

Overview of the IWSLT 2010 Evaluation Campaign

Michael Paul

Marcello Federico

Sebastian Stüker

NICT
Hikaridai 2-2-2,
619-0288 Kyoto, Japan
michael.paul@nict.go.jp

FBK
via Sommarive 18,
38100 Trento, Italy
federico@fbk.eu

KIT
Adenauerring 2,
76131 Karlsruhe, Germany
sebastian.stueker@kit.edu

Abstract

This paper gives an overview of the evaluation campaign results of the 7th *International Workshop on Spoken Language Translation (IWSLT 2010)*¹. This year, we focused on three spoken language tasks: (1) public speeches on a variety of topics (*TALK*) from English to French, (2) spoken dialog in travel situations (*DIALOG*) between Chinese and English, and (3) traveling expressions (*BTEC*) from Arabic, Turkish, and French to English. In total, 28 teams (including 7 first-time participants) took part in the shared tasks, submitting 60 primary and 112 contrastive runs. Automatic and subjective evaluations of the primary runs were carried out in order to investigate the impact of different communication modalities, spoken language styles and semantic context on automatic speech recognition (ASR) and machine translation (MT) system performances.

1. Introduction

The *International Workshop on Spoken Language Translation (IWSLT)* is a yearly, open evaluation campaign for spoken language translation. IWSLT's evaluations are not competition-oriented; their goal is to foster cooperative work and scientific exchange. In this respect, IWSLT proposes challenging research tasks and an open experimental infrastructure for the scientific community working on spoken and written language translation.

Previous IWSLT workshops focused on the establishment of evaluation metrics for multilingual speech-to-speech translation, innovative technologies for the translation of automatic speech recognition results from read-speech and spontaneous-speech input, and monolingual and bilingual dialog conversations [1].

This year, the standard BTEC task was provided for the translation of Arabic and Turkish spoken language text into English. For the first time, *French* was used as an input language for the BTEC task, attracting new groups to participate in this year's event.

As a continuation of last years efforts in translating spoken dialog, the *DIALOG* task focused on task-oriented cross-lingual human dialog in travel situations where the speech data was annotated with dialog and speaker information that could be exploited by the participant to incorporate contextual information into the translation process. For the

DIALOG task, IWSLT participants had to translate both the Chinese and the English outputs of the automatic speech recognizers into English and Chinese, respectively.

The new challenge for this year's evaluation campaign was the translation of public speeches from English to French. The *TALK* task was based on a collection of recordings of public speeches covering a variety of topics, for which high quality transcriptions and translations into several languages are available. This task not only imposes new challenges on the development of MT systems, i.e., on how to deal with unlimited domains, but also on the applicability of standard evaluation protocols for the evaluation of translation results of automatic speech recognition outputs based on reference translations that are segmented differently.

All participants had to submit at least one run (*primary submission*) for each translation task they registered for. The evaluation of the primary runs was carried out using standard automatic evaluation metrics for all translation tasks. In addition to the single-metric scores, all automatic metric scores for the MT output were combined by normalizing each metric score distribution and calculating the average of all the normalized metric scores. Human assessments of translation quality ranking multiple MT systems were also applied for the *DIALOG* and *BTEC* tasks. Based on the evaluation results, the impact of different communication modalities (monologue vs. dialog), spoken language (planned vs. spontaneous) and semantic context (open vs. limited) was investigated.

The outline of the IWSLT 2010 evaluation campaign (its translation tasks and evaluation specifications) are described in detail in Section 2. The evaluation results are summarized and discussed in Section 3.

2. Outline of IWSLT 2010

This year's IWSLT campaign took place during the period of June-September 2010 and featured six different translation tasks that are summarized in Table 1.

In total, 28 research groups (including 7 first-time participants) from all over the world² participated in the event, producing a total of 60 MT engines. Information on the research groups, the utilized translation systems, and transla-

²China: 2, France: 4, Germany: 2, Ireland: 1, Israel: 1, Italy: 1, Japan: 3, Korea: 1, Netherlands: 2, Portugal: 1, Qatar: 1, Singapore: 1, Spain: 3, Tunisia: 1, Turkey: 2, UK: 1, USA: 1

¹<http://iwslt2010.fbk.eu>

Table 1: Translation Tasks

Task	Translation Direction	Participants
<i>TALK</i>	English-French	TT _{EF} 9
<i>DIALOG</i>	English-Chinese	DT _{EC} 11
	Chinese-English	DT _{CE} 11
<i>BTEC</i>	Arabic-English	BT _{AE} 12
	French-English	BT _{FE} 9
	Turkish-English	BT _{TE} 8

tion task participation is summarized in Appendix A. Most participants used phrase-based and syntax-based statistical machine translation (SMT) systems. However, one example-based MT (EBMT) system and various hybrid approaches combining multiple SMT engines or SMT engines with rule-based (RBMT) systems were also exploited.

A detailed description of the translation tasks and the language resources (*supplied corpora*) which were provided to the participating research groups are given in Section 2.1 (TALK), Section 2.2 (DIALOG), and Section 2.3 (BTEC), respectively. The supplied resources were released to the participants three months ahead of the official run submissions period. The official run submission period was limited to two weeks for the BTEC and DIALOG tasks and five weeks for the TALK task. Run submission was carried out via email to the organizers with multiple runs permitted. However, the participant had to specify which runs should be treated as *primary* (evaluation using human assessments and automatic metrics) or *contrastive* (automatic evaluation only). The organizers set-up online evaluation servers for the TALK develop data sets³ as well as the testdata sets⁴ of the BTEC and DIALOG tasks that could be used by the participants to tune their systems (TALK) or carry out additional experiments after the official run submission period (DIALOG, BTEC). The schedule of the evaluation campaign is summarized in Table 2.

Table 2: Evaluation Campaign Schedule

Event	Date
Training/Develop Corpus Release	May 28, 2010
Evaluation Corpus Release	Aug 23, 2010
Translation Result (BTEC/DIALOG)	Sep 6, 2010
Automatic Evaluation Results	Sep 17, 2010
Translation Results (TALK)	Sep 30, 2010
Subjective Evaluation Results	Nov 12, 2010
Workshop	Dec 2-3, 2010

2.1. TALK Task

The new challenge of this year was the translation of public speeches from English to French. The so-called TALK task was based on the TED⁵ talks collection, a Web repository of recordings of public speeches, mostly held in English, cover-

³<http://isl.ira.uka.de/iwslt2010>

⁴https://mastarpj.nict.go.jp/EVAL/IWSLT10/automatic/testset_IWSLT10

⁵<http://www.ted.com>

ing a variety of topics, and for which high quality transcriptions and translations into several languages are available.

The proposed new challenge clearly departs from and completes the application scenarios proposed so far in the IWSLT evaluations. Macroscopic differences between the TALK task and the BTEC and DIALOG tasks are in the assumed communication modality, i.e. monologue vs. dialogue, spoken language style, i.e. planned vs. spontaneous, and semantic context, i.e. open vs. limited.

From a translation point of view, the TALK task is basically a subtitling translation task, in which the ideal translation unit is a single caption as defined by the original transcript. In fact, some word re-ordering across consecutive captions is also permitted in order to accommodate syntactic differences between source and target languages. The wide variety of topics covered by the TED talks has determined the type and volume of training data that has been prepared and released for this challenge. This in fact comprises a small (less than 1 million word) parallel corpus of TED talks and several out-of-domain large parallel corpora including texts from the United Nations, European Parliament, news commentaries, and the Web.

From a speech translation point of view, the problem of processing full transcripts rather than isolated utterances requires handling possible inconsistencies between the speech segmentation introduced by the ASR system and the text segmentation used in the reference transcripts and translations. In particular, this discrepancy impacts when word-graphs produced by the ASR system are used as MT input.

While the significantly larger amount of available training data clearly has an impact on the complexity of the MT systems being developed for the TALK task, the problem of aligning ASR and reference segments also required some important revision of the automatic evaluation method.

These major shifts with respect to the previous IWSLT evaluations are the reasons why we declared that this first evaluation is to be considered as an exercise for establishing reference baselines and appropriate evaluation protocols for future evaluations. Hence, although an evaluation server was set up to compute several translation accuracy metrics, no official ranking of the participants will be reported in this evaluation. As no human evaluation was planned for the new challenge, a different schedule from the other tasks was established in order to ease participation in all of the offered tasks. Before the submission deadline, we received primary submissions by nine teams in total. The majority of these teams also participated in other IWSLT tasks (see Appendix A).

2.1.1. Language Resources

The TALK task is about the translation of speeches taken from the TED website. TED LLC is a nonprofit organization with the declared goal of “disseminating ideas worth spreading”. It regularly organizes two annual conferences in the US and one in UK, in which prominent experts from different fields are invited to give short talks about topics relevant

to the global society. Although TED stands for Technology Entertainment Design, over the years its scope has become much broader, indeed. TED is supported by industrial sponsors and a community of volunteers, which organize similar conferences in other countries of the world and help in creating content for the TED website, namely videos, transcripts and translations of talks. All content is copyrighted and made publicly available under a Creative Commons license. At this time, the TED website hosts around 800 English talks and the TED Open Translation Project has been managing translation of talks into 80 languages. Translations grow at a rapid pace thanks to more than 4,000 volunteers, which have contributed some 12,500 translations in total so far. To ensure quality in the process, all English source transcripts are prepared by professional transcribers and are revised by another translator prior to publication. Languages with the most translated talks at this time are Arabic, Bulgarian, Chinese, French, Italian, Portuguese, and Spanish.

For this IWSLT evaluation, a first parallel corpus of 345 English talks with their French translations was released, named *TED English-French ver.1.1*. The TED corpus and all additional parallel data allowed for this exercise have been made available at the workshop website⁶. Development and test sets, however, were released only to registered participants according to the evaluation schedule. The statistics⁷ of the supplied TALK corpus are summarized in Table 3.

Table 3: Supplied Corpus (TALK)

TALK	data	lang	sent	avg.len	word token	word type
train	(text)	E	86,225	9.8	842,125	31,429
	(text)	F	86,225	10.0	867,963	42,599
dev	(speech)	E	1,368	9.5	12,962	2,687
	(text)	F	1,368	9.3	12,712	3,246
test10	(speech)	E	3,584	9.0	32,155	4,153
	(text)	F	3,584	9.2	33,010	5,571

2.1.2. Task Definition

For the TALK task, participants were requested to translate two input conditions: (1) the reference text that was extracted from the subtitles of the TED talks, and (2) the output from an automatic speech recognition (ASR) system run on the audio of the TED talks selected for the evaluation set. The reference texts were in true case and contained punctuation marks. The segmentation was given by the segmentation of the closed captioning of the TED data. The output from the ASR system was case-sensitive, but did not contain any punctuation marks. The segmentation was obtained automatically from the audio data and thus did not match the reference segmentation of the closed captions. The ASR output provided to the participants consisted of the single best output, a 20-best list, and the word lattices from the recognition system provided in standard lattice format (SLF).

⁶<http://iwslt2010.fbk.eu/node/27>

⁷For details on the additional language resources that were permitted for the TALK task, please refer to <http://www.statmt.org/wmt10/translation-task.html>.

The ASR system used for producing the automatic transcripts was the 2009 KIT English Quero Evaluation system with a language model that was updated with the data from the TED training data. In order to measure the ASR performance we took the subtitles of the talks as provided by the TED website and re-annotated the time boundaries of the sentences to exactly match the speech. The ASR system achieved a case-insensitive word error rate (WER) of 26.4% on the TED development set, and 22.3% on the evaluation set, respectively.

The quality of the results of the automatic translation systems was measured with BLEU, NIST, and TER scores (see Table 9) using one reference translation. The reference translation was taken from the translations provided by the TED open translation project. All translations were supposed to be case-sensitive. Also, regardless of whether the reference transcription of the talks, which contain punctuation marks, or the automatic transcriptions, which do not, were used, the translation systems were supposed to produce punctuation marks which were considered in the automatic measures.

Since the reference translations from the TED website match the segmentation of the reference transcriptions of the talks, the scores for the automatic translation results could be directly computed. This was not the case for the translation of the ASR output, as the segmentation of the ASR output does not match that of the reference translation. We therefore used the method and scoring scripts from [2] which align the automatic translation and reference translation based on the Levenshtein distance first and then computes the automatic translation scores.

2.2. DIALOG Task

As a continuation of last year's efforts [1], the DIALOG task focused on the translation of task-oriented human dialog in travel situations. The speech data was recorded through human interpreters, where native speakers of different languages were asked to complete certain travel-related tasks like hotel reservations using their mother tongue. The translation of the freely-uttered conversation was carried out by human interpreters. The obtained speech data was annotated with dialog and speaker information. In total, 11 research groups participated in this year's DIALOG task (see Appendix A).

2.2.1. Language Resources

The DIALOG task was carried out using the Spoken Language Databases (SLDB) corpus, a collection of human-mediated cross-lingual dialogs in travel situations. Similar to last year, bilingual Chinese-English dialogs were provided to the participants for the training of the MT systems. In addition, the Chinese/English parts of the BTEC corpus (see Section 2.3), were provided to the participants of the DIALOG task and could be used as additional training bitext.

Linguistic tools like word segmentation tools, parsers, etc., were allowed to preprocess the supplied corpus, but par-

Table 4: Supplied Corpus (DIALOG)

BTEC data	lang	sent	avg.len	word token	word type
train	(text) C	19,972	7.4	148,224	8,408
	(text) E	19,972	7.7	153,178	7,294
dev	(speech) C	1,495	9.4	14,002	3,409
	(ref) E	15,029	10.3	139,212	6,176
	(speech) E	506	6.2	3,119	840
	(ref) C	3,542	7.1	25,037	1,665
	(text) C	1,741	5.5	9,666	2,920
	(ref) E	20,762	6.8	141,262	6,306

SLDB data	lang	sent	avg.len	word token	word type
train	(text) C	10,061	8.9	89,110	3,734
	(text) E	10,061	11.8	118,648	3,271
dev	(dialog) C	200	9.3	1,859	377
	(ref) E	800	9.8	7,829	418
	(dialog) E	210	11.8	2,474	403
	(ref) C	840	11.2	9,379	621
	(speech) C	750	5.1	3,818	633
	(ref) E	5,208	6.6	33,827	1,387
	(speech) E	749	5.5	4,146	454
	(ref) C	5,243	6.5	34,693	1,265
test09	(dialog) C	405	11.3	4,562	653
	(ref) E	1,620	13.7	22,253	886
	(dialog) E	393	11.0	4,321	569
	(ref) C	1,572	12.0	18,789	875
test10	(dialog) C	532	8.2	4,361	900
	(ref) E	2,128	13.3	28,384	1,636
	(dialog) E	453	11.0	5,004	870
	(ref) C	1,812	11.2	20,314	1,470

ticipants were asked to declare their usage in the system description paper and to measure the impact of these tools on the system performance. No additional parallel or monolingual corpora or word-lists were permitted to be used for the primary run. However, in order to motivate participants to explore the effects of additional language resources, the organizers also accepted contrastive runs based on additional resources.

Table 4 summarizes the characteristics of the Chinese (C) and English (E) training (*train*), development (*dev*) and evaluation (*test*) data sets. For evaluation purposes, two data sets, i.e., the testset of IWSLT 2009 (*test09*) and the new testset of IWSLT 2010 (*test10*) were used. The first two columns specify the given data set and its type. Besides the source language text (“text”) and target language reference translation (“ref”) resources, all data sets consist of the ASR output and manual transcriptions of the respective *dialog* or *speech* recordings of language *lang*. The number of sentences are given in the “*sent*” column, and the “*avg.len*” column shows the average number of words per training sentence, where the word segmentation for the source language was the one given by the output of the ASR engines without punctuation marks. “*word token*” refers to the number of words in the corpus and “*word type*” refers to the vocabulary size.

For the automatic evaluation of development data sets, 7 (16) reference translations for the SLDB (BTEC) were also included in the supplied corpus. For the DIALOG testset data sets, up to 4 reference translations were available.

2.2.2. Task Definition

For the DIALOG task, participants were asked to translate two input conditions: (1) the automatic speech recognition (ASR) outputs, i.e., word lattices (SLF), N-best lists (NBEST) and 1-best (1BEST) speech recognition results, and (2) the correct recognition results (CRR), i.e., text input without speech recognition errors.

For both input conditions, the input text contained neither case nor punctuation information. However, the reference translations were in true case and contained punctuation marks. Therefore, the participants had to recover case/punctuation information for the MT output run submissions. Instructions⁸ on how to build a baseline tool for case/punctuation insertions using the *SRI Language Modeling Toolkit* [3] was also provided.

Participants of the DIALOG task had to translate both the English ASR outputs into Chinese and the Chinese ASR outputs into English, whereby they could choose the ASR output condition (SLF, NBEST, or 1BEST) that best suits their MT system. Translation of the CRR text input was mandatory for all participants.

The ASR systems used to create the ASR outputs were the Chinese and English ATRASR systems provided by NICT [4]. The recognition accuracies for the DIALOG testdata sets are summarized in Table 5. Besides the ASR output files (lattices, 20-BEST and 1-BEST lists), tools to extract larger NBEST lists were also provided to the participants.

Table 5: Speech Recognition Accuracy (DIALOG)

DIALOG	lang	word (%)		sentence (%)	
		lattice	1BEST	lattice	1BEST
test09	C	92.67	81.46	64.63	39.12
	E	89.58	82.20	50.13	37.15
test10	C	89.36	83.29	61.60	54.64
	E	89.06	81.11	44.30	34.81

2.3. BTEC Task

In order to (1) enable small groups and even newcomers to the field of machine translation to join the evaluation campaign and (2) provide a testbed for new ideas for spoken language translation techniques, a BTEC translation task focusing on frequently used utterances in the domain of travel conversations was provided for the translation of Arabic (A), French (F) and Turkish (T) spoken language text into English (E). In total, 20 research groups took part in at least one of the three BTEC translation tasks, submitting 12 primary runs for Arabic-English, 9 primary runs for French-English, and 8 primary runs for Turkish-English.

2.3.1. The BTEC Corpus

The BTEC task was carried out using the Basic Travel Expression Corpus (BTEC), a multilingual speech corpus con-

⁸http://mastarpj.nict.go.jp/IWSLT2009/downloads/case+punc_tool_using_SRILM.instructions.txt

Table 6: Supplied Corpus (BTEC)

BTEC	data	lang	sent	avg.len	word token	word type
train	(text)	A	19,972	8.0	158,926	18,154
	(text)	F	19,972	9.5	189,665	10,735
	(text)	T	19,972	7.0	139,514	20,106
	(text)	E	19,972	9.1	182,627	8,344
dev	(text)	A	2,508	6.3	15,797	3,875
	(ref)	E	35,238	8.1	284,612	5,609
	(text)	F	1,512	7.5	11,409	2,244
	(ref)	E	24,192	8.1	196,806	4,660
	(text)	T	1,006	5.7	5,766	2,083
	(ref)	E	16,096	8.1	130,518	3,712
test09	(text)	A	469	6.1	2,875	1,099
	(text)	F	469	7.8	3,642	976
	(text)	T	469	5.8	2,741	1,115
	(ref)	E	3,283	8.4	27,507	1,739
test10	(text)	A	464	6.4	2,953	1,180
	(text)	F	464	7.7	3,582	1,004
	(text)	T	464	5.8	2,710	1,149
	(ref)	E	3,248	8.4	27,183	1,580

taining tourism-related sentences similar to those that are usually found in phrasebooks for tourists going abroad.

The participants were supplied with a training corpus of 20K sentence pairs which covered the same sentence IDs for all translation directions. In addition, the testsets of previous IWSLT evaluation campaigns were also provided to the participants and could be used to improve the MT system performance for the respective translation tasks. In contrast to the DIALOG task, the supplied corpus of the BTEC task was in true case and contained punctuation marks. The corpus statistics are summarized in Table 6.

2.3.2. Task Definition

The translation input condition of all BTEC tasks consisted of correct recognition results, i.e., text input, for Arabic, Turkish, and French. The target language for all BTEC tasks was English. The monolingual and bilingual language resources that were allowed for training the translation engines for the primary runs were limited to the supplied corpus. All other BTEC language resources besides the ones for the given language pair were treated as additional language resources.

Similar to the DIALOG task, the evaluation specifications for the BTEC task were defined as case-sensitive with punctuation marks (*case+punc*). Tokenization scripts were applied automatically to all run submissions prior to evaluation. In addition, automatic evaluation scores were also calculated for case-insensitive (lower-case only) MT outputs with punctuation marks removed (*no_case+no_punc*).

2.4. Evaluation Specifications

In this section, we summarize the subjective and automatic evaluation metrics used to assess the translation quality of the primary run submissions.

2.4.1. Subjective Evaluation

Human assessments of translation quality were carried out using the *Ranking* metrics. For the *Ranking* evaluation, hu-

man graders were asked to “rank each whole sentence translation from Best to Worst relative to the other choices (ties are allowed)” [5]. The *Ranking* evaluation was carried out using a web-browser interface and graders had to order up to five system outputs by assigning a grade between 5 (*best*) and 1 (*worse*). This year’s evaluations were carried out by paid evaluation experts, i.e., three graders for each of the target languages. The *Ranking* scores were obtained as the average number of times that a system was judged better than any other system. In addition, normalized ranks (*NormRank*) on a per-judge basis using the method of [6] were calculated for each run submission. The *Ranking* metric was applied to all submitted primary runs of all translation tasks.

Similar to last year’s IWSLT edition [1], the difference of each MT system and the system that obtained the highest *Ranking* score (*BestRankDiff*) was calculated and used to define an alternative method to rank MT systems of a given translation task.

In addition, human assessments of the overall translation quality of a single MT system were carried out with respect to the *Fluency* and *Adequacy* of the translation. *Fluency* indicates how the evaluation segment sounds to a native speaker of the target language. For *Adequacy*, the evaluator was presented with the source language input as well as a “gold standard” translation and had to judge how much of the information from the original translation was expressed in the translation [7]. The *Fluency* and *Adequacy* judgments consisted of one of the grades listed in Table 7. The evaluation of both metrics, *Fluency* and *Adequacy*, was carried out separately using a web-browser tool. For each input sentence, the MT translation outputs of the respective systems were displayed on one screen and judgments were done by selecting one of the possible grades for each MT output.

In addition to the above standard metrics, a modified version of the *Adequacy* metrics (*Dialog*) that takes into account information beyond the current input sentence was applied to the translation results of the DIALOG task in order to judge a given MT output in the context of the respective dialog. For the *Dialog* assessment, the evaluators were presented with the history of previously uttered sentences, the input sentence and the “gold standard” translation. The evaluator had to read the dialog history first and then had to judge how much of the information from the reference translation is expressed in the translation in the context of the given dialog history by assigning one of the *Dialog* grades listed in Table 7. In cases where parts of the information were omitted in the current translation, but they could be understood in the context of the given dialog, such omission would not result in a lower *Dialog* score.

Due to high evaluation costs, the *Fluency*, *Adequacy*, and *Dialog* assessments were limited to the top-ranked MT system for each translation task according to the *Ranking* evaluation results. In addition, the translation results of each translation task were *pooled*, i.e., in cases of identical translations of the same source sentence by multiple engines, the pooled

Table 7: Human Assessment

Fluency		Adequacy / Dialog	
5	Flawless Chinese/English	5	All Information
4	Good Chinese/English	4	Most Information
3	Non-native Chinese/English	3	Much Information
2	Disfluent Chinese/English	2	Little Information
1	Incomprehensible	1	None

translation was graded only once, and the respective rank was assigned to all MT engines with the same output.

For the final metric scores, each system score is calculated as the *median* of the assigned grades. All paid graders took part in a dry-run evaluation exercise prior to this year’s evaluation period in order to get used to the evaluation metrics as well as the browser-based graphical user interfaces.

2.4.2. Grader Consistency

In order to investigate the degree of grading consistency between the human evaluators, we calculated *Fleiss’ kappa coefficient* κ , which measures the agreement between two raters who each classify N items into C mutually exclusive categories taking into account the agreement occurring by chance. It is calculated as:

$$\kappa = \frac{\Pr(a) - \Pr(e)}{1 - \Pr(e)},$$

where $\Pr(a)$ is the relative observed agreement among graders, and $\Pr(e)$ is the hypothetical probability of chance agreement. If the raters are in complete agreement then $\kappa = 1$. If there is no agreement among the raters (other than what would be expected by chance) then $\kappa \leq 0$. The interpretation of the κ values according to [8] is given in Table 8.

Table 8: Interpretation of κ Coefficient

κ	Interpretation
< 0	No agreement
0.0 – 0.20	Slight agreement
0.21 – 0.40	Fair agreement
0.41 – 0.60	Moderate agreement
0.61 – 0.80	Substantial agreement
0.81 – 1.00	Almost perfect agreement

2.4.3. Automatic Evaluation

The automatic evaluation of run submissions was carried out using the standard automatic evaluation metrics listed in Table 9. Besides the NIST metrics, all automatic evaluation metric scores listed in Appendix C are given as percent figure (%). For the DT_{EC} translation task, F1 scores calculated based on the *unigram precision* and *recall* system-level figures of each MT systems are used instead of the METEOR metric scores.

In addition to the single-metric scores of each MT output, the average of all automatic evaluation scores (**z-avg**) is calculated as follows. In the first step, all metric scores are normalized so that the score distribution of the respective metric has a zero mean and unit variance (*z-transform*). In

Table 9: Automatic Evaluation Metrics

BLEU:	the geometric mean of n-gram precision by the system output with respect to reference translations. Scores range between 0 (worst) and 1 (best) [9] → ‘mteval-v13a.pl’
NIST:	a variant of BLEU using the arithmetic mean of weighted n-gram precisions. Scores are positive with 0 being the worst possible [10] → ‘mteval-v13a.pl’
METEOR:	calculates unigram overlaps between a translation and reference texts taking into account various levels of matches (<i>exact, stem, synonym</i>). Scores range between 0 (worst) and 1 (best) [11] → ‘meteor-v1.0’
GTM:	measures the similarity between texts by using a unigram-based F-measure. Scores range between 0 (worst) and 1 (best) [12] → ‘gtm-v1.4’
WER:	<i>Word Error Rate</i> : the edit distance between the system output and the closest reference translation. Scores are positive with 0 being the best possible [13]
PER:	Position independent word error rate: a variant of WER that disregards word ordering [14]
TER:	<i>Translation Edit Rate</i> : a variant of WER that allows phrasal shifts [15] → ‘tercom-0.7.25’

the second step, the obtained z-scores of a given MT system are averaged to obtain the final z-avg system score [1].

2.4.4. Statistical Significance of Evaluation Results

In order to decide whether the translation output on the document-level of one MT engine is significantly better than another, we used the *bootStrap* method that (1) performs a random sampling with replacement from the evaluation data set, (2) calculates the respective evaluation metric score of each engine for the sampled test sentences and the difference between the two MT system scores, (3) repeats the sampling/scoring step iteratively, and (4) applies the *Student’s t-test* at a significance level of 95% confidence to test whether the score differences are significant [16]. In this year’s evaluation, 2000 iterations were used for the analysis of the automatic evaluation results.

2.4.5. Correlation between Evaluation Metrics

Correlations between different metrics were calculated using the *Spearman rank correlation coefficient* ρ which is a non-parametric measure of correlation that assesses how well an arbitrary monotonic function can describe the relationship between two variables without making any assumptions about the frequency distribution of the variables. It is calculated as:

$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)},$$

where d_i is the difference between the rank of the system i and n is the number of systems.

3. Main Findings of IWSLT 2010

The subjective evaluation results of IWSLT 2010 are summarized in Appendix B. In addition to the MT outputs provided by the participants, the organizers used an online MT server to translate the testset data sets. The online system (*online*) represents a state-of-the-art general-domain MT system that differs from the participating MT systems in two aspects: (1) its language resources are not limited to the supplied corpora and (2) its parameters are not optimized using in-domain data. Its purpose is to investigate the applicability of a baseline system with unlimited language resources to the spoken language translation tasks investigated by the IWSLT evaluation campaign. Section B.1 illustrates the overall performance of the best MT system for each translation task and the online system in terms of the *Fluency*, *Adequacy*, and *Dialog* metric scores. The *Ranking*, *NormRank*, and *BestRankDiff* metric results of all MT systems participating in the DIALOG and BTEC tasks are given in Section B.2 and Section B.3, respectively.

The automatic evaluation results of two testsets (*test10* and *test09*) are given in Appendix C for two different subsets of the evaluation data: (1) the subset of testset sentence IDs used for human assessment where the scores are given as the mean score of the significance test described in Section 2.4.4 (see Section C.1), and (2) the full testset translated by the participants where the scores were obtained by the online evaluation server⁹ (see Section C.2). The MT systems are ordered according to the *z-avg* score of the metric combination for the *case+punc* evaluation specifications that achieved the highest rank correlation coefficients toward the subjective *Ranking* evaluation metric. If system performances *do not* differ significantly according to the *bootStrap* method, horizontal lines between two MT engines in the MT engine ranking tables are omitted. For each translation task, the highest (lowest) scores of the respective evaluation metric are highlighted in **boldface** (*italics*).

Finally, Appendix D summarizes the rank correlation coefficients of subjective and automatic evaluation results.

3.1. TALK Task

This section summarizes the main features of the systems that have been developed by the nine participants of the TALK task. This information is derived from the system descriptions provided by each team.

All participants approached the exercise with phrase-based statistical MT relying on linear combination of feature functions. In particular, seven teams employed the Moses decoder, *limsi* used an *n*-gram decoder, and *kit* used an in-house phrase-based decoder.

The best, median and worst BLEU scores of the primary submissions are given in Table 10 for two evaluation specifications: (1) case-sensitive with punctuations tokenized

Table 10: Automatic Evaluation (TALK)

Eval Spec.	Input Cond.	BLEU (%)		
		Best	Median	Worst
<i>case+punc</i>	Text	29.90	25.02	24.24
	ASR	16.34	15.68	12.13
<i>no_case+no_punc</i>	Text	29.98	26.42	23.32
	ASR	20.27	18.75	16.43

(“*case+punc*”) and (2) case-insensitive with punctuations removed (“*no_case+no_punc*”).

The important difference between the best scores achieved in the text and ASR conditions can be explained by the relatively high word-error-rate in the ASR transcripts, i.e., more than one word in every five was wrongly transcribed, and the impact of missing punctuation and letter case information in the ASR transcripts. The specific impact of speech recognition errors can be in part measured by comparing the scores computed on rich and on plain outputs; that is, BLEU scores considering letter case and punctuation (*case+punc*) versus BLEU scores disregarding such information (*no_case+no_punc*). While under the text input condition there is basically no difference between the rich and plain output evaluation, a more significant difference is observed under the ASR condition: the best system score changes from 16.34 to 20.27. As the best scores of the text-plain and ASR-plain conditions are from the same system, we can infer that the impact of speech recognition errors on BLEU scores is around 32%.

Concerning the systems that were developed for the evaluation, most of the participants focused on data filtering, data selection, and model adaptation. The reason for this in general was to find effective ways to make use of the large amount of out-of-domain parallel data that was provided. These approaches account indeed for most of the improvements claimed by the participants over their baselines.

Data filtering methods to extract reliable parallel data from the training data were reported by *fbk*, *kit*, *lium*, and *tubitak*. Data selection to extract parallel data relevant or close to the TALK task was applied by *fbk*, *iti-upv*, *lium*, and *ntt*. Model combination techniques were applied to the language and translation models to weight the contributions of different data sources. In particular, LM interpolation was applied by *fbk*, *kit*, *lig*, *limsi*, *lium*, *mit* and *ntt*. Interpolation of translation models was applied by *fbk* and *mit*. In contrast, *kit* applied a fixed combination scheme to merge two phrase tables. Parameter tuning of the scoring functions in the log-linear model was performed mostly with MERT, with the exception of *mit*, which also reports results with the MIRA algorithm, and *iti-upv*, which compares MERT with a new Bayesian adaptation method.

Concerning the introduction of novel feature functions, *kit* integrated a bilingual LM in its phrase-based decoder, and *limsi* introduced a re-ordering POS-based LM in its *n*-gram model. System combination was applied by only *mit*, while the use of additional resources (Wikipedia) for language modeling was only explored by *limsi*.

⁹*test10*: https://mastarpj.nict.go.jp/EVAL/IWSLT10/automatic/testset_IWSLT10
test09: https://mastarpj.nict.go.jp/EVAL/IWSLT10/automatic/testset_IWSLT09

Work to cope with issues related to the ASR input condition ranged from pre/post-processing methods that handle true casing and punctuation (all participants) up to training a specific MT system that processes ASR word-graphs (*lium*). From the reports of the participants, it seems that developing specific systems for each input condition definitely has rewards in terms of performance.

3.2. DIALOG Task

For the DIALOG task, eleven primary run submissions were submitted. Five participants (*inesc-id*, *postech*, *tubitak*, *uva-illc*, *uva-isca*) employed a single-engine phrase-based SMT approach based on the Moses decoder to translate the bilingual task-oriented human dialogs between Chinese and English. However, the majority of the participants (*dcu*, *i2r*, *ict*, *iti-upv*, *msra*, *nict*) made use of a hybrid MT system architecture combining two or more phrase-based SMT (*PBSMT*), hierarchical phrase-based SMT (*HPBSMT*), or syntax-based SMT (*SBSMT*) engines.

In particular, a standard phrase-based SMT system based on the Moses toolkit was combined with (a) an ITG-based SBSMT system by *iti-upv*, (b) an in-house PBSMT (*Lavender*) and SBSMT system (*Tranyu*) by *i2r*, (c) two HPBSMT systems (SAMT, CCG-based SMT) by *dcu*, (d) two in-house SBSMT systems (*SuperSilenus*, *TemBruin*) and a HPBSMT system (*John*) by *ict*, (e) in-house implementations of PBSMT and HPBSMT, two extended versions of the previous systems using a dependency tree language model, an SBSMT system, and a Treelet-based SMT system by *msra*. *Nict* combined in-house implementations of a PBSMT system (*CleopATRa*) and a HPBSMT system (*Linparse*).

Concerning system combination, median string computation (*iti-upv*), rescoring of combined n-best lists (*i2r*, *ict*, *nict*), and confusion network decoding techniques (*dcu*, *i2r*, *msra*, *nict*) were used. Moreover, new techniques investigated by this year's participants to improve system performance on the DIALOG task include: (1) the paraphrasing of the training data to address the data sparseness problem (*dcu*), (2) the handling of ASR errors using word-to-pinyin conversion (*ict*), confusion network decoding (*iti-upv*) and reranking of ASR output prior to decoding (*msra*), (3) the integration of multiple segmentation schemes for Chinese (*ict*, *nict*, *postech*), (4) the source language side re-ordering via tree induction (*uva-illc*), (5) the combination of multiple word alignment methods (*i2r*, *inesc-id*, *msra*), (6) the incorporation of syntactic constraints (*dcu*, *i2r*, *msra*), and (7) the exploitation of contextual information of the given dialog (*nict*, *uva-isca*). Experiments involving additional resources beyond the supplied corpus were conducted by *iti-upv*.

The human assessment results for the IWSLT 2010 DIALOG testset based on the system ranking evaluation are summarized in Appendix B.2 for all participating MT systems. The *NormRank* scores achieved for the CRR input condition are much higher than the ones obtained for the translation of the ASR output for both translation directions. Moreover, the

translation quality of the English-Chinese (EC) MT systems is higher than the Chinese-English (CE) MT systems for the majority of the participating teams.

Comparing the *Ranking* and *NormRank* results, quite different MT system rankings are obtained for DT_{CE} , especially in the case of the ASR output translation task. In contrast, the DT_{EC} systems are ranked very similarly with minor differences for systems in the mid-range. Both metrics, however, agree at least on the top-ranked MT system for both translation directions.

The more stable rankings for the DT_{EC} vs. the DT_{CE} systems and the CRR vs. the ASR input condition indicate that the reliability of human assessment grading depends to some extent on the overall translation quality of the MT system outputs. For humans, it is more difficult to distinguish between MT systems with relatively lower translation quality, but it is easier to identify the best performing systems.

As an alternative ranking method, we investigated the gain that the best performing system achieved over the other systems on sentence-level. For each MT system, we calculated the ratio of translations that were ranked worse and those that were ranked better than the top-ranked system for a subset of around 300 translations where both systems were judged together. The results summarized in Appendix B.3 show that much smaller gains were achieved by the best system for the ASR output condition (DT_{CE} : 18%~43%, DT_{EC} : 6%~40%) compared to the CRR translation results (DT_{CE} : 35%~68%, DT_{EC} : 37%~71%). Moreover, the difference was much lower for DT_{EC} than for DT_{CE} .

The MT systems ranked most consistently for both translation directions are *ict*, *msra*, and *nict*. In addition, much higher ranks for the ASR vs. the CRR input condition were achieved by *i2r* and *postech* for DT_{CE} and *iti-upv* for DT_{EC} .

In order to get an idea of the absolute translation quality of this year's participating MT systems, *Fluency/Adequacy* (isolated sentences) and *Dialog* (within the context of the given dialog) assessments were carried out for the best ranked *ict* system and the *online* system outputs.

The results listed in Appendix B.1 confirm that the translation quality of the DT_{EC} systems is much higher than the DT_{CE} systems for both input conditions and all subjective evaluation metrics (Fluency: +0.45~0.57, Adequacy: +0.40~0.81, Dialog: +0.29~0.62). However, the Fluency/Adequacy scores are relatively low for the ASR input condition (2.4/2.9 out of 5) for both translation directions and the DT_{CE} translations of the correct recognition results (2.9 out of 5). On the other hand, the best DT_{EC} system achieved moderate scores of 3.6/3.7. Moreover, the lower human assessment scores of the *online* system for both translation directions indicate that current state-of-the-art general-domain MT system have difficulties in handling ill-formed inputs like noisy speech (ASR errors) or spontaneous language styles (ungrammatical constructions).

Comparing the Adequacy and Dialog results, consistently higher scores (DT_{CE} : +0.3, DT_{EC} : +0.2~0.3) for both input

conditions were achieved when the context of the dialog was taken into consideration. This confirms the findings of last year’s evaluation campaign on the same task (DT_{CE} : +0.3, DT_{EC} : +0.1) and indicates that much information necessary to understand a given translation is provided by the history of previously uttered sentences. Therefore, evaluation metrics for the translation of task-oriented dialogs should not be carried out on a sentence-by-sentence basis, but within the context of the given dialog.

The automatic evaluation results confirm the findings of the subjective evaluation, i.e., the scores for the noisy ASR input condition are lower than for the CRR inputs for both translation directions. However, the differences are much larger for the DT_{EC} task (BLEU: +7.1, TER: -6.8, GTM: +6.7) compared to the DT_{CE} task (BLEU: +2.1, TER: -1.9, GTM: +3.4), indicating a higher negative impact of ASR errors on the translation of English input sentences. Looking at the speech recognition results listed in Table 5, we can see that the word-level recognition accuracies of English and Chinese ASR engines are quite similar for both lattices and 1BEST recognition results. However, on sentence-level, the accuracy figures for English (lattice: 44.3%, 1best: 34.8%) are far worse than those for Chinese (lattice: 61.6%, 1best: 54.6%), which underlines the importance of handling ASR recognition errors in the context of the whole input sentence and the preceding dialog. More gains are to be expected for dealing with n-best lists or even lattice input than single best recognition hypotheses.

In addition, we compared the automatic evaluation results obtained for the *test09* testset for all participants that took part in both the 2009 (see [1], Appendix D.2) and the 2010 (see Appendix C.2.2) evaluation campaigns. The results showed that the majority of systems were able to improve their system performance for all automatic evaluation metrics based on last year’s experiences, thus confirming the progress over time made by the best performing systems.

3.3. BTEC Task

A total of 29 MT engines were developed by the 20 participants of the BTEC tasks, with 12, 9, and 8 primary run submissions for the translation of Arabic, French and Turkish spoken language text into English, respectively.

The majority of the participants (14 teams) focused on phrase-based SMT approaches. In addition, an example-based approach was used by *tau* and an n-gram-based SMT approach was used by *dsic-upv*. *Tottori* combined a pattern-based MT approach with a standard phrase-based MT approach. Moreover, hybrid MT approaches combining phrase-based and hierarchical phrase-based SMT systems were investigated by *dcu*, *lig*, and *rwth*. Besides the Moses decoder, in-house phrase-based SMT decoders were employed by *apptek*, *kit*, and *nict*. An open-source hierarchical phrase-based SMT system (*Jane*) was used by *rwth*.

One of the main points of interest of this year’s BTEC task was the identification of word segmentation that helps im-

prove translation performance. Especially for Arabic, many segmentation schemes were explored, including BAMA (*qmul*, *tau*), MADA (*dcu*), ASVM (*lig*), AMIRA (*fbk*) and several in-house segmenters (*greyc*, *miracl*). In addition, the integration of multiple segmentation schemes into the translation process was investigated by *apptek* and *rwth*. For Turkish, the MORFESSOR segmentation toolkit was used by several participants including *qmul* and *tubitak*, and in-house segmenters were applied by *apptek*, *fbk*, and *limsi*. In contrast, the morphological analysis of the French input data was limited to simple tokenization preprocessing for most of the submitted primary runs. However, a stemming approach to reduce the data sparseness problem was applied by *kit* and a learning approach focusing on collocation segmentations was investigated by *upc*.

Other techniques exploited by the participants to improve system performance on the BTEC task include: (1) the combination of multiple word alignment methods (*apptek*, *mit*), (2) a phrase training method using forced alignment (*rwth*), (3) the incorporation of neural network language models (*dsic-upv*), (4) the application of new reordering models covering part-of-speech-based reordering (*apptek*, *kit*), short distance morpheme reordering (*limsi*) and dynamic distortion (*qmul*), (5) the incorporation of syntactic constraints (*rwth*), (6) the handling of unknown words (*apptek*, *limsi*, *qmul*), and (7) system combination techniques based on confusion network decoding (*lig*, *mit*, *rwth*).

The human assessment results for the BTEC task are summarized in Appendix B. The Ranking results listed in Appendix B.2 showed that the *online* system slightly outperformed the participating MT systems. This indicates the potential of using general-domain MT systems for the translation of spoken language input text that does not contain recognition errors. However, the gains were quite small, despite the fact that the MT systems of the IWSLT participants were trained on only 20k bitext.

Fluency/Adequacy grades were obtained for the best ranked MT systems of each translation task (*apptek* for BT_{AE} , *dsic-upv* for BT_{FE} , and *tubitak* for BT_{TE}) and the *online* system. The highest scores were achieved for BT_{FE} (4.0/4.3 out of 5), followed by BT_{TE} (3.7/4.0 out of 5), and BT_{AE} (3.3/3.6 out of 5).

The pair-wise comparison of each MT system with the *online* system listed in Appendix B.3 revealed that the difference in translation performance at sentence level for the BT_{FE} task is very small, i.e., less than 9% of the testset sentences were translated better by the *online* system compared to the majority of the participant’s MT systems. In particular, around 30% of the testset sentences were ranked equally, 31.6%~35.5% were ranked worse, and 35.3%~40.0% were ranked better. The gains for the BT_{TE} task are slightly higher, i.e. up to 18% of the testset sentences. For the BT_{AE} , however, the *online* system is outperformed by the *tubitak* and the *mit* systems, gaining +6.6%/+2.3%, respectively.

Similar to the DIALOG task, the comparison of the auto-

matic evaluation scores obtained for the *test09* evaluation data set that were submitted by participants who also took part in last year’s shared task also confirmed that progress is being made over time for the BTEC task. For the BT_{AE} task, significant gains (BLEU: +1.1, METEOR: +0.8, TER: -0.2, GTM: +1.0, NIST: +0.5) were achieved by the *mit* system against last year’s best system, a joined submission of *mit* and *tubitak*. For the BT_{TE} task, last year’s best performing system combination (*mit+tubitak*) would not be outperformed by this year’s participants. However, the individual system performance of the *tubitak*, *fbk*, and *apptek* systems improved by +1.0~8.7% BLEU, +4.2~8.4% METEOR, -0.3~2.1% TER, +0.2~4.2% GTM, and +0.2~0.6 NIST points.

3.4. Evaluation Metric Correlation

In order to get an idea of how closely the human assessment and automatic evaluation metrics are related, the *Spearman rank correlation coefficients* are summarized in Appendix D.

For each translation task, the MT system ranking obtained for the subjective *Ranking*, *NormRank*, *BestRankDiff* metrics and all investigated automatic evaluation metrics including the *z-avg* metric combination method are compared. For the DIALOG task, the correlation coefficients for ASR and CRR translation results are calculated separately.

The results show that the highest correlation to subjective evaluation metrics is obtained for the *z-avg* metric for the majority of the investigated translation tasks. In contrast to last year’s evaluation campaign where the *z-avg* score was calculated as the average of all investigated automatic evaluation metrics, this year we calculated the *z-avg* score for all possible combinations and selected the metric subset that achieved the highest correlation for each translation task separately. The selected metric combinations are summarized in Appendix D.

However, the optimal subset and correlation coefficients largely depend on the translation task. For the DT_{CE} and DT_{EC} tasks, the highest correlation was achieved for the *BestRankDiff (Ranking)* metric when the *online* system is included in (excluded from) the MT system rankings. For the BTEC tasks, in general, the *z-avg* scores correlates best with the *NormRank* metric.

3.5. Grader Consistency

Each sentence was evaluated by three human judges. Due to different levels of experience and background of the evaluators, variations in judgments were to be expected. Besides the *inter-grader* consistency, we also calculated the *intra-grader* consistency using 100 randomly selected evaluation pages that had to be graded a second time. Concerning the *intra-grader* and *inter-grader* consistencies, the κ coefficients are given in Table 11.

The obtained overall *intra-grader* κ coefficients were high. Substantial agreement coefficients were obtained for the *Ranking* metrics for all translation tasks. Concerning the human assessment in terms of *Fluency/Adequacy*, substan-

Table 11: Grader Consistency

Metric	Intra-Grader κ			Inter-Grader κ		
	DT_{EC}	DT_{CE}	BT_{*E}	DT_{EC}	DT_{CE}	BT_{*E}
Ranking	0.69	0.66	0.78	0.51	0.43	0.59
Fluency	0.61	0.54	0.75	0.34	0.27	0.47
Adequacy	0.62	0.57	0.60	0.39	0.18	0.39
Dialog	0.60	0.52	–	0.38	0.26	–

tial agreement was achieved for all BTEC tasks and moderate agreement for the DIALOG tasks.

Concerning the *inter-grader* consistency, the κ coefficients are much lower for the *Fluency/Adequacy/Dialog* metrics achieving only fair agreements for the DIALOG tasks and moderate agreement for the BTEC tasks. However, Moderate to substantial agreements were achieved for the *Ranking* metrics resulting in a high reliability of this years human assessment results.

4. Conclusions

This year’s workshop provided a testbed for verifying the quality of state-of-the-art speech-to-speech translation technologies for the translation of different communication modalities, spoken language styles and semantic context.

The standard BTEC task of IWSLT 2010 focused on the translation of frequently used utterances in the domain of travel conversations from Arabic, French, and Turkish into English. The analysis of the 29 MT system results submitted by 20 teams showed that even in a resource-limited setting, good translation performances can be achieved for the BTEC task, providing a valuable testbed to investigate new ideas for spoken language translation techniques. Due to the similarity of the source and target language, the French-English task proved to be the easiest task, achieving the highest subjective and automatic evaluation scores. For Turkish and Arabic, word segmentation issues seem to be crucial in order to deal with the significant amount of unknown words contained in this year’s testset and to achieve high quality translation performance.

The DIALOG task was a repetition of last year’s Challenge Task. The participants had to translate a collection of task-oriented dialogs in travel situations for both translation directions (Chinese-English and vice versa) using two input conditions, i.e., automatic speech recognition outputs containing recognition errors and text input without speech recognition errors. The automatic and subjective evaluation of the 11 primary run submissions of the IWSLT 2010 testset resulted in lower scores compared to last year’s testset due to the higher translation complexity of the IWSLT 2010 testset. However, the comparison of the system outputs of last year’s testset submitted by participants that took part in both the 2009 and 2010 evaluation campaigns showed an improvement in automatic evaluation scores, indicating the progress made over time on the DIALOG task. Many new techniques to improve translation quality were investigated in this year’s shared task, including the paraphrasing of training data to

reduce the data sparseness problem, the integration of multiple word segmentation schemes to reduce the problem of unknown words, the handling of ASR errors to cope with noisy input, the introduction of syntactic constraints into hierarchical system to improve grammaticality of the MT output, and system combination techniques to overcome short comes of specific machine translation approaches. In addition, the application of a new evaluation metric taking into account information beyond the current input sentence to judge the quality of a translation in the context of a dialog resulted in new insights into the requirements of the translation and evaluation of human conversations that will help to advance the current state of the art in speech-to-speech translation.

This year for the first time, we ran an evaluation exercise on the translation of talks from English to French. The TALK task was organized around a collection of recordings, transcriptions, and translations of real public speeches covering a variety of topics. In addition to domain-specific parallel data, participants could try to exploit a fairly large amount of out-of-domain training data. The task required translating both manually and automatically generated transcripts. ASR transcripts were provided by the task organizers. Nine teams participated in this first exercise and our analysis of the results confirmed for us that the proposed task is definitely sound, original, interesting, and sufficiently complex. Future work will consider improving the ASR input condition by providing automatic transcripts of better quality and in larger quantity.

5. Acknowledgements

The authors would like to thank all the people involved in the preparation of this workshop and the subjective evaluation task. In particular, we would like to thank Christian Girardi for collecting the TED training data, and Jan Niehues and Teresa Hermann for their help in collecting the TED development and evaluation data and in setting up the evaluation server. Thanks also go to Kevin Kilgour for providing the language model for the ASR system for the TED task and Hermann Ney for providing us with the RWTH sentence segmentation tool for scoring the ASR output of the TALK task. In addition, we'd like to thank Shigeki Matsuda for preparing the speech data sets and generating the ASR outputs of the DIALOG task. Special thanks to the TUBITAK-UEDIN and CEA teams, for providing us with the Turkish and French data sets and to Chris Callison-Burch for letting us use the browser-interface scripts of the subjective *Ranking* metrics. In addition, we thank all the paid experts and volunteers who carried out the human assessment of the translation outputs. We also thank the program committee members for reviewing a large number of MT system descriptions and technical paper submissions. Last, but not least, we thank all the research groups for their active participation in the IWSLT evaluation campaign and for making the IWSLT workshop a success.

6. References

- [1] M. Paul, "Overview of the IWSLT 2009 Evaluation Campaign," in *Proc. of IWSLT*, Tokyo, Japan, 2009, pp. 1–18.
- [2] E. Matusov, G. Leusch, O. Bender, and H. Ney, "Evaluating Machine Translation Output with Automatic Sentence Segmentation," in *Proc. of IWSLT*, Pittsburgh, PA, 2005, pp. 148–154.
- [3] A. Stolcke, "Srlm: an extensible language modeling toolkit," in *Proceedings of the International Conference on Spoken Language Processing*, Denver (Colorado), 2002.
- [4] S. Matsuda, T. Jitsuhiro, K. Markov, and S. Nakamura, "ATR Parallel Decoding Based Speech Recognition System Robust to Noise and Speaking Styles," *IEEE Transactions on Information and Systems*, vol. E89-D(3), pp. 989–997, 2006.
- [5] C. Callison-Burch, C. Fordyce, P. Koehn, C. Monz, and J. Schroeder, "(Meta-) Evaluation of Machine Translation," in *Proc. of the Second Workshop on SMT*. Prague, Czech Republic: ACL, 2007, pp. 136–158.
- [6] J. Blatz, E. Fitzgerald, G. Foster, S. Gandrabur, C. Goutte, A. Kulesza, A. Sanchis, and N. Ueffing, "Confidence estimation for statistical machine translation," in *Final Report of the JHU Summer Workshop*, 2003.
- [7] J. S. White, T. O'Connell, and F. O'Mara, "The ARPA MT evaluation methodologies: evolution, lessons, and future approaches," in *Proc of the AMTA*, 1994, pp. 193–205.
- [8] J. R. Landis and G. G. Koch, "The measurement of observer agreement for categorical data," *Biometrics*, vol. 33 (1), pp. 159–174, 1977.
- [9] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "BLEU: a method for automatic evaluation of machine translation," in *Proc. of the 40th ACL*, Philadelphia, USA, 2002, pp. 311–318.
- [10] G. Doddington, "Automatic evaluation of machine translation quality using n-gram co-occurrence statistics," in *Proc. of the HLT 2002*, San Diego, USA, 2002, pp. 257–258.
- [11] A. Lavie and A. Agarwal, "METEOR: An automatic metric for mt evaluation with high levels of correlation with human judgments," in *Proc. of ACL Workshop on SMT*, Prague, Czech Republic, 2007, pp. 228–231.
- [12] J. P. Turian, L. Shen, and I. D. Melamed, "Evaluation of machine translation and its evaluation," in *Proc. of the MT Summit IX*, New Orleans, USA, 2003, pp. 386–393.
- [13] S. Niessen, F. J. Och, G. Leusch, and H. Ney, "An evaluation tool for machine translation: Fast evaluation for machine translation research," in *Proc. of the 2nd LREC*, Athens, Greece, 2000, pp. 39–45.
- [14] F. J. Och, "Minimum error rate training in smt," in *Proc. of the 41st ACL*, Sapporo, Japan, 2003, pp. 160–167.

- [15] M. Snover, B. Dorr, R. Schwartz, L. Micciulla, and J. Makhoul, "A study of translation edit rate with targeted human annotation," in *Proc. of the AMTA*, Cambridge and USA, 2006, pp. 223–231.
- [16] Y. Zhang, S. Vogel, and A. Waibel, "Interpreting Bleu/NIST Scores: How Much Improvement do We Need to Have a Better System?" in *Proc of the LREC*, 2004, pp. 2051–2054.
- [17] E. Matusov and S. Köprü, "AppTek's APT Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 29–36.
- [18] H. Almaghout, J. Jiang, and A. Way, "The DCU Machine Translation Systems for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 37–44.
- [19] F. Zamora-Martinez, M. J. Castro-Bleda, and H. Schwenk, "N-gram-based Machine Translation enhanced with Neural Networks for the French-English BTEC-IWSLT'10 task," in *Proc. of IWSLT*, Paris, France, 2010, pp. 45–52.
- [20] A. Bisazza, I. Klasanis, M. Cettolo, and M. Federico, "FBK @ IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 53–58.
- [21] J. Gosme, W. Mekki, F. Debili, Y. Lepage, and N. Lucas, "The GREYC/LLACAN Machine Translation Systems for the IWSLT 2010 Campaign," in *Proc. of IWSLT*, Paris, France, 2010, pp. 59–65.
- [22] X. Duan, R. E. Banchs, J. Lang, D. Xiong, A. Aw, M. Zhang, and H. Li, "I²R Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 67–72.
- [23] H. Xiong, J. Xie, H. Yu, K. Liu, W. Luo, H. Mi, Y. Liu, Y. Lü, and Q. Liu, "The ICT Statistical Machine Translation Systems for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 73–79.
- [24] W. Ling, T. Luís, J. Graça, L. Coheur, and I. Trancoso, "The INESC-ID Machine Translation System for the IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 81–84.
- [25] G. Gascó, V. Alabau, J. Andrés-Ferrer, J. González-Rubio, M.-A. Rocha, G. Sanchis-Trilles, F. Casacuberta, J. González, and J.-A. Sánchez, "ITI-UPV system description for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 85–92.
- [26] J. Niehues, M. Mediani, T. Herrmann, M. Heck, C. Herff, and A. Waibel, "The KIT Translation system for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 93–98.
- [27] L. Besacier, H. Afli, T. N. D. Do, H. Blanchon, and M. Potet, "LIG Statistical Machine Translation Systems for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 99–104.
- [28] A. Allauzen, J. M. Crego, I. D. El-Kahlout, H.-S. Le, G. Wisniewski, and F. Yvon, "LIMSI @ IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 105–112.
- [29] A. Rousseau, L. Barrault, P. Deléglise, and Y. Estève, "LIUM's Statistical Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 113–117.
- [30] I. T. Khemakhem, S. Jamoussi, and A. B. Hammadou, "The MIRACL Arabic-English Statistical Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 119–125.
- [31] W. Shen, T. Anderson, R. Slyh, and A. R. Aminzadeh, "The MIT/LL-AFRL IWSLT-2010 MT System," in *Proc. of IWSLT*, Paris, France, 2010, pp. 127–134.
- [32] C.-H. Li, N. Duan, Y. Zhao, S. Liu, L. Cui, M.-Y. Hwang, A. Axelrod, J. Gao, Y. Zhang, and L. Deng, "The MSRA Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 135–138.
- [33] C.-L. Goh, T. Watanabe, M. Paul, A. Finch, and E. Sumita, "The NICT Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 139–146.
- [34] K. Sudoh, K. Duh, and H. Tsukada, "NTT Statistical Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 147–152.
- [35] H. Na and J.-H. Lee, "The POSTECH's Statistical Machine Translation System for the IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 153–156.
- [36] S. Yahyaei and C. Monz, "The QMUL System Description for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 157–162.
- [37] S. Mansour, S. Peitz, D. Vilar, J. Wuebker, and H. Ney, "The RWTH Aachen Machine Translation system for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 163–168.
- [38] K. Bar and N. Dershowitz, "Tel Aviv University's System Description for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 169–174.
- [39] J. Murakami, T. Nishimura, and M. Tokuhisa, "Statistical Pattern-Based Machine Translation with Statistical French-English Machine Translation," in *Proc. of IWSLT*, Paris, France, 2010, pp. 175–182.
- [40] C. Mermer, H. Kaya, and M. U. Doğan, "The TÜBİTAK-UEKAE Statistical Machine Translation System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 183–188.
- [41] C. Henríquez, M. R. Costa-jussá, V. Daudaravicius, R. E. Banchs, and J. Marino, "UPC-BMIC-VDU system description for the IWSLT 2010: testing several collocation segmentations in a phrase-based SMT system," in *Proc. of IWSLT*, Paris, France, 2010, pp. 189–195.
- [42] M. Khalilov and K. Sima'an, "The ILLC-UvA SMT System for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 197–203.
- [43] S. Martzoukos and C. Monz, "The UvA System Description for IWSLT 2010," in *Proc. of IWSLT*, Paris, France, 2010, pp. 205–208.

Appendix A. MT System Overview

Research Group	MT System Description	Type	System	Tasks
Apptek, Inc. (Turkey)	AppTek's APT Machine Translation System for IWSLT 2010 [17]	PBSMT	apptek	BT _{AE} , BT _{TE}
Carnegie Mellon University, Qatar Campus (Qatar)	Morphology-to-Syntax Alignment for Factored Phrase-based SMT	PBSMT	cmu_qatar	BT _{TE}
Dublin City University, School of Computing (Ireland)	The DCU Machine Translation Systems for IWSLT 2010 [18]	Hybrid	dcu	DT _{CE} , BT _{AE}
Universidad CEU-Cardenal Herrera & Politecnica de Valencia (Spain)	N-gram-based Machine Translation enhanced with Neural Networks for the French-English BTEC-IWSLT'10 task [19]	NBSMT	dsic-upv	BT _{FE}
Fondazione Bruno Kessler, Ricerca Scientifica e Tecnologica (Italy)	FBK @ IWSLT 2010 [20]	PBSMT	fbk	TT _{EF} , BT _{AE} , BT _{TE}
University of Caen Basse-Normandie, GREYC (France)	The GREYC/LLACAN Machine Translation Systems for the IWSLT 2010 Campaign [21]	PBSMT	greyc	BT _{AE}
Institute for Infocomm Research (Singapore)	I ² R Machine Translation System for IWSLT 2010 [22]	Hybrid	i2r	DT _{CE}
Chinese Academy of Sciences, Institute of Computing Technology (China)	The ICT Statistical Machine Translation Systems for the IWSLT 2010 [23]	Hybrid	ict	DT _{CE}
Instituto de Engenharia de Sistemas e Computadores Investigacao e Desenvolvimento (Portugal)	The INESC-ID Machine Translation System for the IWSLT 2010 [24]	PBSMT	inesc-id	DT _{CE} , BT _{FE}
Universidad Politécnica de Valencia, Instituto Universitario Mixto de Tecnología Informática (Spain)	ITI-UPV system description for IWSLT 2010 [25]	Hybrid	iti-upv	TT _{EF} , DT _{CE}
Karlsruhe Institute of Technology, interACT (Germany)	The KIT Translation system for IWSLT 2010 [26]	PBSMT	kit	TT _{EF} , BT _{FE}
University J. Fourier, GETALP, LIG (France)	LIG Statistical Machine Translation Systems for IWSLT 2010 [27]	Hybrid	lig	TT _{EF} , BT _{AE}
LIMSI-CNR (France)	LIMSI @ IWSLT 2010 [28]	PBSMT	limsi	TT _{EF} , BT _{TE}
University of Le Mans, LIUM (France)	LIUM's Statistical Machine Translation System for IWSLT 2010 [29]	PBSMT	lium	TT _{EF}
MIRACL Laboratory (Tunisia)	The MIRACL Arabic-English Statistical Machine Translation System for IWSLT 2010 [30]	PBSMT	miracl	BT _{AE}
MIT Lincoln Laboratory (USA)	The MIT/LL-AFRL IWSLT-2010 MT System [31]	PBSMT	mit	TT _{EF} , BT _{AE} , BT _{FE} , BT _{TE}
Microsoft Research Asia, Natural Language Computing (China)	The MSRA Machine Translation System for IWSLT 2010 [32]	Hybrid	msra	DT _{CE}
National Institute of Information and Communications Technology (Japan)	The NICT Translation System for IWSLT 2010 [33]	Hybrid	nict	DT _{CE} , BT _{FE}
NTT Comm. Science Labs (Japan)	NTT Statistical Machine Translation System for IWSLT 2010 [34]	PBSMT	ntt	TT _{EF}
Pohang University of Science and Technology (Korea)	The POSTECH's Statistical Machine Translation System for the IWSLT 2010 [35]	PBSMT	postech	DT _{CE}

EBMT : Example-based MT

PBSMT : Phrase-based SMT

NBSMT : Ngram-based SMT

HPSMT : Hierarchical Phrase-based SMT

Hybrid : Hybrid MT

: MT system description paper is not included in the proceedings.

Research Group	MT System Description	Type	System	Tasks
Queen Mary, University of London (United Kingdom)	The QMUL System Description for IWSLT 2010 [36]	PBSMT	qmul	BT _{AE} , BT _{FE} , BT _{TE}
Rheinisch Westfälische Technische Hochschule (Germany)	The RWTH Aachen Machine Translation system for IWSLT 2010 [37]	Hybrid	rwth	BT _{AE}
Tel Aviv University (Israel)	Tel Aviv University's System Description for IWSLT 2010 [38]	EBMT	tau	BT _{AE}
Tottori University (Japan)	Statistical Pattern-Based Machine Translation with Statistical French-English Machine Translation [39]	Hybrid	tottori	BT _{FE}
TÜBİTAK-UEKAE (Turkey)	The TÜBİTAK-UEKAE Statistical Machine Translation System for IWSLT 2010 [40]	PBSMT	tubitak	TT _{EF} , DT _{CE} , BT _{AE} , BT _{TE}
Universitat Politècnica de Catalunya (Spain)	UPC-BMIC-VDU system description for the IWSLT 2010: testing several collocation segmentations in a phrase-based SMT system [41]	PBSMT	upc	BT _{FE}
University Amsterdam, Institute for Logic Language and Computation (Netherlands)	The ILLC-UvA SMT System for IWSLT 2010 [42]	PBSMT	uva-illc	DT _{CE}
University Amsterdam, Intelligence Systems Lab (Netherlands)	The UvA System Description for IWSLT 2010 [43]	PBSMT	uva-isca	DT _{CE} , BT _{AE} , BT _{FE} , BT _{TE}

EBMT : Example-based MT

NBSMT : Ngram-based SMT

Hybrid : Hybrid MT

PBSMT : Phrase-based SMT

HPSMT : Hierarchical Phrase-based SMT

: MT system description paper is not included in the proceedings.

Appendix B. Human Assessment

B.1. Fluency / Adequacy / Dialog

(best = 5.0, ..., worst = 1.0)

- Only the top-ranked (*NormRank*) primary run submissions (cf. Appendix B.2.) were evaluated.
- *Fluency* indicates how the evaluation segment sounds to a native speaker of the target language.
- *Adequacy* indicates how much of the information from the reference translation was expressed in the MT output.
- *Dialog* is an adequacy assessment taking into account the context of the given dialog.

(testset_IWSLT10)

DIALOG	MT	Fluency	Adequacy	Dialog
DT _{CE}	ict.ASR online	2.41 1.75	2.42 1.84	2.72 2.07
	ict.CRR online	2.94 1.94	2.93 2.05	3.31 2.35
DT _{EC}	ict.ASR online	2.86 2.19	2.83 2.34	3.11 2.59
	ict.CRR online	3.61 2.41	3.74 2.62	3.93 2.88

BTEC	MT	Fluency	Adequacy
BT _{AE}	apptek online	3.43 3.28	3.48 3.56
	dsic-upv online	3.91 4.02	4.05 4.30
BT _{TE}	tubitak online	3.50 3.69	3.74 3.99

B.2. Ranking

(**Ranking**: best = 1.0, ..., worst = 0.0) (**NormRank**: best = 5.0, ..., worst = 1.0)

- The *Ranking* scores are the average numbers of times that a system was judged better than any other system.
- The *NormRank* scores are normalized ranks on a per-judge basis using the method of [6].

DIALOG

DT_{CE} (ASR)

MT	Ranking	MT	NormRank
ict	0.5928	ict	3.52
nict	0.5197	nict	3.35
i2r	0.4524	i2r	3.17
online	0.4442	msra	3.12
msra	0.4392	inesc-id	3.05
iti-upv	0.3966	online	3.05
inesc-id	0.3850	iti-upv	3.00
uva-illc	0.3788	uva-illc	2.96
postech	0.3558	tubitak	2.93
dcu	0.3439	postech	2.93
tubitak	0.3420	dcu	2.83
uva-isca	0.0736	uva-isca	2.10

DT_{EC} (ASR)

MT	Ranking	MT	NormRank
ict	0.5875	ict	3.56
i2r	0.5347	nict	3.44
nict	0.5316	i2r	3.44
msra	0.4929	msra	3.24
iti-upv	0.4730	iti-upv	3.22
postech	0.4670	postech	3.19
inesc-id	0.4670	inesc-id	3.14
online	0.4467	tubitak	3.10
tubitak	0.4296	online	3.02
dcu	0.3145	dcu	2.63
uva-illc	0.2819	uva-illc	2.52
uva-isca	0.0307	uva-isca	1.51

DT_{CE} (CRR)

MT	Ranking	MT	NormRank
ict	0.7212	ict	3.84
nict	0.5720	nict	3.43
i2r	0.5147	i2r	3.29
msra	0.5145	msra	3.26
online	0.4746	online	3.10
dcu	0.4011	inesc-id	3.00
inesc-id	0.3911	dcu	2.91
iti-upv	0.3769	iti-upv	2.89
postech	0.3284	tubitak	2.82
tubitak	0.3245	postech	2.80
uva-illc	0.2483	uva-illc	2.62
uva-isca	0.0766	uva-isca	2.02

DT_{EC} (CRR)

MT	Ranking	MT	NormRank
ict	0.7607	ict	4.07
i2r	0.5614	i2r	3.50
nict	0.5233	nict	3.38
postech	0.4980	postech	3.24
msra	0.4867	tubitak	3.21
tubitak	0.4776	msra	3.19
online	0.4577	inesc-id	2.99
inesc-id	0.4308	online	2.95
iti-upv	0.4086	iti-upv	2.90
dcu	0.3688	dcu	2.75
uva-illc	0.2986	uva-illc	2.48
uva-isca	0.0311	uva-isca	1.38

BTEC

BT_{AE}

MT	Ranking	MT	NormRank
online	0.4863	apptek	3.34
apptek	0.4485	mit	3.34
mit	0.4396	online	3.30
rwth	0.4020	rwth	3.23
qmul	0.3991	dcu	3.23
dcu	0.3889	qmul	3.22
fbk	0.3438	fbk	3.00
lig	0.3300	lig	2.91
miracl	0.2967	miracl	2.87
uva-isca	0.2588	uva-isca	2.78
tau	0.2535	tubitak	2.63
greyc	0.2529	greyc	2.58
tubitak	0.2249	tau	2.57

BT_{FE}

MT	Ranking	MT	NormRank
online	0.4114	online	3.24
tottori	0.3482	dsic-upv	3.13
kit	0.3256	kit	3.13
dsic-upv	0.3248	tottori	3.11
mit	0.3135	mit	3.09
inesc-id	0.3069	inesc-id	3.08
upc	0.3057	upc	3.08
nict	0.3046	nict	3.03
qmul	0.2794	qmul	2.94
uva-isca	0.1437	uva-isca	2.19

BT_{TE}

MT	Ranking	MT	NormRank
online	0.4437	online	3.22
tubitak	0.3378	tubitak	3.13
mit	0.3160	mit	3.05
fbk	0.3137	fbk	3.04
apptek	0.3118	apptek	3.01
limsi	0.2923	limsi	2.89
qmul	0.2724	qmul	2.87
cmu_qatar	0.2697	cmu_qatar	2.79
uva-isca	0.2432	uva-isca	2.72

B.3 Difference To System With Best Ranking Score

(best = 0.0, ..., worst = 1.0)

- The *BestRankDiff* scores are the ratio of translations that the system with the highest *Ranking* score (MT^{top}) gains to the respective system, i.e. $BestRankDiff = \frac{|translations\ ranked\ worse\ than\ MT^{top}| - |translations\ ranked\ better\ than\ MT^{top}|}{number\ of\ translations\ ranked\ together}$.
- The systems are ordered according to the *BestRankDiff* ratios.

DIALOG

DT_{CE} (ASR)

ict	BestRankDiff	Better	Same	Worse
i2r	0.1757	0.2969	0.2303	0.4727
msra	0.1788	0.3196	0.1818	0.4985
nict	0.1953	0.3007	0.2031	0.4960
postech	0.3278	0.2450	0.1821	0.5728
online	0.3333	0.2371	0.1924	0.5704
inesc-id	0.3712	0.2006	0.2274	0.5719
tubitak	0.3880	0.2276	0.1567	0.6156
uva-ille	0.4169	0.2149	0.1530	0.6319
dcu	0.4308	0.1897	0.1897	0.6205
iti-upv	0.4361	0.2021	0.1595	0.6382
uva-isca	<i>0.7927</i>	0.0493	0.1085	0.8421

DT_{EC} (ASR)

ict	BestRankDiff	Better	Same	Worse
nict	0.0574	0.3869	0.1685	0.4444
iti-upv	0.1363	0.3371	0.1893	0.4734
msra	0.1558	0.3290	0.1861	0.4848
tubitak	0.1966	0.3138	0.1757	0.5104
i2r	0.2034	0.3290	0.1385	0.5324
inesc-id	0.2321	0.3080	0.1517	0.5401
postech	0.2672	0.3017	0.1293	0.5689
online	0.3378	0.2702	0.1216	0.6081
dcu	0.3909	0.2510	0.1069	0.6419
uva-illc	0.6460	0.1150	0.1238	0.7610
uva-isca	<i>0.9033</i>	0.0210	0.0546	0.9243

DT_{CE} (CRR)

ict	BestRankDiff	Better	Same	Worse
msra	0.3464	0.2105	0.2324	0.5570
nict	0.3843	0.2313	0.1529	0.6156
i2r	0.4334	0.1931	0.1802	0.6266
inesc-id	0.4771	0.2033	0.1161	0.6804
online	0.5296	0.1857	0.0988	0.7154
iti-upv	0.5530	0.1704	0.1060	0.7234
tubitak	0.5627	0.1578	0.1214	0.7206
dcu	0.6228	0.1315	0.1140	0.7543
postech	0.6837	0.1209	0.0744	0.8046
uva-illc	0.7456	0.1052	0.0438	0.8508
uva-isca	<i>0.8915</i>	0.0283	0.0518	0.9198

DT_{EC} (CRR)

ict	BestRankDiff	Better	Same	Worse
i2r	0.3733	0.2360	0.1545	0.6094
nict	0.4553	0.2008	0.1428	0.6562
msra	0.5144	0.1893	0.1069	0.7037
postech	0.5590	0.1590	0.1227	0.7181
tubitak	0.5603	0.1767	0.0862	0.7370
online	0.5840	0.1681	0.0796	0.7522
iti-upv	0.6290	0.1169	0.1370	0.7459
inesc-id	0.6331	0.1310	0.1048	0.7641
dcu	0.7167	0.1115	0.0600	0.8283
uva-illc	0.8171	0.0622	0.0583	0.8793
uva-isca	<i>0.9723</i>	0.0000	0.0276	0.9723

BTEC

BT_{AE}

online	BestRankDiff	Better	Same	Worse
apptek	-0.0659	0.4258	0.2142	0.3598
mit	-0.0231	0.4035	0.2159	0.3804
qmul	0.0169	0.3785	0.2259	0.3954
dcu	0.0716	0.3253	0.2776	0.3970
rwth	0.0831	0.3490	0.2188	0.4321
miracl	0.1531	0.3213	0.2042	0.4744
lig	0.2159	0.3017	0.1804	0.5177
fbk	0.2682	0.2760	0.1796	0.5442
uva-isca	0.2988	0.2486	0.2039	0.5474
greyc	0.3577	0.2323	0.1775	0.5900
tubitak	0.3668	0.2100	0.2130	0.5769
tau	<i>0.4330</i>	0.1784	0.2099	0.6115

BT_{FE}

online	BestRankDiff	Better	Same	Worse
dsic-upv	0.0082	0.3476	0.2965	0.3558
mit	0.0246	0.3550	0.2653	0.3796
tottori	0.0270	0.3260	0.3209	0.3530
upc	0.0539	0.3252	0.2956	0.3791
kit	0.0630	0.3257	0.2854	0.3887
nict	0.0658	0.3309	0.2722	0.3967
qmul	0.0782	0.3435	0.2346	0.4217
inesc-id	0.0842	0.3157	0.2842	0.4000
uva-isca	<i>0.4595</i>	0.1742	0.1919	0.6338

BT_{TE}

online	BestRankDiff	Better	Same	Worse
tubitak	0.0086	0.3683	0.2547	0.3769
mit	0.0424	0.3663	0.2247	0.4088
fbk	0.0730	0.3483	0.2303	0.4213
limsi	0.1071	0.3464	0.2000	0.4535
apptek	0.1243	0.3204	0.2346	0.4448
cmu_qatar	0.1853	0.3239	0.1666	0.5093
qmul	0.1858	0.3097	0.1946	0.4955
uva-isca	<i>0.2890</i>	0.2630	0.1840	0.5520

BTEC Turkish-English (BT_{TE})

"case+punc" evaluation							CRR	"no case+no punc" evaluation						
BLEU	METEOR	WER	PER	TER	GTM	NIST		BLEU	METEOR	WER	PER	TER	GTM	NIST
57.63	81.66	31.37	25.95	24.35	78.12	8.743	tubitak	55.78	78.50	36.06	28.46	27.94	76.23	9.124
57.70	80.22	31.61	25.57	24.99	78.76	8.612	fbk	55.56	77.15	36.65	28.91	28.49	76.18	8.822
60.21	79.45	30.16	25.93	23.44	77.87	7.781	mit	59.18	76.07	34.59	29.26	26.58	75.81	7.728
52.97	76.95	35.55	29.43	27.71	74.80	7.750	limsi	50.75	73.02	41.14	33.01	31.84	72.13	7.811
52.64	77.92	36.66	29.45	27.87	74.73	7.748	apptek	50.87	74.20	42.92	32.64	31.96	73.60	7.979
53.54	75.50	36.20	30.93	29.06	73.93	7.532	qmul	50.99	71.44	42.17	35.95	32.81	71.11	7.507
49.42	77.09	38.29	31.62	30.26	74.60	7.922	online	46.71	73.17	44.59	34.88	34.78	72.89	8.178
49.06	75.78	39.42	32.30	31.40	74.03	7.880	cmu_qatar	46.48	71.17	45.88	37.02	36.20	71.35	8.090
36.07	68.28	48.91	42.21	35.70	63.97	5.899	uva-isca	38.08	62.46	53.37	43.50	41.34	64.09	6.341

Appendix D. Evaluation Metric Correlation

- The correlation between evaluation metrics are measured using the *Spearman's rank correlation coefficient* $\rho \in [-1.0, 1.0]$ with $\rho = 1.0$ if all systems ranked in same order, $\rho = -1.0$ if all systems ranked in reverse order and $\rho = 0.0$ if no correlation exists
- Z-avg is the average system score of the best z-transformed automatic evaluation metric subset obtained for the respective translation task. The z-avg scores are given for all MT systems including (w/ online) and excluding (w/o online) the online translation system).
- The automatic evaluation metrics that correlate best with the respective human assessments are marked in boldface

(testset_IWSLT10)

DT _{CE} (ASR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	METEOR	WER	PER	TER	GTM	NIST
Ranking	0.7272	(0.9090)	0.5174	0.6503	-0.7062	-0.6713	0.6433	-0.7202	0.3986
NormRank	0.8216	(0.9340)	0.6363	0.7587	-0.7937	-0.7412	0.7027	-0.7972	0.5104
BestRankDiff	0.6909	(0.7939)	0.6909	0.6181	-0.6727	-0.5181	0.4818	-0.6818	0.3727

DT _{CE} (CRR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	METEOR	WER	PER	TER	GTM	NIST
Ranking	0.7342	(0.8454)	0.6433	0.7412	-0.6923	-0.6853	0.7272	-0.6853	0.5314
NormRank	0.7762	(0.9000)	0.6783	0.7832	-0.7272	-0.7272	0.7622	-0.7202	0.5804
BestRankDiff	0.8741	(0.9727)	0.7902	0.8391	-0.8531	-0.8111	0.8111	-0.8461	0.6153

DT _{EC} (ASR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	F1	WER	PER	TER	GTM	NIST
Ranking	0.9510	(0.9636)	0.9300	0.9160	-0.9300	-0.8881	0.9160	-0.9090	0.9160
NormRank	0.9702	(0.9613)	0.9562	0.9423	-0.9458	-0.9248	0.9423	-0.9318	0.9423
BestRankDiff	0.9090	(0.9151)	0.8545	0.8909	-0.8818	-0.8909	0.8545	-0.8636	0.8818

DT _{EC} (CRR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	F1	WER	PER	TER	GTM	NIST
Ranking	0.8881	(0.9272)	0.8601	0.8461	-0.8601	-0.8601	0.8041	-0.8041	0.8391
NormRank	0.8951	(0.9090)	0.8601	0.8531	-0.8531	-0.8741	0.8181	-0.7972	0.8601
BestRankDiff	0.9230	(0.9727)	0.8951	0.8881	-0.9090	-0.8951	0.8391	-0.8601	0.8741

BT_{AE} (CRR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	METEOR	WER	PER	TER	GTM	NIST
Ranking	0.8241	(0.9370)	0.6483	0.7802	-0.6483	-0.6153	0.7582	-0.5329	0.8351
NormRank	0.9203	(0.9755)	0.8049	0.8928	-0.7994	-0.7774	0.8598	-0.7225	0.9258
BestRankDiff	0.8516	(0.9090)	0.6923	0.8461	-0.6923	-0.6648	0.7692	-0.6098	0.8626

BT_{FE} (CRR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	METEOR	WER	PER	TER	GTM	NIST
Ranking	0.5393	(0.8333)	-0.1242	0.2727	0.0393	0.1272	0.5636	0.1151	0.5515
NormRank	0.5090	(0.7916)	-0.0242	0.3212	-0.0333	0.0606	0.5272	0.0242	0.4787
BestRankDiff	0.5151	(0.8833)	0.2818	0.4909	-0.3060	-0.2242	0.4303	-0.2969	0.5272

BT_{TE} (CRR)	z-avg		single metrics						
	(w/ online)	(w/o online)	BLEU	METEOR	WER	PER	TER	GTM	NIST
Ranking	0.8333	(1.0000)	0.4523	0.7380	-0.5238	-0.6250	0.7857	-0.4523	0.8095
NormRank	0.8333	(1.0000)	0.4523	0.7380	-0.5238	-0.6250	0.7857	-0.4523	0.8095
BestRankDiff	0.7857	(0.9285)	0.3333	0.6666	-0.3809	-0.5654	0.7380	-0.3333	0.8095

(Z-avg Metric Combinations)

Task	Ranking	NormRank	BestRankDiff
DT_{CE} (ASR) (CRR)	WER, TER, NIST BLEU, METEOR	METEOR, TER BLEU, METEOR	WER, TER METEOR, WER, PER, TER
DT_{EC} (ASR) (CRR)	F1, NIST F1, NIST	F1, NIST F1, PER	F1, PER BLEU, NIST
BT_{AE} (CRR)	GTM, NIST	GTM, NIST	BLEU, METEOR, NIST
BT_{FE} (CRR)	METEOR, GTM	METEOR, GTM	METEOR, GTM
BT_{TE} (CRR)	METEOR, NIST	METEOR, NIST	METEOR, NIST