

# Results of the WMT18 Metrics Shared Task

**Qingsong Ma**  
Tencent-MIG  
AI Evaluation & Test Lab  
qingsong.mqs@gmail.com

**Ondřej Bojar**  
Charles University  
MFF ÚFAL  
bojar@ufal.mff.cuni.cz

**Yvette Graham**  
Dublin City University  
ADAPT  
graham.yvette@gmail.com

## Abstract

This paper presents the results of the WMT18 Metrics Shared Task. We asked participants of this task to score the outputs of the MT systems involved in the WMT18 News Translation Task with automatic metrics. We collected scores of 10 metrics and 8 research groups. In addition to that, we computed scores of 8 standard metrics (BLEU, SentBLEU, chrF, NIST, WER, PER, TER and CDER) as baselines. The collected scores were evaluated in terms of system-level correlation (how well each metric's scores correlate with WMT18 official manual ranking of systems) and in terms of segment-level correlation (how often a metric agrees with humans in judging the quality of a particular sentence relative to alternate outputs). This year, we employ a single kind of manual evaluation: direct assessment (DA).

## 1 Introduction

Accurate machine translation (MT) evaluation is important for measuring improvements in system performance. Human evaluation can be costly and time consuming, and it is not always available for the language pair of interest. Automatic metrics can be employed as a substitute for human evaluation in such cases, metrics that aim to measure improvements to systems quickly and at no cost to developers. In the usual set-up, an automatic metric carries out a comparison of MT system output translations and human-produced reference translations to produce a single overall

score for the system.<sup>1</sup> Since there exists a large number of possible approaches to producing quality scores for translations, it is sensible to carry out a meta-evaluation of metrics with the aim to estimate their accuracy as a substitute for human assessment of translation quality. The Metrics Shared Task<sup>2</sup> of WMT annually evaluates the performance of automatic machine translation metrics in their ability to provide a substitute for human assessment of translation quality.

Again, we keep the two main types of metric evaluation unchanged from the previous years. In *system-level* evaluation, each metric provides a quality score for the whole translated test set (usually a set of documents, in fact). In *segment-level* evaluation, a score is assigned by a given metric to every individual sentence.

The underlying texts and MT systems come from the News Translation Task (Bojar et al., 2018, denoted as Findings 2018 in the following). The texts were drawn from the news domain and involve translations to/from Chinese (zh), Czech (cs), German (de), Estonian (et), Finnish (fi), Russian (ru), and Turkish (tr), each paired with English, making a total of 14 language pairs.

A single form of golden truth of translation quality judgement is used this year:

- In *Direct Assessment* (DA) (Graham et al., 2016), humans assess the quality of a given MT output translation by comparison with a reference translation (as opposed to the source and reference). DA is the new standard used in WMT

<sup>1</sup>The availability of a reference translation is the key difference between our task and *MT quality estimation*, where no reference is assumed.

<sup>2</sup><http://www.statmt.org/wmt18/metrics-task.html>, starting with Koehn and Monz (2006) up to Bojar et al. (2017)

News Translation Task evaluation, requiring only monolingual evaluators.

As in last year’s evaluation, the official method of manual evaluation of MT outputs is no longer “relative ranking” (RR, evaluating up to five system outputs on an annotation screen relative to each other) as this was changed in 2017 to DA. For system-level evaluation, we thus use the Pearson correlation  $r$  of automatic metrics with DA scores. For segment-level evaluation, we re-interpret DA judgements as relative comparisons and use Kendall’s  $\tau$  as a substitute, see below for details and references.

Section 2 describes our datasets, i.e. the sets of underlying sentences, system outputs, human judgements of translation quality and also participating metrics. Sections 3.1 and 3.2 then provide the results of system and segment-level metric evaluation, respectively. We discuss the results in Section 4.

## 2 Data

This year, we provided the task participants with one test set along with reference translations and outputs of MT systems. Participants were free to choose which language pairs they wanted to participate and whether they reported system-level, segment-level scores or both.

### 2.1 Test Sets

We use the following test set, i.e. a set of source sentences and reference translations:

**newstest2018** is the test set used in WMT18 News Translation Task (see Findings 2018), with approximately 3,000 sentences for each translation direction (except Chinese and Estonian which have 3,981 and 2,000 sentences, resp.). **newstest2018** includes a single reference translation for each direction.

### 2.2 Translation Systems

The results of the Metrics Task are likely affected by the actual set of MT systems participating in a given translation direction. For instance, if all of the systems perform similarly, it will be more difficult, even for humans, to distinguish between the quality of

translations. If the task includes a wide range of systems of varying quality, however, or systems are quite different in nature, this could in some way make the task easier for metrics, with metrics that are more sensitive to certain aspects of MT output performing better.

This year, the MT systems included in the Metrics Task were:

**News Task Systems** are machine translation systems participating in the WMT18 News Translation Task (see Findings 2018).<sup>3</sup>

**Hybrid Systems** are created automatically with the aim of providing a larger set of systems against which to evaluate metrics, as in Graham and Liu (2016). Hybrid systems were created for newstest2018 by randomly selecting a pair of MT systems from all systems taking part in that language pair and producing a single output document by randomly selecting sentences from either of the two systems. In short, we create 10K hybrid MT systems for each language pair.

Excluding the hybrid systems, we ended up with 149 systems across 14 language pairs.

### 2.3 Manual MT Quality Judgments

Direct Assessment (DA) was employed as the “golden truth” to evaluate metrics again this year. The details of this method of human evaluation is provided in two sections for system-level evaluation (Section 2.3.1) and segment-level evaluation (Section 2.3.2).

The DA manual judgements were provided by MT researchers taking part in WMT tasks, a number of in-house human evaluators at Amazon and crowd-sourced workers on Amazon Mechanical Turk.<sup>4</sup> Only judgements from workers who passed DA’s quality control mechanism were included in the final datasets used to compute system and segment-level scores employed as a gold standard in the Metrics Task.

<sup>3</sup>One system for tr-en was unfortunately omitted from the first run of human evaluation in the News Translation Task and due to time constraints was subsequently omitted from the Metrics Task evaluation, Alibaba-Ensemble.

<sup>4</sup><https://www.mturk.com>

### 2.3.1 System-level Manual Quality Judgments

In the system-level evaluation, the goal is to assess the quality of translation of an MT system for the whole test set. Our manual scoring method, DA, nevertheless proceeds sentence by sentence, aggregating the final score as described below.

**Direct Assessment (DA)** This year the translation task employed monolingual direct assessment (DA) of translation adequacy (Graham et al., 2013; Graham et al., 2014a; Graham et al., 2016). Since sufficient levels of agreement in human assessment of translation quality are difficult to achieve, the DA setup simplifies the task of translation assessment (conventionally a bilingual task) into a simpler monolingual assessment. In addition, DA avoids bias that has been problematic in previous evaluations introduced by assessment of several alternate translations on a single screen, where scores for translations had been unfairly penalized if often compared to high quality translations (Bojar et al., 2011). DA therefore employs assessment of individual translations in isolation from other outputs.

Translation adequacy is structured as a monolingual assessment of similarity of meaning where the target language reference translation and the MT output are displayed to the human assessor. Assessors rate a given translation by how adequately it expresses the meaning of the reference translation on an analogue scale corresponding to an underlying 0-100 rating scale.<sup>5</sup>

Large numbers of DA human assessments of translations for all 14 language pairs included in the News Translation Task were collected from researchers and from workers on Amazon’s Mechanical Turk, via sets of 100-translation hits to ensure sufficient repeat assessments per worker, before application of strict quality control measures to filter out assessments from poor performers.

In order to iron out differences in scoring strategies attributed to distinct human assessors, human assessment scores for translations were standardized according to an indi-

<sup>5</sup>The only numbering displayed on the rating scale are extreme points 0 and 100%, and three ticks indicate the levels of 25, 50 and 75%.

vidual judge’s overall mean and standard deviation score. Final scores for MT systems were computed by firstly taking the average of scores for individual translations in the test set (since some were assessed more than once), before combining all scores for translations attributed to a given MT system into its overall adequacy score. The gold standard for system-level DA evaluation is thus what is denoted “Ave z” in Findings 2018 (Bojar et al., 2018).

Finally, although it was necessary to apply a sentence length restriction in WMT human evaluation prior to the introduction of DA, the simplified DA setup does not require restriction of the evaluation in this respect and no sentence length restriction was applied in DA WMT18.

### 2.3.2 Segment-level Manual Quality Judgments

Segment-level metrics have been evaluated against DA annotations for the newstest2018 test set. This year, a standard segment-level DA evaluation of metrics, where each translation is assessed a minimum of 15 times, was unfortunately not possible due to insufficient number of judgements collected. DA judgements were therefore converted to relative ranking judgements (daRR) to produce results. This is the same strategy as carried out for some out-of-English language pairs in last year’s evaluation.

**daRR** When we have at least two DA scores for translations of the same source input, it is possible to convert those DA scores into a relative ranking judgement, if the difference in DA scores allows conclusion that one translation is better than the other. In the following, we will denote these re-interpreted DA judgements as “daRR”, to distinguish it clearly from the “RR” golden truth used in the past years.

Since the analogue rating scale employed by DA is marked at the 0-25-50-75-100 points, the difference in DA scores we employ to distinguish translations that are better/worse than one another is 25 points. Note that we rely on judgements collected from known-reliable volunteers and crowd-sourced workers who passed DA’s quality control mechanism. Any inconsistency that could arise from re-

	DA>1	Ave	DA pairs	DARR
cs-en	2,491	3.6	13,223	5,110
de-en	2,995	11.4	192,702	77,811
en-cs	1,586	4.9	15,311	5,413
en-de	2,150	5.3	47,041	19,711
en-et	1,035	13.6	90,755	32,202
en-fi	1,481	5.3	30,613	9,809
en-ru	2,954	6.2	54,260	22,181
en-tr	707	3.4	4,750	1,358
en-zh	3,915	6.5	86,286	28,602
et-en	2,000	11.2	118,066	56,721
fi-en	2,972	5.4	39,127	15,648
ru-en	2,916	4.9	31,361	10,404
tr-en	2,991	4.5	24,325	8,525
zh-en	3,952	7.2	97,474	33,357

Table 1: Number of judgements for DA converted to DARR data; “DA>1” is the number of source input sentences in the manual evaluation where at least two translations of that same source input segment received a DA judgement; “Ave” is the average number of translations with at least one DA judgement available for the same source input sentence; “DA pairs” is the number of all possible pairs of translations of the same source input resulting from “DA>1”; and “DARR” is the number of DA pairs with an absolute difference in DA scores greater than the 25 percentage point margin.

liance on DA judgements collected from low quality crowd-sourcing, for example, is thus prevented.

From the complete set of human assessments collected for the News Translation Task, all possible pairs of DA judgements attributed to distinct translations of the same source were converted into DARR better/worse judgements. Distinct translations of the same source input whose DA scores fell within 25 percentage points (which could have been deemed equal quality) were omitted from the evaluation of segment-level metrics. Conversion of scores in this way produced a large set of DARR judgements for all language pairs, shown in Table 1 due to combinatorial advantage of extracting DARR judgements from all possible pairs of translations of the same source input.

**Kendall’s Tau-like Formulation for daRR** We measure the quality of metrics’ segment-level scores against the DARR golden truth using a Kendall’s Tau-like formulation, which is an adaptation of the conventional Kendall’s Tau coefficient. Since we do not have a total order ranking of all translations we use to evaluate metrics, it is not possible to apply conventional Kendall’s Tau given the current DARR human evaluation setup (Graham et al., 2015). Vazquez-Alvarez and Huckvale (2002) also note that a genuine pairwise comparison is likely to lead to more stable results for segment-level metric evaluation.

Our Kendall’s Tau-like formulation,  $\tau$ , is as follows:

$$\tau = \frac{|Concordant| - |Discordant|}{|Concordant| + |Discordant|} \quad (1)$$

where *Concordant* is the set of all human comparisons for which a given metric suggests the same order and *Discordant* is the set of all human comparisons for which a given metric disagrees. The formula is not specific with respect to ties, i.e. cases where the annotation says that the two outputs are equally good.

The way in which ties (both in human and metric judgement) were incorporated in computing Kendall  $\tau$  has changed across the years of WMT Metrics Tasks. Here we adopt the version used in the last years’ WMT17 DARR evaluation (but not earlier). For a detailed discussion on other options, see also Macháček and Bojar (2014).

Whether or not a given comparison of a pair of distinct translations of the same source input,  $s_1$  and  $s_2$ , is counted as a concordant (Conc) or discordant (Disc) pair is defined by the following matrix:

		Metric		
		$s_1 < s_2$	$s_1 = s_2$	$s_1 > s_2$
Human	$s_1 < s_2$	Conc	Disc	Disc
	$s_1 = s_2$	—	—	—
	$s_1 > s_2$	Disc	Disc	Conc

In the notation of Macháček and Bojar (2014), this corresponds to the setup used in WMT12 (with a different underlying method of manual judgements, RR):

WMT12		Metric		
		<	=	>
Human	<	1	-1	-1
	=	X	X	X
	>	-1	-1	1

The key differences between the evaluation used in WMT14–WMT16 and evaluation used in WMT17 and WMT18 are (1) the move from RR to daRR and (2) the treatment of ties.<sup>6</sup> In the years 2014–2016, ties in metrics scores were not penalized. With the move to daRR, where the quality of the two candidate translations is deemed substantially different and no ties in human judgements arise, it makes sense to penalize ties in metrics’ predictions in order to promote discerning metrics.

Note that the penalization of ties makes our evaluation asymmetric, dependent on whether the metric predicted the tie for a pair where humans predicted < or >. It is now important to interpret the meaning of the comparison identically for humans and metrics. For error metrics, we thus reverse the sign of the metric score prior to the comparison with human scores: higher scores have to indicate better translation quality. In WMT18, we did this for ITER and the original authors did this for CHARACTER.

To summarize, the WMT18 Metrics Task for segment-level evaluation:

- excludes all human ties (this is already implied by the construction of daRR from DA judgements),
- counts metric’s ties as a *Discordant* pairs,
- ensures that error metrics are first converted to the same orientation as the human judgements, i.e. higher score indicating higher translation quality.

We employ bootstrap resampling (Koehn, 2004; Graham et al., 2014b) to estimate confidence intervals for our Kendall’s Tau formulation, and metrics with non-overlapping 95% confidence intervals are identified as having statistically significant difference in performance.

<sup>6</sup>Due to an error in the write-up for WMT17 (errata to follow), this second change was not properly reflected in the paper, only in the evaluation scripts.

## 2.4 Participants of the Metrics Shared Task

Table 2 lists the participants of the WMT18 Shared Metrics Task, along with their metrics. We have collected 10 metrics from a total of 8 research groups.

The following subsections provide a brief summary of all the metrics that participated. The list is concluded by our baseline metrics in Section 2.4.9.

As in last year’s task, we asked participants whose metrics are publicly available to provide links to where the code can be accessed. Table 3 provides links for metrics that participated in WMT18 that are publicly available for download.

### 2.4.1 BEER

BEER (Stanojević and Sima’an, 2015) is a trained evaluation metric with a linear model that combines features sub-word feature indicators (character n-grams) and global word order features (skip bigrams) to get language agnostic and fast to compute evaluation metric. BEER has participated in previous years of the evaluation task.

### 2.4.2 Blend

BLEND incorporates existing metrics to form an effective combined metric, employing SVM regression for training and DA scores as the gold standard. For to-English language pairs, incorporated metrics include 25 lexical based metrics and 4 other metrics. Since some lexical based metrics are simply different variants of the same metric, there are only 9 kinds of lexical based metrics, namely BLEU, NIST, GTM, METEOR, ROUGE, Ol, WER, TER and PER. 4 other metrics are CHARACTER, BEER, DPMF and ENTF.

BLEND has participated in the Metrics Task in WMT17. This year, BLEND follows its setup in WMT17, but enlarges the training data since there are some data available in WMT17. For to-English language pairs, there are 9280 sentences as training data, while 1620 sentences for English-Russian (en-ru). Experiments show the performance of BLEND can be improved if the training data increases.

BLEND is flexible to be applied to any language pairs if incorporated metrics support the

Metric	Seg-level	Sys-level	Hybrids	Participant
BEER	•	⊙	⊙	ILLC – University of Amsterdam (Stanojević and Sima’an, 2015)
BLEND	•	⊙	⊙	Tencent-MIG-AI Evaluation & Test Lab (Ma et al., 2017)
CHARACTER	•	•	•	RWTH Aachen University (Wang et al., 2016a)
ITER	•	•	★	Jadavpur University (Panja and Naskar, 2018)
METEOR++	•	⊙	⊙	Peking University (Guo et al., 2018)
RUSE	•	⊙	⊙	Tokyo Metropolitan University (Shimanaka et al., 2018)
UHH_TSKM	•	⊙	⊙	(Duma and Menzel, 2017)
YISI-*	•	⊙	⊙	NRC (Lo, 2018)

Table 2: Participants of WMT18 Metrics Shared Task. “•” denotes that the metric took part in (some of the language pairs) of the segment- and/or system-level evaluation and whether hybrid systems were also scored. “⊙” indicates that the system-level and hybrids are implied, simply taking arithmetic average of segment-level scores. “★” indicates that the original ITER system-level scores should be calculated as the *micro-average* of segment-level scores but we calculate them as simple macro-averaged for the hybrid systems. See the ITER paper for more details.

BEER	<a href="https://github.com/stanojevic/beer">https://github.com/stanojevic/beer</a>
BLEND	<a href="https://github.com/qingsongma/blend">https://github.com/qingsongma/blend</a>
CHARACTER	<a href="https://github.com/rwth-i6/Character">https://github.com/rwth-i6/Character</a>
RUSE	<a href="https://github.com/Shi-ma/RUSE">https://github.com/Shi-ma/RUSE</a>
YISI-0, incl. -1 and -1_srl	<a href="http://chikiu-jackie-lo.org/home/index.php/yisi">http://chikiu-jackie-lo.org/home/index.php/yisi</a>
Baselines:	<a href="http://github.com/moses-smt/mosesdecoder">http://github.com/moses-smt/mosesdecoder</a>
BLEU, NIST	<code>scripts/generic/mteval-v13a.pl</code>
CDER, PER, TER, WER	<code>mert/evaluator</code> (“Moses scorer”)
SENTBLEU	<code>mert/sentence-bleu</code>
CHRf, CHRf+	<a href="https://github.com/m-popovic/chrF">https://github.com/m-popovic/chrF</a>

Table 3: Metrics available for public download that participated in WMT18. Most of the baseline metrics are available with Moses, relative paths are listed.

specific language pair and DA scores are available.

### 2.4.3 CharacTer

CHARACTER (Wang et al., 2016b; Wang et al., 2016a), identical to the 2016 setup, is a character-level metric inspired by the commonly applied translation edit rate (TER). It is defined as the minimum number of character edits required to adjust a hypothesis, until it completely matches the reference, normalized by the length of the hypothesis sentence. CHARACTER calculates the character-level edit distance while performing the shift edit on word level. Unlike the strict matching criterion in TER, a hypothesis word is considered to match a reference word and could be shifted, if the edit distance between them is below a threshold value. The Levenshtein distance between the reference and

the shifted hypothesis sequence is computed on the character level. In addition, the lengths of hypothesis sequences instead of reference sequences are used for normalizing the edit distance, which effectively counters the issue that shorter translations normally achieve lower TER.

Similarly to other character-level metrics, CHARACTER is applied to non-tokenized outputs and references, which also holds for this year’s submission.

This year tokenization was carried out for en-ru hypotheses and reference before calculating the scores, since this results in large improvements in terms of correlations. For other language pairs a tokenizer was not used for pre-processing. A python library was used for calculating the Levenshtein distance, so that the metric is now about 7 times faster than before.

#### 2.4.4 ITER

ITER (Panja and Naskar, 2018) is an improved Translation Edit/Error Rate (TER) metric. In addition to the basic edit operations in TER (insertion, deletion, substitution and shift), ITER also allows stem matching and uses optimizable edit costs and better normalization.

Note that for segment-level evaluation, we reverse the sign of the score, so that better translations get higher scores. For system-level confidence, we calculate the system-level scores for hybrids systems slightly differently than the original ITER definition would require. We use the unweighted arithmetic average of segment-level scores (macro-average) whereas ITER would use the micro-average.

#### 2.4.5 meteor++

METEOR++ (Guo et al., 2018) is metric based on Meteor (Denkowski and Lavie, 2014), adding explicating treatment of “copy-words”, i.e. words that are likely to be preserved across all paraphrases of a sentence in a given language.

#### 2.4.6 RUSE

RUSE (Shimanaka et al., 2018) is a perception regressor based on three types of sentence embeddings: Infersent, Quick-Thought and Universal Sentence Encoder, designed with the aim to utilize global sentence information that cannot be captured by local features based on character or word n-grams. The sentence embeddings come from pre-trained models and the regression itself is trained on past manual judgements in WMT shared tasks.

#### 2.4.7 UHH\_TSKM

UHH\_TSKM (Duma and Menzel, 2017) is a non-trained metric utilizing kernel functions, i.e. methods for efficient calculation of overlap of substructures between the candidate and the reference translations. The metric uses both sequence kernels, applied on the tokenized input data, together with tree kernels, that exploit the syntactic structure of the sentences. Optionally, the match can also be performed for the candidate and a pseudo-reference (i.e. a translation by another MT system) or for the source sentence and the

candidate back-translated into the source language.

#### 2.4.8 YiSi-0, YiSi-1 and YiSi-1\_srl

The YiSi metrics (Lo, 2018) are recently proposed semantic MT evaluation metrics inspired by MEANT\_2.0 (?). Specifically, YiSi-1 is identical to MEANT\_2.0-NOSRL which featured in the WMT17 Metrics Task.

YiSi-1 also successfully served in the parallel corpus filtering task. Some details are provided in the system description paper (?).

YiSi-1 measures the relative lexical semantic similarity (weighted word embeddings cosine similarity aggregated into  $n$ -grams similarity) of the candidate and reference translations, optionally taking the shallow semantic structure (“srl”) into account. YiSi-0 is a degenerate resource-free version using the longest common character substring, instead of word embeddings cosine similarity, to measure the word similarity of the candidate and reference translations.

#### 2.4.9 Baseline Metrics

As mentioned by Bojar et al. (2016), Metrics Task occasionally suffers from “loss of knowledge” when successful metrics participate only in one year.

We attempt to avoid this by regularly evaluating also a range of “baseline metrics”:

- **Mteval.** The metrics BLEU (Papineni et al., 2002) and NIST (Dodgington, 2002) were computed using the script `mteval-v13a.pl`<sup>7</sup> that is used in the OpenMT Evaluation Campaign and includes its own tokenization. We run `mteval` with the flag `--international-tokenization` since it performs slightly better (Macháček and Bojar, 2013).
- **Moses Scorer.** The metrics TER (Snover et al., 2006), WER, PER and CDER (Leusch et al., 2006) were produced by the Moses scorer, which is used in Moses model optimization. To tokenize the sentences, we used the standard tokenizer script as available in Moses toolkit. When tokenizing, we also convert all outputs to lowercase.

<sup>7</sup><http://www.itl.nist.gov/iad/mig/tools/>

Since Moses scorer is versioned on Github, we strongly encourage authors of high-performing metrics to add them to Moses scorer, as this will ensure that their metric can be easily included in future tasks.

- **SentBLEU.** The metric SENTBLEU is computed using the script `sentence-bleu`, a part of the Moses toolkit. It is a smoothed version of BLEU that correlates better with human judgements for segment-level. Standard Moses tokenizer is used for tokenization.
- **chrF** The metrics CHRF and CHRF+ (Popović, 2015; Popović, 2017) are computed using their original Python implementation.

We run `chrF++.py` with the parameters `-nw 0 -b 3` to obtain the **chrF** score and with `-nw 0 -b 1` to obtain the CHRF+ score. Note that CHRF intentionally removes all spaces before matching the  $n$ -grams, detokenizing the segments but also concatenating words.

We originally planned to use the CHRF implementation which was recently made available in Moses Scorer but it mishandles Unicode characters for now.

The baselines serve in system and segment-level evaluations as customary: BLEU, TER, WER, PER and CDER for system-level only; SENTBLEU for segment-level only and CHRF for both.

**Chinese word segmentation** is unfortunately not supported by the tokenization scripts mentioned above. For scoring Chinese with baseline metrics, we thus preprocessed MT outputs and reference translations with the script `tokenizeChinese.py`<sup>8</sup> by Shujian Huang, which separates Chinese characters from each other and also from non-Chinese parts.

For computing system-level and segment-level scores, the same scripts were employed as in last year’s Metrics Task as well as for generation of hybrid systems from the given hybrid descriptions.

<sup>8</sup><http://hdl.handle.net/11346/WMT17-TVXH>

### 3 Results

We discuss system-level results for news task systems in Section 3.1. The segment-level results are in Section 3.2.

#### 3.1 System-Level Results

As in previous years, we employ the Pearson correlation ( $r$ ) as the main evaluation measure for system-level metrics. The Pearson correlation is as follows:

$$r = \frac{\sum_{i=1}^n (H_i - \bar{H})(M_i - \bar{M})}{\sqrt{\sum_{i=1}^n (H_i - \bar{H})^2} \sqrt{\sum_{i=1}^n (M_i - \bar{M})^2}} \quad (2)$$

where  $H_i$  are human assessment scores of all systems in a given translation direction,  $M_i$  are the corresponding scores as predicted by a given metric.  $\bar{H}$  and  $\bar{M}$  are their means respectively.

Since some metrics, such as BLEU, for example, aim to achieve a strong positive correlation with human assessment, while error metrics, such as TER aim for a strong negative correlation, after computation of  $r$  for metrics, we compare metrics via the absolute value of a given metric’s correlation with human assessment.

Table 4 provides the system-level correlations of metrics evaluating translation of newstest2018 into English while Table 5 provides the same for out-of-English language pairs. The underlying texts are part of the WMT18 News Translation test set (newstest2018) and the underlying MT systems are all MT systems participating in the WMT18 News Translation Task with the exception of a single tr-en system not included in the initial human evaluation run.

As recommended by Graham and Baldwin (2014), we employ Williams significance test (Williams, 1959) to identify differences in correlation that are statistically significant. Williams test is a test of significance of a difference in dependent correlations and therefore suitable for evaluation of metrics. Correlations not significantly outperformed by any other metric for the given language pair are highlighted in bold in Tables 4 and 5.

Since pairwise comparisons of metrics may be also of interest, e.g. to learn which metrics



	cs-en	de-en	et-en	fi-en	ru-en	tr-en	zh-en
$n$	5	16	14	9	8	5	14
Correlation	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $
BEER	<b>0.958</b>	0.994	<b>0.985</b>	<b>0.991</b>	0.982	0.870	<b>0.976</b>
BLEND	<b>0.973</b>	0.991	0.985	<b>0.994</b>	<b>0.993</b>	<b>0.801</b>	<b>0.976</b>
BLEU	<b>0.970</b>	0.971	<b>0.986</b>	0.973	0.979	<b>0.657</b>	<b>0.978</b>
CDER	<b>0.972</b>	0.980	<b>0.990</b>	0.984	0.980	<b>0.664</b>	<b>0.982</b>
CHARACTER	<b>0.970</b>	<b>0.993</b>	0.979	0.989	<b>0.991</b>	<b>0.782</b>	0.950
ITER	<b>0.975</b>	0.990	0.975	<b>0.996</b>	0.937	<b>0.861</b>	<b>0.980</b>
METEOR++	<b>0.945</b>	0.991	0.978	0.971	<b>0.995</b>	0.864	0.962
NIST	<b>0.954</b>	0.984	0.983	0.975	0.973	<b>0.970</b>	0.968
PER	<b>0.970</b>	0.985	<b>0.983</b>	<b>0.993</b>	0.967	0.159	0.931
RUSE	<b>0.981</b>	<b>0.997</b>	<b>0.990</b>	<b>0.991</b>	<b>0.988</b>	<b>0.853</b>	<b>0.981</b>
TER	<b>0.950</b>	0.970	<b>0.990</b>	0.968	0.970	0.533	<b>0.975</b>
UHH_TSKM	0.952	0.980	<b>0.989</b>	0.982	0.980	0.547	<b>0.981</b>
WER	<b>0.951</b>	0.961	<b>0.991</b>	0.961	0.968	0.041	<b>0.975</b>
YiSi-0	0.956	<b>0.994</b>	0.975	0.978	<b>0.988</b>	<b>0.954</b>	0.957
YiSi-1	<b>0.950</b>	0.992	0.979	0.973	<b>0.991</b>	<b>0.958</b>	0.951
YiSi-1_SRL	<b>0.965</b>	<b>0.995</b>	0.981	0.977	<b>0.992</b>	<b>0.869</b>	0.962

newstest2018

Table 4: Absolute Pearson correlation of to-English system-level metrics with DA human assessment in newstest2018; correlations of metrics not significantly outperformed by any other for that language pair are highlighted in bold; ensemble metrics are highlighted in gray.

	en-cs	en-de	en-et	en-fi	en-ru	en-tr	en-zh
$n$	5	16	14	12	9	8	14
Correlation	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $
BEER	<b>0.992</b>	<b>0.991</b>	<b>0.980</b>	<b>0.961</b>	<b>0.988</b>	<b>0.965</b>	0.928
BLEND	–	–	–	–	<b>0.988</b>	–	–
BLEU	0.995	0.981	0.975	<b>0.962</b>	0.983	0.826	0.947
CDER	0.997	0.986	<b>0.984</b>	<b>0.964</b>	<b>0.984</b>	0.861	0.961
CHARACTER	<b>0.993</b>	<b>0.989</b>	0.956	<b>0.974</b>	0.983	0.833	<b>0.983</b>
ITER	0.915	<b>0.984</b>	<b>0.981</b>	<b>0.973</b>	0.975	0.865	–
NIST	<b>0.999</b>	0.986	<b>0.983</b>	0.949	<b>0.990</b>	0.902	0.950
PER	0.991	0.981	0.958	0.906	<b>0.988</b>	0.859	0.964
TER	<b>0.997</b>	<b>0.988</b>	<b>0.981</b>	<b>0.942</b>	<b>0.987</b>	0.867	<b>0.963</b>
WER	<b>0.997</b>	<b>0.986</b>	<b>0.981</b>	<b>0.945</b>	0.985	0.853	0.957
YiSi-0	0.973	0.985	0.968	0.944	<b>0.990</b>	<b>0.990</b>	0.957
YiSi-1	<b>0.987</b>	0.985	<b>0.979</b>	0.940	<b>0.992</b>	<b>0.976</b>	0.963
YiSi-1_SRL	–	<b>0.990</b>	–	–	–	–	0.952

newstest2018

Table 5: Absolute Pearson correlation of out-of-English system-level metrics with DA human assessment in newstest2018; correlations of metrics not significantly outperformed by any other for that language pair are highlighted in bold; ensemble metrics are highlighted in gray.

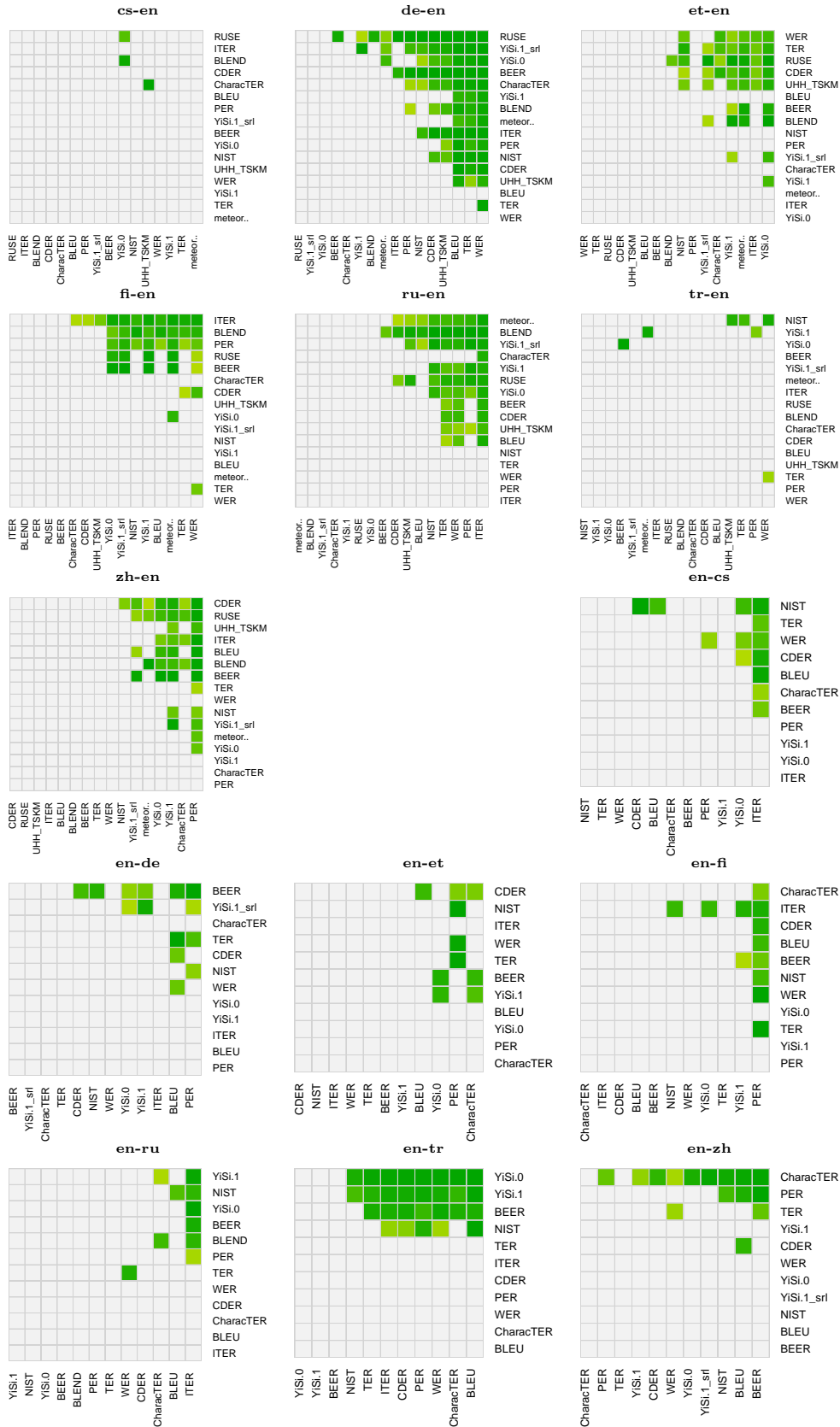


Figure 1: System-level metric significance test results for DA human assessment in newstest2018; green cells denote a statistically significant increase in correlation with human assessment for the metric in a given row over the metric in a given column according to Williams test.

	<b>cs-en</b>	<b>de-en</b>	<b>et-en</b>	<b>fi-en</b>	<b>ru-en</b>	<b>tr-en</b>	<b>zh-en</b>
$n$	10K	10K	10K	10K	10K	10K	10K
Correlation	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $
BEER	0.9497	0.9927	0.9831	0.9824	0.9755	0.7234	0.9677
BLEND	0.9646	0.9904	0.9820	0.9853	0.9865	0.7243	0.9686
BLEU	0.9557	0.9690	0.9812	0.9618	0.9719	0.5862	0.9684
CDER	0.9642	0.9797	0.9876	0.9764	0.9739	0.5767	0.9733
CHARACTER	0.9595	0.9919	0.9754	0.9791	0.9841	0.6798	0.9424
ITER	0.9656	0.9904	0.9746	<b>0.9885</b>	0.9429	0.7420	<b>0.9780</b>
METEOR++	0.9367	0.9898	0.9753	0.9621	<b>0.9892</b>	0.7871	0.9541
NIST	0.9419	0.9816	0.9804	0.9655	0.9650	0.8622	0.9589
PER	0.9369	0.9820	0.9782	0.9834	0.9550	0.0433	0.9233
RUSE	<b>0.9736</b>	<b>0.9959</b>	0.9879	0.9829	0.9820	0.7796	0.9734
TER	0.9419	0.9699	0.9882	0.9599	0.9635	0.4495	0.9670
UHH_TSKM	0.9429	0.9794	0.9869	0.9738	0.9734	0.4433	0.9717
WER	0.9420	0.9612	<b>0.9892</b>	0.9534	0.9618	0.0720	0.9667
YiSi-0	0.9465	0.9925	0.9719	0.9694	0.9817	0.8629	0.9495
YiSi-1	0.9425	0.9909	0.9758	0.9641	0.9846	<b>0.8810</b>	0.9429
YiSi-1_SRL	0.9565	0.9940	0.9783	0.9682	0.9860	0.7850	0.9540
<b>newstest2018 Hybrids</b>							

Table 6: Absolute Pearson correlation of to-English system-level metrics with DA human assessment for 10K hybrid super-sampled systems in newstest2018; ensemble metrics are highlighted in gray.

	<b>en-cs</b>	<b>en-de</b>	<b>en-et</b>	<b>en-fi</b>	<b>en-ru</b>	<b>en-tr</b>	<b>en-zh</b>
$n$	10K	10K	10K	10K	10K	10K	10K
Correlation	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $	$ r $
BEER	0.9903	<b>0.9891</b>	0.9775	0.9587	0.9864	0.9327	0.9251
BLEND	—	—	—	—	0.9861	—	—
BLEU	0.9931	0.9774	0.9706	0.9582	0.9767	0.7963	0.9414
CDER	0.9949	0.9842	0.9809	0.9605	0.9821	0.8322	0.9564
CHARACTER	0.9902	0.9862	0.9495	0.9627	0.9814	0.7752	<b>0.9784</b>
ITER	0.8649	0.9778	<b>0.9817</b>	<b>0.9664</b>	0.9650	0.8724	—
NIST	<b>0.9967</b>	0.9839	0.9797	0.9436	0.9877	0.8703	0.9442
PER	0.9865	0.9787	0.9545	0.9044	0.9862	0.8289	0.9500
TER	0.9948	0.9861	0.9770	0.9391	0.9845	0.8373	0.9591
WER	0.9944	0.9842	0.9772	0.9418	0.9829	0.8239	0.9537
YiSi-0	0.9713	0.9829	0.9648	0.9422	0.9879	<b>0.9530</b>	0.9513
YiSi-1	0.9851	0.9826	0.9761	0.9384	<b>0.9893</b>	0.9418	0.9572
YiSi-1_SRL	—	0.9881	—	—	—	—	0.9479
<b>newstest2018 Hybrids</b>							

Table 7: Absolute Pearson correlation of out-of-English system-level metrics with DA human assessment for 10K hybrid super-sampled systems in newstest2018; ensemble metrics are highlighted in gray.

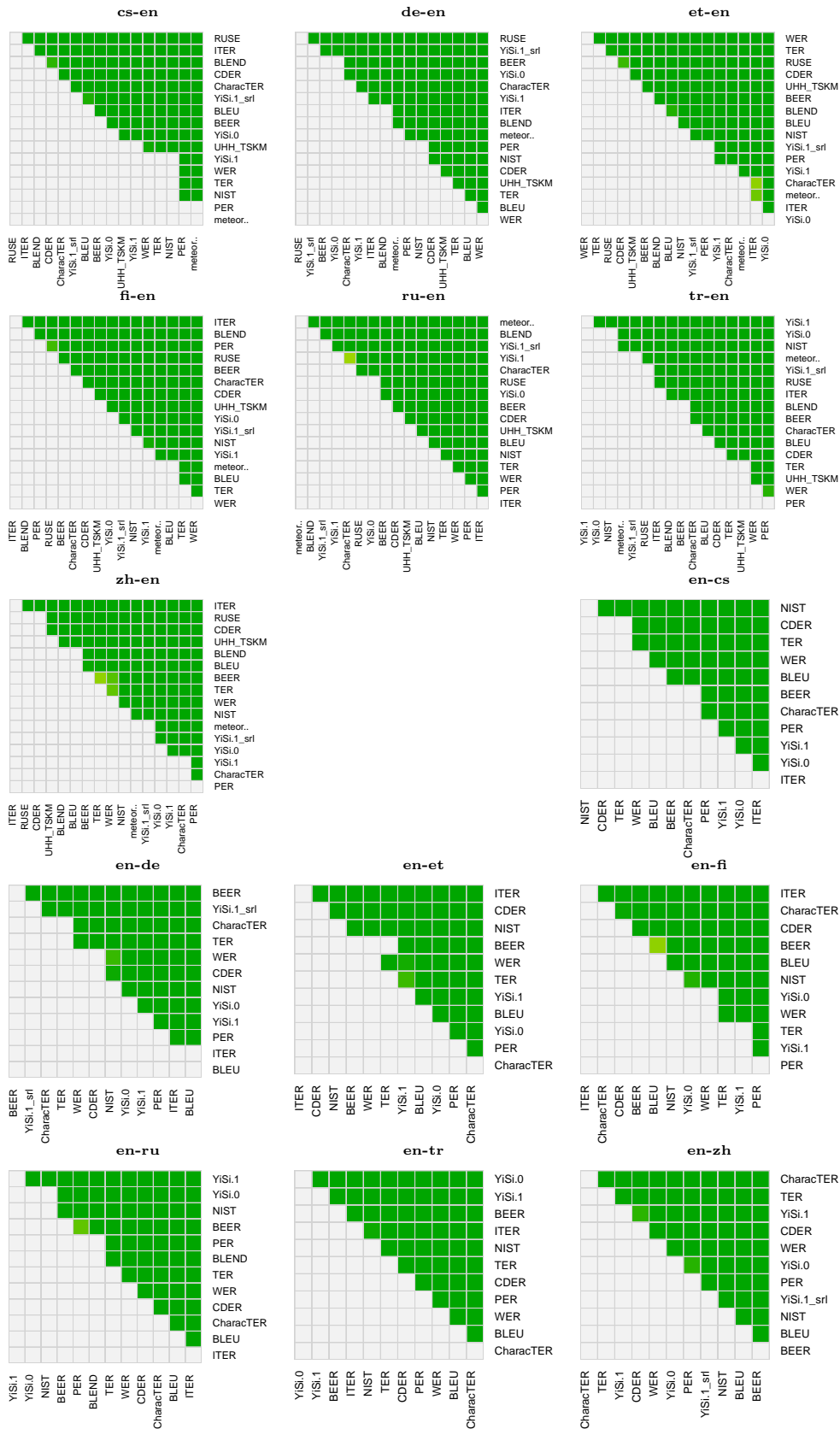


Figure 2: System-level metric significance test results for 10K hybrid systems (DA human evaluation) from newstest2018; green cells denote a statistically significant increase in correlation with human assessment for the metric in a given row over the metric in a given column according to Williams test.

significantly outperform the most widely employed metric BLEU, we include significance test results for every competing pair of metrics including our baseline metrics in Figure 1.

The sample of systems we employ to evaluate metrics is often small, as few as five MT systems for cs-en, for example. This can lead to inconclusive results, as identification of significant differences in correlations of metrics is unlikely at such a small sample size. Furthermore, Williams test takes into account the correlation between each pair of metrics, in addition to the correlation between the metric scores themselves, and this latter correlation increases the likelihood of a significant difference being identified.

To strengthen the conclusions of our evaluation, we include significance test results for large hybrid-super-samples of systems (Graham and Liu, 2016). 10K hybrid systems were created per language pair, with corresponding DA human assessment scores by sampling pairs of systems from WMT18 News Translation Task, creating hybrid systems by randomly selecting each candidate translation from one of the two selected systems. Similar to last year, not all metrics participating in the system-level evaluation submitted metric scores for the large set of hybrid systems. Fortunately, taking a simple average of segment-level scores is the proper aggregation method for almost all metrics this year, so where needed, we provided scores for hybrids ourselves, see Table 2.

Correlations of metric scores with human assessment of the large set of hybrid systems are shown in Tables 6 and 7, where again metrics not significantly outperformed by any other are highlighted in bold. Figure 2 then provides significance test results for hybrid super-sampled correlations for all pairs of competing metrics for a given language pair.

### 3.2 Segment-Level Results

Segment-level evaluation relies on the manual judgements collected in the News Translation Task evaluation. This year, we were unable to follow the methodology outlined in Graham et al. (2015) for evaluation of segment-level metrics because the sampling of sentences did not provide sufficient number of assessments of the same segment. We therefore convert

pairs of DA scores for competing translations to DARR better/worse preferences and employ a Kendall’s Tau formulation as described in Section 2.3.2.

Results of the segment-level human evaluation for translations sampled from the News Translation Task are shown in Tables 8 and 9, where metric correlations not significantly outperformed by any other metric are highlighted in bold. Head-to-head significance test results for differences in metric performance are included in Figure 3.

## 4 Discussion

### 4.1 Obtaining Human Judgements

Human data was collected in the usual way, a portion via crowd-sourcing and the remaining from researchers who mainly committed their time contribution to the manual evaluation as they had submitted a system in that language pair. Evaluation of translations employed the DA set-up and it again successfully acquired sufficient judgments to evaluate systems. As in the previous years, hybrid super-sampling proved very effective and allowed to obtain conclusive results of system-level evaluation even for language pairs where as few as 5 MT systems participated. We should however note that hybrid systems are constructed by randomly mixing sentences coming from different MT systems. As soon as document-level evaluation becomes relevant (which we anticipate in the next evaluation campaign already), this style of hybridization is susceptible to breaking cross-sentence references in MT outputs and may no longer be applicable.

In the case of segment-level evaluation, the optimal human evaluation data was unfortunately not available due to resource constraints. Conversion of document-level data held as a substitute for segment-level DA scores. These scores are however not optimal for evaluation of segment-level metrics and we would like to return to DA’s standard segment-level evaluation in future, where a minimum of 15 human judgments of translation quality are collected per translation and combined to get highly accurate scores for translations.

	cs-en	de-en	et-en	fi-en	ru-en	tr-en	zh-en
Human Evaluation	DARR	DARR	DARR	DARR	DARR	DARR	DARR
$n$	5,110	77,811	56,721	15,648	10,404	8,525	33,357
Correlation	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$
BEER	0.295	0.481	0.341	0.232	<b>0.288</b>	<b>0.229</b>	<b>0.214</b>
BLEND	<b>0.322</b>	<b>0.492</b>	0.354	0.226	<b>0.290</b>	<b>0.232</b>	<b>0.217</b>
CHARACTER	0.256	0.450	0.286	0.185	0.244	0.172	<b>0.202</b>
ITER	0.198	0.396	0.235	0.128	0.139	-0.029	0.144
METEOR++	0.270	0.457	0.329	0.207	0.253	0.204	0.179
RUSE	<b>0.347</b>	<b>0.498</b>	<b>0.368</b>	<b>0.273</b>	<b>0.311</b>	<b>0.259</b>	<b>0.218</b>
SENTBLEU	0.233	0.415	0.285	0.154	0.228	0.145	0.178
UHH_TSKM	0.274	0.436	0.300	0.168	0.235	0.154	0.151
YiSi-0	0.301	0.474	0.330	0.225	<b>0.294</b>	0.215	<b>0.205</b>
YiSi-1	<b>0.319</b>	0.488	0.351	0.231	<b>0.300</b>	<b>0.234</b>	<b>0.211</b>
YiSi-1_SRL	<b>0.317</b>	0.483	0.345	0.237	<b>0.306</b>	<b>0.233</b>	<b>0.209</b>

**newstest2018**

Table 8: Segment-level metric results for to-English language pairs in newstest2018: absolute Kendall’s Tau formulation of segment-level metric scores with DA scores; correlations of metrics not significantly outperformed by any other for that language pair are highlighted in bold; ensemble metrics are highlighted in gray.

	en-cs	en-de	en-et	en-fi	en-ru	en-tr	en-zh
Human Evaluation	DARR	DARR	DARR	DARR	DARR	DARR	DARR
$n$	5,413	19,711	32,202	9,809	22,181	1,358	28,602
Correlation	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$	$\tau$
BEER	<b>0.518</b>	<b>0.686</b>	<b>0.558</b>	<b>0.511</b>	<b>0.403</b>	<b>0.374</b>	0.302
BLEND	–	–	–	–	<b>0.394</b>	–	–
CHARACTER	0.414	0.604	0.464	0.403	0.352	<b>0.404</b>	<b>0.313</b>
ITER	0.333	0.610	0.392	0.311	0.291	0.236	–
SENTBLEU	0.389	0.620	0.414	0.355	0.330	0.261	<b>0.311</b>
YiSi-0	0.471	0.661	0.531	0.464	<b>0.394</b>	<b>0.376</b>	<b>0.318</b>
YiSi-1	<b>0.496</b>	<b>0.691</b>	<b>0.546</b>	<b>0.504</b>	<b>0.407</b>	<b>0.418</b>	<b>0.323</b>
YiSi-1_SRL	–	<b>0.696</b>	–	–	–	–	<b>0.310</b>

**newstest2018**

Table 9: Segment-level metric results for out-of-English language pairs in newstest2018: absolute Kendall’s Tau formulation of segment-level metric scores with DA scores; correlations of metrics not significantly outperformed by any other for that language pair are highlighted in bold; ensemble metrics are highlighted in gray.

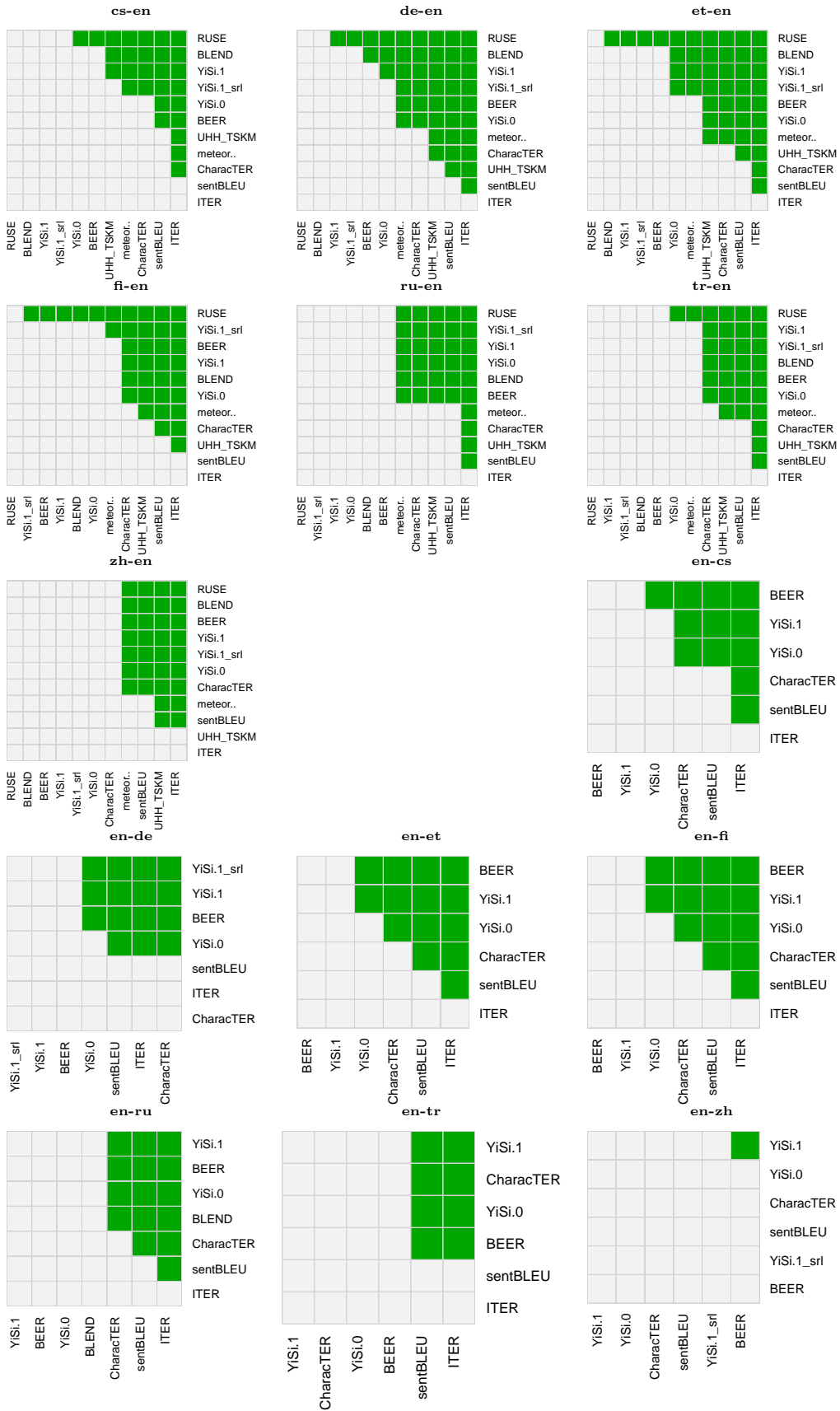


Figure 3: DARR segment-level metric significance test results for all language pairs (newstest2018): Green cells denote a significant win for the metric in a given row over the metric in a given column according to bootstrap resampling.

## 4.2 Overall Metric Performance

As always, the observed performance of metrics depends on the underlying texts and systems that participate in the News Translation Task. Two new metrics, RUSE and YISI stand out as metrics that achieve highest correlation in the system level evaluation in more than one language pair according to the hybrid evaluation, and perform great across all their language pairs on average. ITER also performs very well in en-et, en-fi, zh-en and several other languages but fails for en-ru and en-es, which drags its overall performance down.

Both YISI and RUSE are based on neural networks (YISI via word and phrase embeddings, RUSE via sentence embeddings). This is a new trend compared to the last year evaluation where the best performance was reached by character-level (not deep) metrics BEER, CHRFB (and its variants) and CHARACTER.

It is however important to note that the results of performance aggregated over language pairs are not particularly stable across years. In the last year’s evaluation, NIST seemed worse than TER. The overall results is the opposite this year and NIST even ranks slightly better than RUSE in terms of average system-level correlation across languages.

Overall, the reported figures confirm the observation from the past years that system-level metrics can achieve correlations above 0.9 but even the best ones can fall to 0.7 or 0.8 for some language pairs. Kendall’s Tau achieved by segment-level metrics are generally lower, in the range of 0.25–0.4. The best metrics in their best language pairs can reach up to 0.69 of segment-level correlations with humans. This capping could be possibly in part attributed to the sub-optimal human evaluation data, DA judgements converted to relative ranking.

Two metrics that stand out as performing consistently well are RUSE for evaluation of into-English translation and YISI-1\* for out-of-English. Overall, YISI\*, BEER, CHARACTER, RUSE, and BLEND (in this order) outperform SENTBLEU.

All of the “winners” in this years campaign are publicly available, which is very good for their prospective wider adoption. If participants could put the additional effort of adding

their code to Moses scorer, this would guarantee their long-term inclusion in the Metrics Task.

## 5 Conclusion

This paper summarizes the results of WMT18 shared task in machine translation evaluation, the Metrics Shared Task. Participating metrics were evaluated in terms of their correlation with human judgment at the level of the whole test set (system-level evaluation), as well as at the level of individual sentences (segment-level evaluation). For the former, best metrics reach over 0.95 Pearson correlation or better across several language pairs. Correlations varied more than usual between 0.2 and 0.7 in terms of segment-level metrics Kendall’s  $\tau$  results.

## Acknowledgments

Results in this shared task would not be possible without tight collaboration with organizers of the WMT News Translation Task.

This study was supported in parts by the grants 18-24210S of the Czech Science Foundation, ADAPT Centre for Digital Content Technology ([www.adaptcentre.ie](http://www.adaptcentre.ie)) at Dublin City University funded under the SFI Research Centres Programme (Grant 13/RC/2106) co-funded under the European Regional Development Fund, and Charles University Research Programme “Progress” Q18+Q48.

## References

- Ondřej Bojar, Miloš Ercegovčević, Martin Popel, and Omar Zaidan. 2011. 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. Association for Computational Linguistics.
- Ondřej Bojar, Christian Federmann, Barry Haddow, Philipp Koehn, Matt Post, and Lucia Specia. 2016. Ten Years of WMT Evaluation Campaigns: Lessons Learnt. In *Proceedings of the LREC 2016 Workshop “Translation Evaluation – From Fragmented Tools and Data Sets to an Integrated Ecosystem”*, pages 27–34, Portorož, Slovenia, 5.
- Ondřej Bojar, Yvette Graham, and Amir Kamran. 2017. Results of the WMT17 metrics shared



- task. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Tasks Papers*, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Ondřej Bojar, Jiří Mírovský, Kateřina Rysová, and Magdaléna Rysová. 2018. Evald reference-less discourse evaluation for wmt18. In *Proceedings of the Third Conference on Machine Translation*, Belgium, Brussels, October. Association for Computational Linguistics.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, Baltimore, Maryland, USA, June. Association for Computational Linguistics.
- George Doddington. 2002. Automatic Evaluation of Machine Translation Quality Using N-gram Co-occurrence Statistics. In *Proceedings of the Second International Conference on Human Language Technology Research, HLT '02*, pages 138–145, San Francisco, CA, USA. Morgan Kaufmann Publishers Inc.
- Melania Duma and Wolfgang Menzel. 2017. UHH submission to the WMT17 metrics shared task. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Tasks Papers*, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Yvette Graham and Timothy Baldwin. 2014. Testing for Significance of Increased Correlation with Human Judgment. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 172–176, Doha, Qatar, October. Association for Computational Linguistics.
- Yvette Graham and Qun Liu. 2016. Achieving Accurate Conclusions in Evaluation of Automatic Machine Translation Metrics. In *Proceedings of the 15th Annual Conference of the Association for Computational Linguistics: Human Language Technologies*, San Diego, CA. Association for Computational Linguistics.
- Yvette Graham, Timothy Baldwin, Alistair Mofat, and Justin Zobel. 2013. Continuous Measurement Scales in Human Evaluation of Machine Translation. In *Proceedings of the 7th Linguistic Annotation Workshop & Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.
- Yvette Graham, Timothy Baldwin, Alistair Mofat, and Justin Zobel. 2014a. Is Machine Translation Getting Better over Time? In *Proceedings of the 14th Conference of the European Chapter of the Association for Computational Linguistics*, pages 443–451, Gothenburg, Sweden, April. Association for Computational Linguistics.
- Yvette Graham, Nitika Mathur, and Timothy Baldwin. 2014b. Randomized significance tests in machine translation. In *Proceedings of the ACL 2014 Ninth Workshop on Statistical Machine Translation*, pages 266–274. Association for Computational Linguistics.
- Yvette Graham, Nitika Mathur, and Timothy Baldwin. 2015. Accurate Evaluation of Segment-level Machine Translation Metrics. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics Human Language Technologies*, Denver, Colorado.
- Yvette Graham, Timothy Baldwin, Alistair Mofat, and Justin Zobel. 2016. Can machine translation systems be evaluated by the crowd alone. *Natural Language Engineering*, FirstView:1–28, 1.
- Yinuo Guo, Chong Ruan, and Junfeng Hu. 2018. Meteor++: Incorporating copy knowledge into machine translation evaluation. In *Proceedings of the Third Conference on Machine Translation*, Belgium, Brussels, October. Association for Computational Linguistics.
- Philipp Koehn and Christof Monz. 2006. Manual and Automatic Evaluation of Machine Translation Between European Languages. In *Proceedings of the Workshop on Statistical Machine Translation, StatMT '06*, pages 102–121, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In *Proc. of Empirical Methods in Natural Language Processing*, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.
- Gregor Leusch, Nicola Ueffing, and Hermann Ney. 2006. CDER: Efficient MT Evaluation Using Block Movements. In *In Proceedings of EACL*, pages 241–248.
- Chi-kiu Lo. 2018. The NRC metric submission to the WMT18 metric and parallel corpus filtering shared task. In *Arxiv*.
- Qingsong Ma, Yvette Graham, Shugen Wang, and Qun Liu. 2017. Blend: a novel combined MT metric based on direct assessment — casict-dcu submission to WMT17 metrics task. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Tasks Papers*, Copenhagen, Denmark, September. Association for Computational Linguistics.

- Matouš Macháček and Ondřej Bojar. 2014. Results of the WMT14 metrics shared task. In *Proceedings of the Ninth Workshop on Statistical Machine Translation*, pages 293–301, Baltimore, MD, USA. Association for Computational Linguistics.
- Matouš Macháček and Ondřej Bojar. 2013. Results of the WMT13 Metrics Shared Task. In *Proceedings of the Eighth Workshop on Statistical Machine Translation*, pages 45–51, Sofia, Bulgaria, August. Association for Computational Linguistics.
- Joybrata Panja and Sudip Kumar Naskar. 2018. Iter: Improving translation edit rate through optimizable edit costs. In *Proceedings of the Third Conference on Machine Translation*, Belgium, Brussels, October. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: A Method for Automatic Evaluation of Machine Translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics*, ACL '02, pages 311–318.
- Maja Popović. 2015. chrF: character n-gram F-score for automatic MT evaluation. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, Lisboa, Portugal, September. Association for Computational Linguistics.
- Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation, Volume 2: Shared Tasks Papers*, Copenhagen, Denmark, September. Association for Computational Linguistics.
- Hiroki Shimanaka, Tomoyuki Kajiwara, and Mamoru Komachi. 2018. Ruse: Regressor using sentence embeddings for automatic machine translation evaluation. In *Proceedings of the Third Conference on Machine Translation*, Belgium, Brussels, October. Association for Computational Linguistics.
- Matthew Snover, Bonnie Dorr, Richard Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *In Proceedings of Association for Machine Translation in the Americas*, pages 223–231.
- Miloš Stanojević and Khalil Sima'an. 2015. BEER 1.1: ILLC UvA submission to metrics and tuning task. In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, Lisboa, Portugal, September. Association for Computational Linguistics.
- Yolanda Vazquez-Alvarez and Mark Huckvale. 2002. The reliability of the ITU-t p.85 standard for the evaluation of text-to-speech systems. In *Proc. of ICSLP - INTERSPEECH*.
- Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016a. Character: Translation edit rate on character level. In *ACL 2016 First Conference on Machine Translation*, pages 505–510, Berlin, Germany, August.
- Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016b. Character: Translation Edit Rate on Character Level. In *Proceedings of the First Conference on Machine Translation*, Berlin, Germany, August. Association for Computational Linguistics.
- Evan James Williams. 1959. *Regression analysis*, volume 14. Wiley New York.