, DES is time and space inefficient, especially when it comes to system scenarios that, by virtue of the flexibility introduced by FLAVIA, may be exceedingly complex. This makes DES impractical to study the system performance operating new or improved mechanisms.

A different solution technique, called Traffic-Centric Modelling (TCM), is thus necessary to make the evaluation of FLAVIA scenarios more practical by significantly reducing time and resource usage. An extensive model of an initial LTE single cell scenario was therefore developed to illustrate the possible resource reduction gains achievable by using TCM, while still achieving accurate performance estimates.

### 4.2.2 Extensible model

Following the TCM principle of viewing the system from the perspective of that which requires service, i.e., the traffic, as opposed to that which serves, an underlying semi-Markov process (SMP) is used to capture the experience of the unit of traffic (t-unit) to be served in the system. In the LTE system, we consider both dynamic connection- and packet-level activity. Therefore two t-units were nominated and consequently two SMPs were defined.

Then, the effects of the system hardware and software are considered from the t-unit's perspective, defining the state transitions of each SMP according to the experience the t-unit may have in the system.

Since the system is very complex, no mathematical formulae are determined for these transition probabilities nor the sojourn times. Such formulations would require that a great deal of oversimplifying assumptions be made to achieve a tractable analytic solution.

Instead, an execution algorithm is developed to drive the underlying ergodic SMP models to convergence. The steady-state vector can then be calculated from the converged transition probabilities and finally, the required performance metrics can be calculated through formulations of the SMP (experience) parameters solved.

The resulting model is highly extensible since the execution algorithm resembles an activity-scanning simulator, while the SMPs can be adjusted by including or removing transitions in its definition.
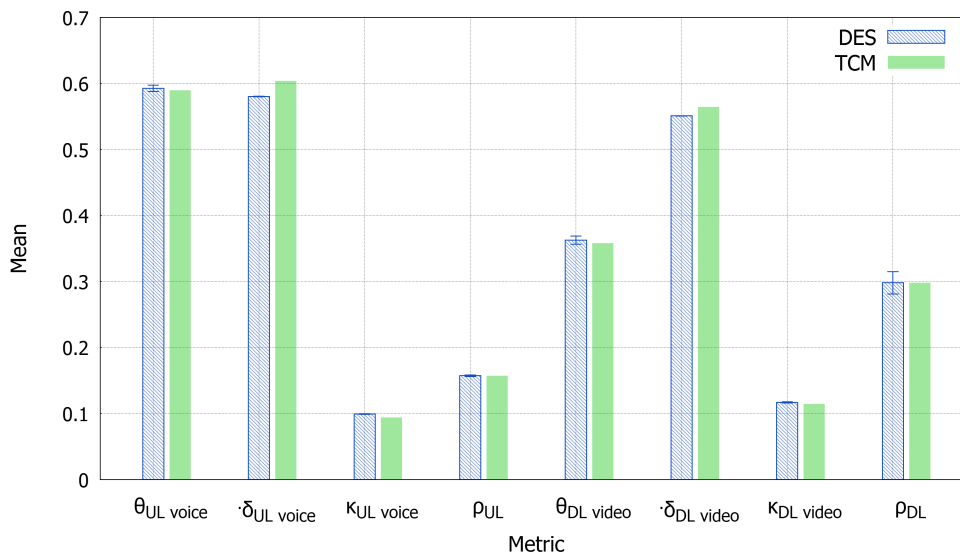
The scenario considers a single LTE cell with one eNB and a fixed number of UEs. Both uplink and downlink are modelled operating in FDD mode. Information for three different service types are transported along these links, namely web (best-effort service), voice (CBR, delay-sensitive) and downlink streaming video (VBR, delay-sensitive).

Both connection- and packet-level have dynamic behaviour. Therefore, connection admission control and scheduler modules are implemented. Such components are easily interchangeable, mimicking the flexibility feature of the FLAVIA architecture. By developing this model further, it should thus be possible to develop a FLAVIA modelling framework that closely resembles the entire FLAVIA architecture.

In the scenario, imperfect channel conditions were considered, allowing for adaptive modulation of data transmission as well.

### 4.2.3 Results

It has been shown, by comparing a DES model of the same scenario to that of the TCM model, that the TCM results are relatively very accurate, as shown in Figure 37. For normalized values of eight typical performance measures, the percentage error is less than 1 percent, with the exception of two particular cases where the delay is 2.36 and 1.36 percent, which is still significantly low.



*Figure 44: Performance metric mean values estimated by DES and TCM models*

Furthermore, it is also possible to estimate entire distributions, such as the uplink voice delay distribution, shown in Figure 38. The R-squared goodness-of-fit is very close to unity (all being greater than 0.99) for all the distributions estimated, when compared to the distribution estimated by the DES model. The R-squared value for the uplink voice delay fit shown in Figure 38 is 0.9908.



*Figure 45: Normalized uplink voice delay distribution estimated by DES and TCM models*

The time and space resource usages were recorded for both TCM and DES models. The TCM speedup gained with respect to DES speedup proved to be a run-time reduction of at least 14 times. The entire time-to-solution gained by using TCM, instead of using DES as the solution technique, was at least 550.
The TCM model generated no trace data (even though this would be possible but was not necessary). However, in the two experiments conducted, the DES model generated 138 and 620 gigabytes of trace data, respectively, which was necessary to compute the system performance.
These scenarios experimented with were lightly loaded in terms of the intensity of operations to be executed by the implemented model. However, it is apparent that all TCM gains increase as the intensity increases.
We therefore conclude that we have managed to significantly reduce the time and space resource usages to a practical level, while still achieving accurate performance estimates, by using the TCM solution technique instead of the usual DES technique.

## 4.3 System level simulation overview (NEC)

In FLAVIA WP 5, a system level simulator has been developed which is calibrated with 3GPP LTE RAN1 [47] and RAN2 for the study item HetNet Mobility Enhancements [48]. The system level simulator supports

- Shadow map implementation (large scale parameters)
- Implementation of mobility within LS parameter maps
- Implementation of HO logic
- Implementation of radio link failure (RLF) and handover failure (HOF) logic (incl. measurements such as RSRP/RSRQ, timers, in/out-of-sync tracking, logging, …)
- Scalable complexity (fast fading, scheduling, SINR computation etc. can be parameterized to adjust simulation speed accordingly)
- Multi-threading support
- Full HO parameterization (TTT, A3, …)

The simulator has been used to contribute to the study item small-cell enhancements by providing an overview on the impact of different handover preparation times on the handover failure performance. Within this study item, backhaul latency of up to 60ms for non-ideal backhaul is considered. Considering the round trip time as well as the required processing at the target eNB during a handover, the handover preparation time may be in the order of up to 150ms (60ms latency in each direction plus processing at the target eNB), which will have considerable impact on the handover performance.
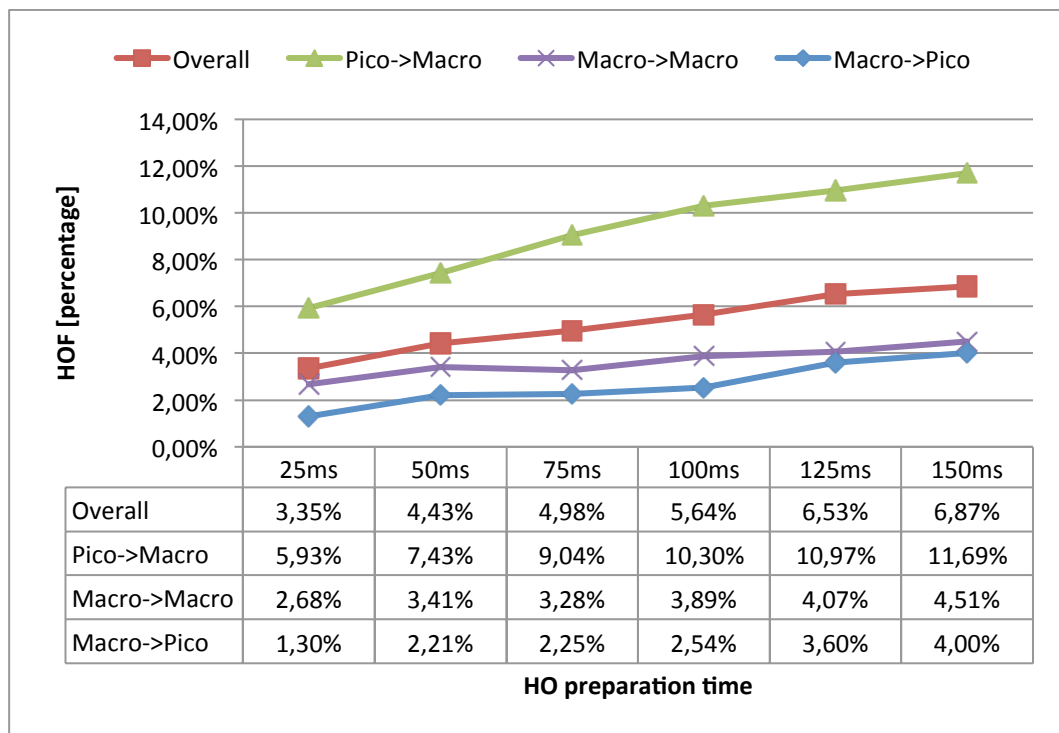
The reason for a HOF, i.e. bad link quality, may appear during the TTT interval, the HO preparation time, or the HO execution time. The TTT interval can be seen as a guard interval to protect from premature handovers and to avoid ping-pong effects. The HO preparation time cannot be reduced but is a parameter which depends on the link between base stations as well as the processing speed at both base stations. Finally, the HO execution time cannot be reduced because it consists of pre-defined steps to access the target cell.

While the TTT interval is a design parameter, the HO preparation time is an architectural parameter, which may not be changed. If the HO preparation time is rather long due to high backhaul latency or long processing delays, the link quality on the link between source cell and user terminal may become sufficiently bad to start T310. In most cases, this will also cause a HOF because the available time to recover the link is rather short. If the link does not recover, the HO command cannot be received and a HOF occurs. Therefore, the longer HO preparation time, the higher the probability to cause a HOF.

In order to evaluate these effects, the handover module as described in deliverable D3.1.2 has been implemented and tested. The corresponding performance results are shown in see Figure 39 and reported to 3GPP

RAN2#81. We can observer that the overall HOF rate and in particular the HOFs for Pico-to-Macro handovers suffer from an increased HO preparation time. Specifically, the pico-to-macro HOF rate increases from 7.43% (baseline assumption for [48]) to about 12% for 150ms (a pessimistic but still possible scenario in the case of non-ideal backhaul). This applies similarly to the overall HOF rate which increases from 4.4% to about 7% within the same value range.



| | 25ms | 50ms | 75ms | 100ms | 125ms | 150ms |
|---|---|---|---|---|---|---|
| Overall | 3,35% | 4,43% | 4,98% | 5,64% | 6,53% | 6,87% |
| Pico->Macro | 5,93% | 7,43% | 9,04% | 10,30% | 10,97% | 11,69% |
| Macro->Macro | 2,68% | 3,41% | 3,28% | 3,89% | 4,07% | 4,51% |
| Macro->Pico | 1,30% | 2,21% | 2,25% | 2,54% | 3,60% | 4,00% |

**HO preparation time**

***Figure 46 Handover failure performance depending on the handover preparation time***

# References

[1]  E. B. T. Israeli and O. Gurewitz, "Experimental assessment of power-save behavior of commercial ieee 802.16e network," *Arxiv,* 2013.

[2]  H. Weingarten, Y. Steinberg and S. Shamai, "The capacity region of the Gaussian multiple-input multiple-output broadcast channel," *IEEE transactions on information theory,* vol. 52, no. 9, pp. 3936--3964, 2006.

[3]  R. Knopp and P. Humblet, "Information capacity and power control in single-cell multiuser communications," in *IEEE International Conference on Communications*, Seattle, 1995.

[4]  T. Yoo and A. Goldsmith, "On the optimality of multiantenna broadcast scheduling using zero-forcing beamforming," *IEEE Journal on Selected Areas in Communications (JSAC),* vol. 24, no. 3, pp. 528--541, 2006.

[5]  R. Zakhour and S. Hanly, "Min-max fair coordinated beamforming via large system analysis," in *IEEE International Symposium on Information Theory Proceedings*, 2011.

[6]  C. Chen and L. Wang, "Enhancing coverage and capacity for multiuser {MIMO} systems by utilizing scheduling}," *IEEE Transactions on Wireless Communications,* vol. 5, no. 5, pp. 1148--1157, 2006.

[7]  K. Jagannathan, S. Borst, P. Whiting and E. Modiano, "Scheduling of multi-antenna broadcast systems with heterogeneous users," *IEEE Journal on Selected Areas in Communications,* vol. 25, no. 7, pp. 1424--1434, 2007.

[8]  L. Kleinrock and M. Scholl, "Packet Switching in Radio Channels:New Conflict-Free Multiple Access Schemes," *IEEE transactions on communications,* vol. 28, no. 7, pp. 1015--1029, 1980.

[9]  E. Telatar, "Capacity of Multi-antenna Gaussian Channels," *European transactions on telecommunications,* vol. 10, no. 6, pp. 585--595, 1999.

[10] V. Girko, "A refinement of the central limit theorem for random determinants," *Theory of Probability & Its Applications,* vol. 42, no. 1, pp. 121--129, 1998.

[11] P. Smith and M. Shafi, "On a Gaussian approximation to the capacity of wireless MIMO systems," in *IEEE International Conference on Communications (ICC)*, 2002.

[12] M. Chiani, M. Win and A. Zanella, "On the capacity of spatially correlated {MIMO} Rayleigh-fading channels," *IEEE transactions on information theory,* vol. 49, no. 10, pp. 2363--2371, 2003.

[13] L. De Haan and A. Ferreira, Extreme value theory: an introduction, New-

York: Springer Verlag, 2006.

[14] M. Leadbetter, Extremes and Related Properties of Random Sequences and Processes, New York: Springer-Verlag, 1983.

[15] S. Coles, An introduction to statistical modeling of extreme values, New York: Springer Verlag, 2001.

[16] W. Choi and J. Andrews, "The capacity gain from intercell scheduling in multi-antenna systems," *IEEE Transactions on Wireless Communications,* vol. 7, no. 2, pp. 714--725, 2008.

[17] O. Kallenberg, Random Measures, 3rd ed., N. Y. Academic Press, Ed., Akademie Verlag, Berlin, , 1983.

[18] R. Smith, "Extreme value analysis of environmental time series: an application," *Statistical Science,* vol. 4, no. 4, pp. 367--377, 1989.

[19] J. Galambos, J. Lechner and E. Simiu, "Extreme value theory and applications," in *Extreme Value Theory and Applications*, Gaithersburg, Maryland, 1994.

[20] X. a. B. R. Qin, "Exploiting multiuser diversity for medium access control in wireless," in *Twenty-Second Annual Joint Conference of the IEEE Computer (INFOCOM)*, 2003.

[21] A. Kochut, A. Vasan, A. Shankar and A. Agrawala, "Sniffing out the correct physical layer capture model in 802.11 b," in *The 12th IEEE International Conference on Network Protocols (ICNP)*, 2004.

[22] C. Reis, R. Mahajan, M. Rodrig, D. Wetherall and J. Zahorjan, "Measurement-based models of delivery and interference in static wireless networks," *ACM SIGCOMM Computer Communication Review,* vol. 36, no. 4, pp. 51--62, 2006.

[23] J. Lee, W. Kim, S. Lee, D. Jo, J. Ryu, T. Kwon and Y. Choi, "An experimental study on the capture effect in 802.11 a networks," in *The second ACM international workshop on Wireless network testbeds, experimental evaluation and characterization*, 2007.

[24] A. C. G. R. E. Biton and O. Gurewitz, "Distributed inter-cell interference mitigation via joint scheduling and power control under noise rise constraints," *Arxiv,* 2012.

[25] E. Coffman Jr., M. Garey and D. Johnson, "Approximation algorithms for bin packing: A survey," *Approximation Algorithms for NP-Hard Problems,* pp. 46-93, 1996.

[26] Wi-Fi Alliance Specification, *Wi-Fi Peer-to-Peer (P2P) Specification v1.1,* 2011.

[27] S. Sesia, I. Toufik and M. Baker, LTE - the UMTS long term evolution:

from theory to practice, Wiley, 2011.

[28] P. Rost, "Robust and Efficient Multi-Cell Cooperation under Imperfect CSI and Limited Backhaul," *accepted for publication in IEEE Transactions on Wireless Communications,* January 2013.

[29] O. S. G. Primolevo and U. Spagnolini, "Effects of imperfect channel state information on the capacity of broadcast OSDMA-MIMO systems," in *in IEEE Workshop on Signal Processing Advances in Wireless Communications, Lisbon, Portugal*, July 2004.

[30] T. Cover, "Broadcast channels," *IEEE Transactions on Information Theory,* Vols. 18, no. 1, p. 2–14, January 1972.

[31] T. Han and K. Kobayashi, "A new achievable rate region for the interference channel," *IEEE Transactions on Information Theory,* Vols. IT-27, no. 1, p. 49–60, January 1981.

[32] P. Rost, "Achievable Net-Rates in Multi-User OFDMA with Partial CSI and Finite Channel Coherence," in *IEEE Vehicular Technology Conference*, Quebec City, Canada, September 2012.

[33] P. Marsch and G. Fettweis, "On downlink network MIMO under a constrained backhaul and imperfect channel knowledge," in *in IEEE Global Communications Conference*, Hawaii, USA,, December 2009.

[34] P. Marsch and G. Fettweis, "Uplink CoMP under a constrained backhaul and imperfect channel knowledge," *IEEE Transactions on Wireless Communications,* Vols. 10, no. 6, p. 1730 – 1742, June 2011.

[35] ITU-R, "Report ITU-R M.2135-1, Guidelines for evaluation of radio interface technologies for IMT-Advanced," 2009.

[36] R. Ahlswede, N. Cai, S.-Y. R. Li and R. W. Yeung, "Network information flow," *IEEE Trans. Inform. Theory,* vol. 46, no. 4, pp. 1204-1216, July 2000.

[37] P. Chou, Y. Wu and K. Jain, "Practical network coding," in *Proceedings of the Annual Allerton Conference on Communication Control and Computing*, 2003.

[38] S. Katti, D. Katabi, W. Hu, H. Rahul and M. Medard, "The importance of being opportunistic: Practical network coding for wireless environments," in *Proc. 43rd Annual Allerton Conference on Communication, Control, and Computing*, 2005.

[39] S. Katti, H. Rahul, W. Hu, D. Katabi, M. M'edard and J. Crowcroft, "XORs in the air: practical wireless network coding," *IEEE/ACM Transactions on Networking (TON),* vol. 16, no. 3, pp. 497-510, 2008.

[40] S. Jaggi, P. Sanders, P. Chou, M. Effros, S. Egner, K. Jain and L.

Tolhuizen, "Polynomial time algorithms for multicast network code construction," *IEEE Transactions on Information Theory,* vol. 51, no. 6, pp. 1973-1982, 2005.

[41] H. Shojania and B. Li, "Random network coding on the iPhone: fact or fiction?," in *Proceedings of the 18th international workshop on Network and operating systems support for digital audio and video*, 2009.

[42] F. Fitzek, M. Pedersen, J. Heide and M. M'edard, "Network coding: applications and implementations on mobile devices," in *Proceedings of the 5th ACM workshop on Performance monitoring and measurement of heterogeneous wireless and wired networks*, 2010.

[43] S. Chieochan and E. Hossain, "Network Coding for Unicast in a WiFi Hotspot: Promises, Challenges, and Testbed Implementation," *Computer Networks,* 2012.

[44] E. B. J. K. Asaf Cohen and O. Gurewitz, "Coded Unicast Downstream Traffic in a Wireless Network: Analysis and WiFi Implementation," *EURASIP Journal on Advances in Signal Processing, accepted.*

[45] F. Berggren and R. Jantti, "Asymptotically fair scheduling on fading channels," in *in IEEE Veh. Techn. Conf.*, Vancouver, Canada, Sep 2002.

[46] J. Choi and S. Bahk, "Cell-throughput analysis of the proportional fair scheduler in the single-cell environment," *IEEE Transactions on Vehicular Technology,* Vols. 56, no. 2, p. 766–778, March 2007.

[47] 3rd Generation Partnership Project, "TR 36.814, Further advancements for E-UTRA physical layer aspects," October 2010.

[48] 3rd Generation Partnership Project, "TS 36.839, Mobility enhancements in heterogeneous networks," September 2012.

[49] J. Zander, "Performance of optimum transmitter power control in cellular radio systems," *IEEE Transactions on Vehicular Technology,* vol. 41, no. 1, pp. 57-62, 1992.

[50] R. D. Yates, "A framework for uplink power control in cellular radio systems," *IEEE Journal on Selected Areas in Communications,* vol. 13, no. 7, pp. 1341-1347, 1995.

[51] V. G. Subramanian, R. A. Berry and R. Agrawal, "Joint Scheduling and Resource Allocation in CDMA Systems," *IEEE Transactions on Information Theory,* vol. 56, no. 5, pp. 2416-2432, 2010.

[52] R1-073224, "Way Forward on Power Control of PUSCH," 3GPP TSG- RAN WG1 49-bis, 2007.

[53] Y.-H. Lin and R. L. Cruz, "Power control and scheduling for interfering links," in *Proc. IEEE Information Theory Workshop*, 2004.

[54] J. Huang, V. G. Subramanian, R. Agrawal and R. Berry, "Joint scheduling and resource allocation in uplink OFDM systems for broadband wireless access networks," *IEEE Journal on Selected Areas in Communications,* vol. 27, no. 2, pp. 226-234, 2009.

[55] H. Holmax and A. Toskala, WCDMA for UMTS: Radio Access for Third Generation Mobile Communications, John Willey, 2004.

[56] T. ElBatt and A. Ephremides, "Joint scheduling and power control for wireless ad-hoc networks," 2002.

[57] R. Cruz and A. Santhanam, "Optimal routing, link scheduling and power control in multihop wireless networks," 2003.

[58] M. Chiang, P. Hande, T. Lan and C. W. Tan, "Power Control in Wireless Cellular Networks," *Found. Trends Netw.,* vol. 2, no. 4, pp. 381-533, 2008.

[59] R. Cendrillon, J. Huang, M. Chiang and M. Moonen, "Autonomous Spectrum Balancing for Digital Subscriber Lines," *Signal Processing, IEEE Transactions on,* vol. 55, no. 8, pp. 4241-4257, 2007.

[60] C. U. Castellanos, D. L. Villa, C. Rosa, K. I. Pedersen, F. D. Calabrese, P.-H. Michaelsen and J. Michel, "Performance of Uplink Fractional Power Control in UTRAN LTE," in *Proc. IEEE Vehicular Technology Conf. VTC Spring 2008*, 2008.

[61] R. Agrawal, V. Subramanian and R. Berry, "Joint Scheduling and Resource Allocation in CDMA Systems," 2004.