



**A Network of Excellence forging the  
Multilingual Europe Technology Alliance**

---

## **Optimal Choice Selection in MT**

---

**Authors:** Christian Federmann, Maite Melero, Josef van Genabith

**Dissemination Level:** Public

**Date:** February 1, 2012



Grant agreement no.	249119
Project acronym	T4ME Net (META-NET)
Project full title	Technologies for the Multilingual European Information Society
Funding scheme	Network of Excellence
Coordinator	Prof. Hans Uszkoreit (DFKI)
Start date, duration	1 February 2010, 36 months
Distribution	Public
Contractual date of delivery	31 January 2012
Actual date of delivery	6 February 2012
Deliverable number	D2.2
Deliverable title	Optimal Choice Selection in MT
Type	Report
Status and version	Draft
Number of pages	16
Contributing partners	BM, DCU, DFKI
WP leader	DCU
Task leader	DFKI
Authors	Christian Federmann, Maite Melero, Josef van Genabith
EC project officer	Hanna Klimek
The partners in META-NET are:	Deutsches Forschungszentrum für Künstliche Intelligenz (DFKI), Germany
	Barcelona Media (BM), Spain
	Consiglio Nazionale Ricerche – Istituto di Linguistica Computazionale “Antonio Zampolli” (CNR), Italy
	Institute for Language and Speech Processing, R.C. “Athena” (ILSP), Greece
	Charles University in Prague (CUP), Czech Republic
	Centre National de la Recherche Scientifique – Laboratoire d’Informatique pour la Mécanique et les Sciences de l’Ingénieur (CNRS), France
	Universiteit Utrecht (UU), Netherlands
	Aalto University (AALTO), Finland
	Fondazione Bruno Kessler (FBK), Italy
	Dublin City University (DCU), Ireland
	Rheinisch-Westfälische Technische Hochschule Aachen (RWTH), Germany
	Jozef Stefan Institute (JSI), Slovenia
Evaluations and Language Resources Distribution Agency (ELDA), France	

For copies of reports, updates on project activities and other META-NET-related information, contact:

DFKI GmbH  
 META-NET  
 Dr. Georg Rehm  
 Alt-Moabit 91c  
 10559 Berlin, Germany

office@meta-net.eu  
 Phone: +49 30 23895-1833  
 Fax: +49 30 23895-1810

Copies of reports and other material can also be accessed via <http://www.meta-net.eu>

© 2012, The Individual Authors

No part of this document may be reproduced or transmitted in any form, or by any means, electronic or mechanical, including photocopy, recording, or any information storage and retrieval system, without permission from the copyright owner.

## **Table of Contents**

1	Executive Summary .....	4
2	Optimal Choice Selection in MT .....	5

## 1 Executive Summary

Machine translation (MT) is an active field of research with many competing paradigms to solve fundamental translation problems. In recent years, an important focus for research has been investigating how hybrid MT as well as MT combination systems including several translation engines can be designed and implemented so that the resulting translations give an improvement over the individual translations. One of the main objectives in our research within the T4ME project is to provide a systematic investigation and exploration of optimal choices in Hybrid Machine Translation.

As a first step towards achieving this objective we have developed a parallel corpus with source data and the output of a number of MT systems, representing carefully selected MT paradigms, annotated with meta-information, capturing aspects of the translation process performed by the different MT systems. Including detailed and heterogeneous system specific information as metadata in the translation output (rather than just providing strings) is intended to provide rich features for machine learning methods to optimise system combination.

As a second step, we have organised a shared task in which participants were requested to build Hybrid/System Combination systems by combining the output of several systems of different types, using the annotated corpus as input. The main focus of the shared task is trying to answer the following question: *Could Hybrid MT algorithms or System Combination techniques benefit from extra information (linguistically motivated, decoding and runtime) from the different systems involved?*

In this deliverable, we describe the annotated corpus we have created and then provide an overview on the participating systems from the shared task as well as a discussion of the results. The deliverable has been submitted for publication in “The Prague Bulletin of Mathematical Linguistics” journal and is attached on the following pages.

# 1 Introduction

Machine translation is an active field of research with many competing paradigms to tackle core translation problems. In recent years, an important focus for research has been investigating how hybrid machine translation engines as well as combination systems including several translation engines can be designed and implemented so that the resulting translations give an improvement over the component parts.

One of the main objectives in our research within the T4ME project is to provide a systematic investigation and exploration of optimal choices in Hybrid Machine Translation supporting Hybrid MT design using sophisticated machine-learning technologies. As a first step towards achieving this objective we have developed a parallel corpus with source data and the output of a number of MT systems, representing carefully selected MT paradigms, annotated with meta-information, capturing aspects of the translation process performed by the different machine translation systems. Including detailed and heterogeneous system specific information as metadata in the translation output (rather than just providing strings) is intended to provide rich features for machine learning methods to optimise combination in hybrid machine translation.

This first version of the corpus is available under <http://www.dfki.de/ml4hmt/> and comprises annotated outputs of five machine translation systems, namely Joshua, Lucy, Metis, Apertium, MaTrEx. The language pairs available from the corpus are: English↔German, English↔Spanish, English↔Czech (all in both directions).

In this paper, we describe the annotated corpus we have created—including the data used to obtain the sample corpus (Section 2), the translation engines applied when building the corpus (Section 3), and the format of the corpus (Section 4)—and then provide an overview on the challenge (Section 5) and give descriptions of the participating systems from the shared task (Section 6). This includes a comparison to a state-of-the-art system combination system. Using automated metrics scores and results from manual evaluation, we discuss the performance of the various systems and their different implementations (Section 7). One interesting result from the shared task is the fact that we observed different systems winning according to the automated metrics and according to the manual evaluation. We conclude by summarising our research results and give an outlook to future work (Section 8).

## 2 Data

### 2.1 Annotation Data

As a source of the data to be included and annotated in the corpus we decided to use the WMT 2008 news test set, which is a set of 2,051 sentences from the news domain translated to all languages of our interest (English, Spanish, German, Czech) and also some others (French, Hungarian). This test set was provided by the organizers of the Third Workshop on Machine Translation (WMT) in 2008 as test data for the shared translation task. The test set from WMT 2010 has been reserved for final tests in the future (if needed).

### 2.2 Training Data-Driven Systems

Some of the MT systems used in this work are data-driven (Joshua and MaTrEx). They require parallel data for translation phrase pair extraction, monolingual data for language modeling, and parallel development data for tuning of system parameters. Originally we intended to use the Europarl corpus [11] for training purposes, but since this widely used parallel corpus did not include Czech at the time of the development, we have opted for the Acquis [23] and News Commentary parallel corpora instead.

### 2.3 JRC-Acquis Multilingual Parallel Corpus

The JRC-Acquis Multilingual Parallel Corpus is an “approximation” of the Acquis Communautaire, the total body of European Union (EU) law applicable in the the EU Member States. It comprises documents that were available in at least in ten of the twenty EU-25 languages (official languages in EU before 2007)

and that additionally existed in at least three of the nine languages that became official languages with the Enlargement of the EU in 2004 (i.e. Czech, Hungarian, Slovak, etc.).

## 2.4 WMT News Commentary Parallel Corpus

The WMT News Commentary Parallel Corpus contains news and commentaries from the Project Syndicate and is provided as training data for the series of WMT translation shared tasks (<http://www.statmt.org/>). Version 10 was released in 2010 and is available in English, French, Spanish, German, and Czech.

## 2.5 Development Data

The development data sets (for tuning the statistical systems) were taken from the WMT 2008 test set, which consists of 2,007 sentences from the news-commentary domain available in English, French, Spanish, German, and Czech.

# 3 System Descriptions

## 3.1 Joshua

**Description** Joshua, system *t1*, [15] is an open-source toolkit for statistical machine translation, providing a full implementation of state-of-the-art techniques making use of synchronous context free grammars (SCFGs). The decoding process features algorithms such as chart-parsing, n-gram language model integration, beam-and cube-pruning and k-best extraction, while training includes suffix-array grammar extraction and minimum error rate training.

**Annotation** In our metadata annotations, we provide the output of the decoding process given the “test set”, as processed by Joshua (SVN revision 1778). The annotation set contains the globally applied feature weights and for each translated sentence: the full output of the produced translation with the highest total score (among the n-best candidates), the language model and translation table scores, the scores from derivation of the sentence (phrase scores) and merging/pruning statistics of the search process. Each translated sentence, represented by a hierarchical phrase, contains zero or more tokens and points to zero or more child phrases. Finally, the word-alignment of each phrase to the source text, using word indices, is available.

## 3.2 Lucy

**Description** The Lucy RBMT system, system *t2*, [1] uses a sophisticated RBMT transfer approach with a long research history. It employs a complex lexicon database and grammars to transform a source into a target language representation. The translation of a sentence is carried out in three major phases: analysis, transfer, and generation.

**Annotation** In addition to the translated target text Lucy provides information about the tree structures that have been created in the three translation phases and which have been used to generate the final translation of the source text. Inside these trees, information about POS, phrases, word lemma information, and word/phrase alignment can be found. In our metadata annotations, we provide a “flattened” representation of the trees. For each token, annotation may contain allomorphs, canonical representations, linguistic categories, or surface string.

### 3.3 Metis

**Description** The Metis system, system *t3*, [25] achieves corpus-based translation on the basis of a monolingual target corpus and a bilingual dictionary only. The bilingual dictionary functions as a flat translation model that provides  $n$  translations for each source word. The most probable translation given the context is then selected by consulting the statistical models built from the target language corpus.

**Annotation** Meta-data information for Metis is extracted from the set of final translations ranked by the Metis search engine. For each translation we obtain the score computed during the search process, together with some linguistic information. The basic linguistic information provided is: lemma, POS tag, and morphological features. Morphological features are grouped under one feature derived from the source token, including gender, number, tense, etc.

### 3.4 Apertium

**Description** Apertium, system *t4*, [20] originated as one of the machine translation engines in the project OpenTrad, funded by the Spanish government. Apertium is a shallow-transfer machine translation system, which uses finite state transducers for all of its lexical transformations, and hidden Markov models for POS tagging or word category disambiguation. Constraint Grammar taggers are also used for some language pairs (e.g., Breton-French).

**Annotation** We use Apertium version 3.2. Our metadata annotation includes tags, lemmas and syntactic information. We have used the following commands (in English-to-Spanish): `en-es-chunker` (for syntax information), `en-es-postchunk` (for tags and lemmas) and `en-es` (for the translation).

### 3.5 MaTrEx

**Description** The MaTrEx machine translation system, system *t5*, [19] is a combination-based multi-engine architecture developed at Dublin City University exploiting aspects of both the Example-based Machine Translation (EBMT) and Statistical Machine Translation (SMT) paradigms. The architecture includes various individual systems: phrase-based, example-based, hierarchical phrase-based, and tree-based MT. For the corpus data produced here we use the standard MOSES PB-SMT system as integrated into MaTrEx.

**Annotation** Sentence translations provided by MaTrEx in this work were obtained using the MOSES PB-SMT system decomposing the source side to phrases (n-grams), finding their translation and composing them to a target language sentence which has the highest score according the PB-SMT model. Meta-data annotations for each sentence translated by MaTrEx include scores from each model and is decomposed into phrases each provided with two additional scores: translation probability and future cost estimate. Additionally information about unknown words is also included.

## 4 Corpus Description

We have developed a new dedicated format derived from XLIFF (XML Localisation Interchange File Format) to represent and store the corpus data. XLIFF is an XML-based format created to standardize localization. It was standardized by OASIS in 2002 and its current specification is v1.2 (<http://docs.oasis-open.org/xliff/xliff-core/xliff-core.html>).

An XLIFF document is composed of one or more `<file>` elements, each corresponding to an original file or source. Each `<file>` element contains the source of the data to be localized and the corresponding localized (translated) data for one locale only. The localizable texts are stored in `<trans-unit>` elements each having a `<source>` element to store the source text and a `<target>` (not mandatory) element to store the translation.

ranked 1st	WER[%]	sBLEU[%]	BLEU
system <i>t1</i>	26.44	14.73	12.80
system <i>t2</i>	37.56	38.63	14.94
system <i>t3</i>	5.85	5.85	8.29
system <i>t4</i>	16.00	23.41	13.34
system <i>t5</i>	58.54	29.95	14.47
sBLEU-combined			18.95
WER-combined			17.62

Table 1: Preliminary investigation for the usability of the corpus for Hybrid MT

We introduced new elements into the basic XLIFF format (in the "metanet" namespace) allowing a wide variety of metadata annotation of the translated texts by different MT systems (tools). The tool information is included in the `<tool>` element appearing in the header of the file. Each tool can have several parameters (model weights) which are described in the `<metanet:weight>`.

Annotation is stored in `<alt-trans>` element within the `<trans-unit>` elements. The `<source>` and `<target>` elements in the `<trans-unit>` elements refer to the source sentence and its reference translation, respectively. The `<source>` and `<target>` elements in the `<alt-trans>` elements specifies the input and output of a particular MT system (tool). Tool-specific scores assigned to the translated sentence are listed in the `<metanet:scores>` element and the derivation of the translation is specified in the `<metanet:derivation>` element. Its content is tool-specific.

The full format specification is available as an XML schema. An example annotation is depicted in Figure 1.

#### 4.1 Usability of the Corpus

There are already strong indications that the corpus is indeed useful for the purpose it has been designed for, i.e. to provide data and in particular rich metadata for advancing Hybrid Machine Translation.

At first, we compare the performance of the contributing systems (*t1–t5*) on the sentence level, using two popular metrics, WER and smoothed-BLEU. In columns 2 and 3 of Table 1, the percentage of the cases that a system gave the best translation for a sentence according the two sentence-level metrics is shown<sup>1</sup>. This indicates that the systems included in the corpus perform complementary to each other. Column 3 indicates the overall BLEU score of each system, whereas the last row indicates what the optimal BLEU performance would be, if a system would be able to choose the best sentence of each system. This indicates the possibilities of improvement by a sentence-selection approach given the corpus. We believe that even higher performance would be possible using more sophisticated system combination methods.

## 5 Challenge Description

The “Shared Task on Applying Machine Learning techniques to optimise the division of labour in Hybrid MT” is an effort to trigger systematic investigation on improving state-of-the-art Hybrid MT, using advanced machine-learning (ML) methodologies. Participants of the challenge are requested to build Hybrid/System Combination systems by combining the output of several MT systems of different types and with very heterogeneous types of metadata information, as provided by the organizers. The main focus of the shared task is trying to answer the following question:

---

<sup>1</sup>Measured over the dev set, ties were allowed.

*Could Hybrid Machine Translation algorithms or System Combination techniques benefit from extra information—such as linguistic or linguistically motivated features, decoding parameters, or runtime annotation—from the different systems involved?*

The participants are given a bilingual development set, aligned at a sentence level. For each sentence, the corresponding *bilingual data set* contains:

- the source sentence,
- the target (reference) sentence, and
- the corresponding multiple output translations from five different systems, based on various machine translation approaches.

## 5.1 Development and Test Sets

We decided to use the WMT 2008 [6] news test set as a source for the annotated corpus. This is a set of 2,051 sentences from the news domain translated to several languages, including English and Spanish but also others. The data was provided by the organizers of the Third Workshop on Machine Translation (WMT) in 2008. This data set was split into our own development set (containing 1,025 sentence pairs) and test set (containing 1,026 sentence pairs).

## 6 Participating Systems

### 6.1 DCU

The system described in [17] presents a system combination module in the MT system MaTrEx (Machine Translation using Examples) developed at Dublin City University. A system combination module deployed by them achieved an improvement of 2.16 BLEU [18] points absolute and 9.2% relative compared to the best single system, which did not use any external language resources. Their system is based on system combination techniques which use a confusion network on top of a Minimum Bayes Risk (MBR) decoder [13].

One interesting, novel point in their submission is that for the given single best translation outputs, they tried to identify which inputs they will consider for the system combination, possibly discarding the worst performing system(s) from the combination input. As a result of this selection process, their BLEU score, from the combination of the four single best systems, achieved 0.48 BLEU points absolute higher than the combination of the five single best systems.

### 6.2 DFKI-A

A system combination approach with a sentence ranking component is presented in [2]. The paper reports on a pilot study on a Hybrid Machine Translation that takes advantage of multilateral system-specific metadata provided as part of the shared task. The proposed solution offers a machine learning approach, resulting in a selection mechanism able to learn and rank and select systems' translation output on the complete sentence level, based on their respective quality.

For training, due to the lack of human annotations, word-level Levenshtein distance has been used as a (minimal) quality indicator, whereas a rich set of sentence features was extracted and selected from the dataset. Three classification algorithms (Naive Bayes, SVM and Linear Regression) were trained and tested on pairwise featured sentence comparisons. The approaches yielded high correlation with original rankings ( $\tau = 0.52$ ) and selected the best translation on up to 54% of the cases.

### 6.3 DFKI-B

The authors of [9] report on experiments that are focused on word substitution using syntactic knowledge. From the data provided by the workshop organisers, they choose one system to provide the “translation backbone”. The Lucy MT system was suited best for this task, as it offers parse trees of both the source and target side, which allows the authors to identify interesting phrases, such as noun phrases, in the source and replace them in the target language output. The remaining four systems are mined for alternate translations on the word level that are potentially substituted into the aforementioned template translation if the system finds enough evidence that the candidate translation is better. Each of these substitution candidates is evaluated concerning a number of factors:

- the part-of-speech of the original translation must match the candidate fragment.
- Additionally they may consider the 1-left and 1-right context.
- Besides the part-of-speech, all translations plus their context are scored with a language model trained on EuroParl.
- Additionally, the different systems may turn up with the same translation, in that case the authors select the candidate with the highest count (“majority voting”).

The authors reported improvements in terms of BLEU score when comparing to the translations from the Lucy RBMT system.

### 6.4 LIUM

Barrault and Lambert submitted results from applying the open-source MANY [4] system on our data set. The MANY system can be decomposed into two main modules.

1. The first one is the alignment module which actually is a modified version of TER<sub>p</sub> [22]. Its role is to incrementally align the hypotheses against a backbone in order to create a confusion network. Each hypothesis acts as backbone, yielding each the corresponding confusion network. Those confusion networks are then connected together to create a lattice.
2. The second module is the decoder. This decoder is based on the token pass algorithm and it accepts as input the lattice previously created. The costs computed in the decoder can be expressed as a weighted sum of the logarithm of feature functions. The following features are considered in decoding:
  - the language model probability, given by a 4-gram language model,
  - a word penalty, which depends on the number of words in the hypothesis,
  - a null-arc penalty, which depends on the number of null arcs crossed in the lattice to obtain the hypothesis, and
  - the system weights: each word receives a weight corresponding to the sum of the weights of all systems which proposed it.

## 7 Evaluation Results

To evaluate the performance of the participating systems, we computed automated scores, namely BLEU, NIST, METEOR [3], PER, Word error rate (WER) and Translation Error Rate (TER) and also performed an extensive, manual evaluation with 3 annotators ranking system combination results for a total of 904 sentences.

## 7.1 Automated Scores

Results from running automated scoring tools on the submitted translations are reported in Table 2. The overall best value for each of the scoring metrics is print in **bold face**. Table 3 presents automated metric scores for the individual systems in the ML4HMT corpus, also computed on the test set. These scores give an indicative baseline for comparison with the system combination results.

## 7.2 Manual Ranking

The manual evaluation is undertaken using the Appraise [8] system; a screenshot of the evaluation interface is shown in Figure 2. Users are shown a reference sentence and the translation output from all four participating systems and have to decide on a ranking in *best-to-worst order*. Table 4 shows the average ranks per system from the manual evaluation, again the best value per column is printed in **bold face**. Table 5 gives the statistical mode per system which is the value that occurs most frequently in a data set.

## 7.3 Inter-annotator Agreement

Next to computing the average rank per system and the statistical mode, we follow [7] and compute  $\kappa$  scores to estimate the inter-annotator agreement. In our manual evaluation campaign, we had  $n = 3$  annotators so computing basic, pairwise annotator agreement is not sufficient—instead, we apply [10] who extends [21] for computing inter-annotator agreement for  $n > 2$ .

**Annotation Setup** As we have mentioned before, we had  $n = 3$  annotators assign ranks to our four participating systems. As ties were not allowed, this means there exist  $4! = 24$  possible rankings per sentence (e.g., *ABCD*, *ABDC*, etc.)<sup>2</sup>. In a second evaluation scenario, we only collected the *1-best* ranking system per sentence, resulting in a total of four categories (A: “*system A ranked 1st*”, etc.). In this second scenario, we can expect a higher annotator agreement due to the reduced number categories. Overall, we collected 904 sentences with an overlap of  $N = 146$  sentences for which all annotators assigned ranks.

**Scott’s  $\pi$**  allows to measure the pairwise annotator agreement for a classification task. It is defined as

$$\pi = \frac{P(A) - P(E)}{1 - P(E)} \quad (1)$$

where  $P(A)$  represents the fraction of rankings on which the annotators agree, and  $P(E)$  is the probability that they agree by chance. Table 6 lists the pairwise agreement of annotators for all four participating systems. Assuming  $P(E) = 0.5$  we obtain an overall agreement  $\pi$  score of

$$\pi = \frac{0.673 - 0.5}{1 - 0.5} = 0.346 \quad (2)$$

which can be interpreted as *fair agreement* following [14]. WMT shared tasks have shown this level of agreement is common for language pairs, where the performance of all systems is rather close to each other, which in our case is indicated by the small difference measured by automatic metrics on the test set (Table 2). The lack of ties, in this case might have meant an extra reason for disagreement, as annotators were forced to distinguish a quality difference which otherwise might have been annotated as “equal”. We have decided to compute Scott’s  $\pi$  scores to be comparable to WMT11 [5].

**Fleiss  $\kappa$**  Next to the  $\pi$  scores, there also exists the so-called  $\kappa$  score. Its basic equation is strikingly similar to (1)

$$\kappa = \frac{P(A) - P(E)}{1 - P(E)} \quad (3)$$

---

<sup>2</sup>Given this huge number of possible categories, we were already expecting resulting  $\kappa$  scores to be low.

System	BLEU	NIST	METEOR	PER	WER	TER
DCU	<b>25.32</b>	<b>6.74</b>	56.82	<b>60.43</b>	<b>45.24</b>	<b>0.65</b>
DFKI-A	23.54	6.59	54.30	61.31	46.13	0.67
DFKI-B	23.36	6.31	<b>57.41</b>	65.22	50.09	0.70
LIUM	24.96	6.64	55.77	61.23	46.17	0.65

Table 2: Automated scores for ML4HMT test set.

System	BLEU	NIST	METEOR	PER	WER
Joshua	19.68	6.39	50.22	47.31	62.37
Lucy	<b>23.37</b>	6.38	<b>57.32</b>	49.23	64.78
Metis	12.62	4.56	40.73	63.05	77.62
Apertium	22.30	6.21	55.45	50.21	64.91
MaTrEx	23.15	<b>6.71</b>	54.13	<b>45.19</b>	<b>60.66</b>

Table 3: Automated scores for baseline systems on ML4HMT test set.

System	Annotator #1	Annotator #2	Annotator #3	Overall
DCU	2.44	2.61	2.51	2.52
DFKI-A	2.50	2.47	2.48	2.48
DFKI-B	<b>2.06</b>	<b>2.13</b>	<b>1.97</b>	<b>2.05</b>
LIUM	2.89	2.79	2.93	2.87

Table 4: Average rank per system per annotator from manual ranking of 904 (overlap=146) translations.

System	Ranked 1st	Ranked 2nd	Ranked 3rd	Ranked 4th	Mode
DCU	62	79	<b>97</b>	62	3rd
DFKI-A	73	65	<b>82</b>	80	3rd
DFKI-B	<b>127</b>	84	47	42	1st
LIUM	38	72	74	<b>116</b>	4th

Table 5: Statistical mode per system from manual ranking of 904 (overlap=146) translations.

Systems	$\pi$ -Score	Systems	$\pi$ -Score	Annotators	$\pi$ -Score
DCU, DFKI-A	0.296	DCU, DFKI-B	0.352	#1,#2	0.331
DCU, LIUM	0.250	DFKI-A, DFKI-B	0.389	#1,#3	0.338
DFKI-A, LIUM	0.352	DFKI-B, LIUM	0.435	#2,#3	0.347

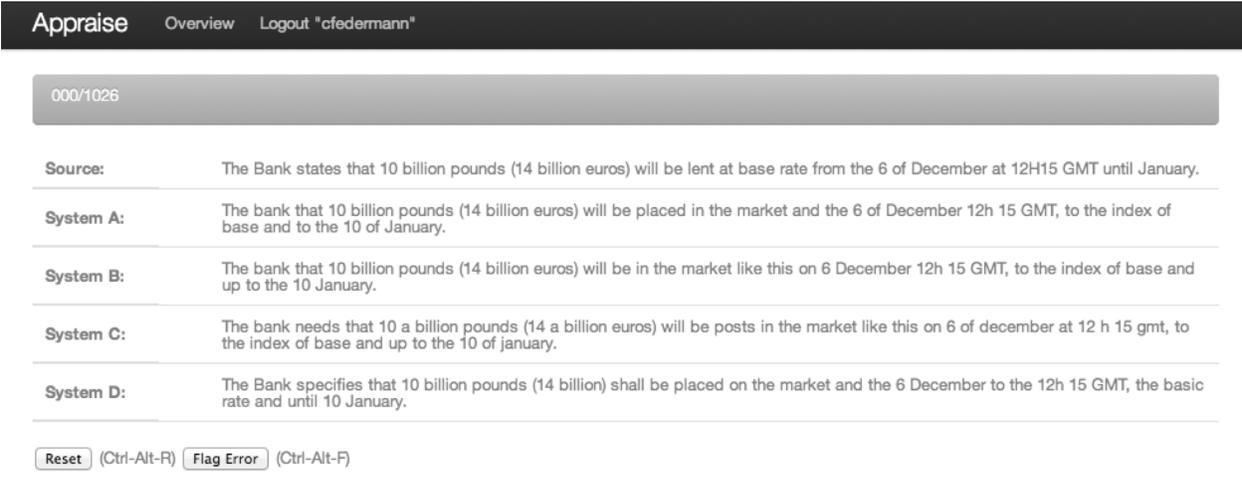
Table 6: Pairwise agreement (using Scott's  $\pi$ ) for all pairs of systems/annotators.

```

<trans-unit id="s71">
  <source xml:lang="es">El paciente fue aislado.</source>
  <target xml:lang="en">The patient was isolated.</target>
  <alt-trans rank="1" tool-id="t3">
    <source xml:lang="es">El paciente fue aislado.</source>
    <target xml:lang="en">The paciente was isolated .</target>
    <metanet:scores>
      <metanet:score type="total" value="-60.4375047559049"/>
    </metanet:scores>
    <metanet:derivation id="s71_t3_r1_d1">
      <metanet:phrase id="s71_t3_r1_d1_p1">
        <metanet:string>The</metanet:string>
        <metanet:annotation type="lemma" value="the"/>
        <metanet:annotation type="pos" value="AT0"/>
        <metanet:annotation type="morph_feat" value=":m:sg:"/>
        <metanet:alignment from="0" to="0"/>
      </metanet:phrase>
    </metanet:derivation>
  </alt-trans>
</trans-unit>

```

Figure 1: Example of annotation from the ML4HMT corpus.



This is the GitHub version of the Appraise evaluation system. Some rights reserved.

Figure 2: Screenshot of the Appraise interface for human evaluation.

with the main difference being the  $\kappa$  score’s support for  $n > 2$  annotators. We compute  $\kappa$  for two configurations. Both are based on  $n = 3$  annotators and  $N = 146$  sentences. They differ in the number of categories that a sentence can be assigned to ( $k$ ).

1. *complete* scenario:  $k = 24$  categories. For this, we obtained a  $\kappa$  score of

$$\kappa_{complete} = \frac{0.1 - 0.054}{1 - 0.054} = 0.049 \quad (4)$$

2. *1-best* scenario:  $k = 4$  categories. Here,  $\kappa$  improved to

$$\kappa_{1-best} = \frac{0.368 - 0.302}{1 - 0.302} = 0.093 \quad (5)$$

It seems that the large number of categories of the *complete* scenario has indeed had an effect on the resulting  $\kappa_{complete}$  score. This is a rather expected outcome, still we report the  $\kappa$  scores for future reference. The *1-best* scenario supports an improved  $\kappa_{1-best}$  score but does not reach the level of agreement observed for the  $\pi$  score.

It seems that DFKI-B was underestimated by BLEU scores, potentially due to its rule-based characteristics. This is a possible reason for the relatively higher inter-annotator agreement when compared with other systems. Also, DCU and LIUM may have low inter-annotator agreement as their background is similar. It is worth noting that METEOR was the only automated metric correlating with results from the manual evaluation.

Due to the fact that  $\kappa$  is not really defined for *ordinal data* (such as rankings in our case), we will investigate other measures for inter-annotator agreement. It might be a worthwhile idea to compute  $\alpha$  scores, as described in [12]. Given the average rank information, statistical mode,  $\pi$  and  $\kappa$  scores, we still think that we have derived enough information from our manual evaluation to support for future discussion.

## 8 Conclusion

We have developed an Annotated Hybrid Sample MT Corpus which is a set of 2,051 sentences translated by five different MT systems<sup>3</sup> (Joshua, Lucy, Metis, Apertium, and MaTrEx) in six translation directions (Czech→English, German→English, Spanish→English, English→Czech, English→German, and English→Spanish) and annotated with various information provided by the MT systems.

Using this resource we have launched the Shared Task on Applying Machine Learning techniques to optimise the division of labour in Hybrid MT (ML4HMT-2011), asking participants to create combined, hybrid translations using machine learning algorithms or other, novel ideas for making best use of the provided ML4HMT corpus data. The language pair for the Shared Task was Spanish→English.

Four participating combination systems, each following a different solution strategy, have been submitted to the shared task. We computed automated metric scores and conducted an extensive manual evaluation campaign to assess the quality of the hybrid translations. Interestingly, the system winning nearly all the automatic scores (DCU) only reached a third place in the manual evaluation. Vice versa, the winning system according to manual rankings (DFKI-B) ranked last place in the automatic metric scores based evaluation, with only one automated metric choosing it as winning system. This clearly indicates that more systematic investigation of hybrid system combination approaches, both on a system level and with regards to the evaluation of such systems’ output, needs to be undertaken.

We have learned from the participants that our ML4HMT corpus is too heterogeneous to be used easily in system combination approaches; hence we will work on an updated version for the next edition of this shared task. Also, we will further focus on the integration of advanced machine learning techniques as these are expected to support better exploitation of our corpus’ data properties.

---

<sup>3</sup>Not all systems available for all language pairs.

## Acknowledgments

This work has been funded under the Seventh Framework Programme for Research and Technological Development of the European Commission through the T4ME contract (grant agreement no.: 249119).

## References

- [1] Juan A. Alonso and Gregor Thurmair. The Compendium Translator System. In *Proceedings of the Ninth Machine Translation Summit*, New Orleans, USA, 2003.
- [2] Eleftherios Avramidis. DFKI System Combination with Sentence Ranking at ML4HMT-2011. In *Proceedings of the International Workshop on Using Linguistic Information for Hybrid Machine Translation (LIHMT 2011) and of the Shared Task on Applying Machine Learning Techniques to Optimise the Division of Labour in Hybrid Machine Translation (ML4HMT)*, Barcelona, Spain, November 2011. META-NET.
- [3] Satanjeev Banerjee and Alon Lavie. METEOR: An Automatic Metric for MT Evaluation with Improved Correlation with Human Judgments. In *Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization*, pages 65–72, Ann Arbor, Michigan, June 2005. Association for Computational Linguistics.
- [4] Loïc Barrault. MANY : Open-Source Machine Translation System Combination. *Prague Bulletin of Mathematical Linguistics, Special Issue on Open Source Tools for Machine Translation*, 93:147–155, 2010.
- [5] Ondrej Bojar, Miloš Ercegovčević, Martin Popel, and Omar Zaidan. A Grain of Salt for the WMT Manual Evaluation. In *Proceedings of the Sixth Workshop on Statistical Machine Translation*, pages 1–11, Edinburgh, Scotland, July 2011. Association for Computational Linguistics.
- [6] Chris Callison-Burch, Cameron Fordyce, Philipp Koehn, Christof Monz, and Josh Schroeder. Further Meta-Evaluation of Machine Translation. In *Proceedings of the Third Workshop on Statistical Machine Translation*, pages 70–106, Columbus, Ohio, June 2008. Association for Computational Linguistics.
- [7] Jean Carletta. Assessing Agreement on Classification Tasks: the kappa Statistic. *Computational Linguistics*, 22:249–254, June 1996.
- [8] Christian Federmann. Appraise: An Open-Source Toolkit for Manual Phrase-Based Evaluation of Translations. In *Proceedings of the Seventh International Conference on Language Resources and Evaluation (LREC'10)*, Valetta, Malta, May 2010.
- [9] Christian Federmann, Yu Chen, Sabine Hunsicker, and Rui Wang. DFKI System Combination using Syntactic Information at ML4HMT-2011. In *Proceedings of the International Workshop on Using Linguistic Information for Hybrid Machine Translation (LIHMT 2011) and of the Shared Task on Applying Machine Learning Techniques to Optimise the Division of Labour in Hybrid Machine Translation (ML4HMT)*, Barcelona, Spain, November 2011. META-NET.
- [10] J.L. Fleiss. Measuring Nominal Scale Agreement among Many Raters. *Psychological Bulletin*, 76(5):378–382, 1971.
- [11] Philipp Koehn. Europarl: A Parallel Corpus for Statistical Machine Translation. In *Proceedings of MT Summit 2005*, 2005.
- [12] Klaus Krippendorff. Reliability in Content Analysis. Some Common Misconceptions and Recommendations. *Human Communication Research*, 30(3):411–433, 2004.

- [13] Shankar Kumar and William Byrne. Minimum Bayes-Risk Word Alignments of Bilingual Texts. In *Proceedings of the ACL-02 conference on Empirical methods in natural language processing - Volume 10*, EMNLP '02, pages 140–147, 2002.
- [14] J.R. Landis and G.G. Koch. Measurement of Observer Agreement for Categorical Data. *Biometrics*, 33(1):159–174, 1977.
- [15] Zhifei Li, Chris Callison-Burch, Chris Dyer, Sanjeev Khudanpur, Lane Schwartz, Wren Thornton, Jonathan Weese, and Omar Zaidan. Joshua: An Open-Source Toolkit for Parsing-Based Machine Translation. In *Proceedings of the Fourth Workshop on Statistical Machine Translation*, pages 135–139, Athens, Greece, March 2009. Association for Computational Linguistics.
- [16] Franz Josef Och. Minimum Error Rate Training in Statistical Machine Translation. In *ACL '03: Proceedings of the 41st Annual Meeting on Association for Computational Linguistics*, pages 160–167, Morristown, NJ, USA. Association for Computational Linguistics.
- [17] Tsuyoshi Okita and Josef van Genabith. DCU Confusion Network-based System Combination for ML4HMT. In *Proceedings of the International Workshop on Using Linguistic Information for Hybrid Machine Translation (LIHMT 2011) and of the Shared Task on Applying Machine Learning Techniques to Optimise the Division of Labour in Hybrid Machine Translation (ML4HMT)*, Barcelona, Spain, November 2011. META-NET.
- [18] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A Method for Automatic Evaluation of Machine Translation. IBM Research Report RC22176(W0109-022), IBM, 2001.
- [19] Sergio Penkale, Rejwanul Haque, Sandipan Dandapat, Pratyush Banerjee, Ankit K. Srivastava, Jinhua Du, Pavel Pecina, Sudip Kumar Naskar, Mikel L. Forcada, and Andy Way. MaTrEx: the DCU MT system for WMT 2010. In *Proceedings of the Joint Fifth Workshop on Statistical Machine Translation and MetricsMATR*, WMT '10, pages 143–148, Stroudsburg, PA, USA, 2010. Association for Computational Linguistics.
- [20] Gema Ramírez-Sánchez, Felipe Sánchez-Martínez, Sergio Ortiz-Rojas, Juan Antonio Pérez-Ortiz, and Mikel L. Forcada. OpenTrad Apertium open-source machine translation system: an opportunity for business and research. In *Proceeding of Translating and the Computer 28 Conference*, London, United Kingdom, November 2006.
- [21] William A. Scott. Reliability of Content Analysis: The Case of Nominal Scale Coding. *The Public Opinion Quarterly*, 19(3):321–325, 1955.
- [22] Matthew G. Snover, Nitin Madnani, Bonnie Dorr, and Richard Schwartz. TER-Plus: paraphrase, semantic, and alignment enhancements to Translation Edit Rate. *Machine Translation*, 23:117–127, September 2009.
- [23] Ralf Steinberger, Bruno Pouliquen, Anna Widiger, Camelia Ignat, Tomaz Erjavec, and Dan Tufis. The JRC-Acquis: A Multilingual Aligned Parallel Corpus with 20+ Languages. In *Proceedings of the Fifth International Conference on Language Resources and Evaluation (LREC'06)*, pages 2142–2147, 2006.
- [24] Andreas Stolcke. SRILM - An Extensible Language Modeling Toolkit. In *the 7th International Conference on Spoken Language Processing (ICSLP) 2002*, Denver, Colorado, 2002.
- [25] Ineke Schuurman Stella Markantonatou Sokratis Sofianopoulos Marina Vassiliou Olga Yannoutsou Toni Badia Maite Melero Gemma Boleda Michael Carl Vincent Vandeghinste, Peter Dirix and Paul Schmidt. Evaluation of a Machine Translation System for Low Resource Languages: METIS-II. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco, May 2008.