

Quality-Aware Decoding for Neural Machine Translation

Anonymous ACL submission

Abstract

Despite the progress in machine translation quality estimation and evaluation in the last years, decoding in neural machine translation (NMT) is mostly oblivious to this and centers around finding the most probable translation according to the model (MAP decoding), approximated with beam search. In this paper, we bring together these two lines of research and propose *quality-aware decoding* for NMT, by leveraging recent breakthroughs in reference-free and reference-based MT evaluation through various inference methods like N -best reranking and minimum Bayes risk decoding. We perform an extensive comparison of various possible candidate generation and ranking methods across four datasets and two model classes and find that quality-aware decoding consistently outperforms MAP-based decoding according both to state-of-the-art automatic metrics (COMET and BLEURT) and to human assessments.

1 Introduction

The most common procedure in neural machine translation (NMT) is to train models using maximum likelihood estimation (MLE) at training time, and to decode with beam search at test time, as a way to approximate maximum-a-posteriori (MAP) decoding. However, several works have questioned the utility of model likelihood as a good proxy for translation quality (Koehn and Knowles, 2017; Ott et al., 2018; Stahlberg and Byrne, 2019; Eikema and Aziz, 2020). In parallel, significant progress has been made in methods for quality estimation and evaluation of generated translations (Specia et al., 2020; Mathur et al., 2020b), but this progress is (by large) not yet reflected in either training or decoding methods. Exceptions such as minimum risk training (Shen et al., 2016; Edunov et al., 2018) come at a cost of more expensive and unstable training, often with modest quality improvements.

An appealing alternative is to modify the decoding procedure only, separating it into two stages:

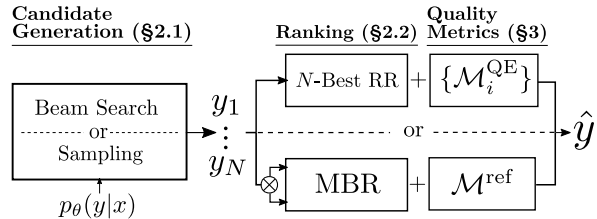


Figure 1: Quality-aware decoding framework. First, translation candidates are *generated* according to the model. Then, using reference-free and/or reference-based MT metrics, these candidates are *ranked*, and the highest ranked one is picked as the final translation.

candidate generation (§2.1; where candidates are generated with beam search or sampled from the whole distribution) and *ranking* (§2.2; where they are scored using a quality metric of interest, and the translation with the highest score is picked). This strategy has been explored in approaches using N -best reranking (Ng et al., 2019; Bhattacharyya et al., 2021) and minimum Bayes risk (MBR) decoding (Shu and Nakayama, 2017a; Eikema and Aziz, 2021; Müller and Sennrich, 2021). While this previous work has exhibited promising results, it has mostly focused on optimizing lexical metrics such as BLEU or METEOR (Papineni et al., 2002; Lavie and Denkowski, 2009), which have limited correlation with human judgments (Mathur et al., 2020a; Freitag et al., 2021a). Moreover, a rigorous apples-to-apples comparison among this suite of techniques and their variants is still missing, even though they share similar building blocks.

Our work fills these gaps by asking the question:

“Can we leverage recent advances in MT quality evaluation to generate better translations? If so, how can we most effectively do so?”

To answer this question, we systematically explore NMT decoding using a suite of ranking procedures. We take advantage of recent state-of-the-art learnable metrics, both reference-based, such as COMET and BLEURT (Rei et al., 2020a; Sel-

lam et al., 2020), and reference-free (also known as *quality estimation*; QE), such as TransQuest and OpenKiwi (Ranasinghe et al., 2020; Kepler et al., 2019). We compare different ranking strategies under a unified framework, which we name **quality-aware decoding** (§3). First, we analyze the performance of decoding using N -best reranking, both *fixed* according to a single metric and *learned* using multiple metrics, where the coefficients for each metric are optimized according to a reference-based metric. Second, we explore ranking using reference-based metrics directly through MBR decoding. Finally, to circumvent the expensive computational cost of the latter when the number of candidates is large, we develop a two-stage ranking procedure, where we use N -best reranking to pick a subset of the candidates to be ranked through MBR decoding. We explore the interaction of these different ranking methods with various candidate generation procedures including beam search, vanilla sampling, and nucleus sampling.

Experiments with two model sizes and four datasets (§4) reveal that while MAP-based decoding appears competitive when evaluating with lexical-based metrics (BLEU and ChrF), the story is very different with state-of-the-art evaluation metrics, where quality-aware decoding shows significant gains, both with N -best reranking and MBR decoding. We perform a human-study to more faithfully evaluate our systems and find that, while performance on learnable metrics is not always predictive of the best system, quality-aware decoding usually results in translations with higher quality than MAP-based decoding.

2 Candidate Generation and Ranking

We start by reviewing some of the most commonly used methods for both candidate generation and ranking under a common lens.

2.1 Candidate Generation

An NMT model defines a probability distribution $p_\theta(y|x)$ over a set of hypotheses \mathcal{Y} , conditioned on a source sentence x , where θ are learned parameters. A translation is typically predicted using MAP decoding, formalized as

$$\hat{y}_{\text{MAP}} = \arg \max_{y \in \mathcal{Y}} \log p_\theta(y|x). \quad (1)$$

In words, MAP decoding searches for the most probable translation under $p_\theta(y|x)$, *i.e.*, the mode

of the model distribution. Finding the exact \hat{y}_{MAP} is intractable since the search space \mathcal{Y} is combinatorially large, thus, approximations like **beam search** (Graves, 2012; Sutskever et al., 2014) are used. However, it has been shown that the translation quality *degrades* for large values of the beam size (Koehn and Knowles, 2017; Yang et al., 2018; Murray and Chiang, 2018; Meister et al., 2020), with the empty string often being the true MAP hypothesis (Stahlberg and Byrne, 2019).

A stochastic alternative to beam search is to *draw samples* directly from $p_\theta(y|x)$ with ancestral sampling, optionally with variants that truncate this distribution, such as top- k sampling (Fan et al., 2018) or p -**nucleus sampling** (Holtzman et al., 2020) – the latter samples from the smallest set of words whose cumulative probability is larger than a predefined value p . Deterministic methods combining beam and nucleus search have also been proposed (Shaham and Levy, 2021).

Unlike beam search, sampling is not a search algorithm nor a decision rule – it is not expected for a single sample to outperform MAP decoding (Eikema and Aziz, 2020). However, samples from the model can still be useful for alternative decoding methods, as we shall see. While beam search focus on high probability candidates, typically similar to each other, sampling allows for more *exploration*, leading to higher candidate *diversity*.

2.2 Ranking

We assume access to a set $\bar{\mathcal{Y}} \subseteq \mathcal{Y}$ containing N candidate translations for a source sentence, obtained with one of the generation procedures described in §2.1. As long as N is relatively small, it is possible to (re-)rank these candidates in a post-hoc manner, such that the best translation maximizes a given metric of interest. We highlight two different lines of work for ranking in MT decoding: first, **N -best reranking**, using reference-free metrics as features; second, **MBR decoding**, using reference-based metrics.

2.2.1 N -best Reranking

In its simplest form (which we call *fixed* reranking), a *single* feature f is used (*e.g.*, an estimated quality score), and the candidate that maximizes this score is picked as the final translation,

$$\hat{y}_{\text{F-RR}} = \arg \max_{y \in \bar{\mathcal{Y}}} f(y). \quad (2)$$

When *multiple* features $[f_1, \dots, f_K]$ are available, one can tune weights $[w_1, \dots, w_K]$ for these features to maximize a given reference-based evaluation metric on a validation set (Och, 2003; Duh and Kirchhoff, 2008) – we call this *tuned* reranking. In this case, the final translation is

$$\hat{y}_{\text{T-RR}} = \arg \max_{y \in \tilde{\mathcal{Y}}} \sum_{k=1}^K w_k f_k(y). \quad (3)$$

2.2.2 Minimum Bayes Risk (MBR) Decoding

While the techniques above rely on *reference-free* metrics for the computation of features, MBR decoding uses *reference-based* metrics to rank candidates. Unlike MAP decoding, which searches for the most probable translation, MBR decoding aims to find the translation that maximizes the expected *utility* (equivalently, that minimizes *risk*, Kumar and Byrne 2002, 2004; Eikema and Aziz 2020). Let again $\tilde{\mathcal{Y}} \subseteq \mathcal{Y}$ be a set containing N hypotheses and $u(y, y^*)$ a utility function measuring the similarity between a hypothesis $y \in \mathcal{Y}$ and a reference $y^* \in \tilde{\mathcal{Y}}$ (e.g. an automatic evaluation metric such as BLEU or COMET). MBR decoding seeks for

$$\hat{y}_{\text{MBR}} = \arg \max_{y^* \in \tilde{\mathcal{Y}}} \underbrace{\mathbb{E}_{Y \sim p_\theta(y|x)} [u(Y, y^*)]}_{\approx \frac{1}{M} \sum_{j=1}^M u(y^{(j)}, y^*)}, \quad (4)$$

where in Eq. 4 the expectation is approximated as a Monte Carlo (MC) sum using model samples $y^{(1)}, \dots, y^{(M)} \sim p_\theta(y|x)$.¹ In practice, the translation with the highest expected utility can be computed for each hypothesis $y^* \in \tilde{\mathcal{Y}}$ by comparing it to all the other hypotheses in the set.

3 Quality-Aware Decoding

While recent works have explored various combinations of candidate generation and ranking procedures for NMT (Lee et al., 2021; Bhattacharyya et al., 2021; Eikema and Aziz, 2021; Müller and Sennrich, 2021), they suffer from two limitations:

- The ranking procedure is usually based on simple lexical-based metrics (BLEU, chrF, METEOR). Although these metrics are well established and inexpensive to compute, they correlate poorly with human judgments at segment level (Mathur et al., 2020b; Freitag et al., 2021b).

¹We also consider the case where $y^{(1)}, \dots, y^{(M)}$ are obtained from nucleus sampling or beam search. Although the original MC estimate is unbiased, these ones are biased.

- Each work independently explores N -best reranking or MBR decoding, making unclear which method produces better translations.

In this work, we hypothesize that using more powerful metrics in the ranking procedure may lead to better quality translations. We propose a unified framework for ranking with both reference-based (§3.1) and reference-free metrics (§3.2), independently of the candidate generation procedure. We explore four methods with different computational costs for a given number of candidates, N .

Fixed N -best Reranker. An N -best reranker using a single reference-free metric (§3.2) as a feature, according to Eq. 2. The computational cost of this reranker is $\mathcal{O}(N \times C_{\mathcal{M}^{\text{QE}}})$, where $C_{\mathcal{M}^{\text{QE}}}$ denotes the cost of running an evaluation with a metric \mathcal{M}^{QE} .

Tuned N -best Reranker. An N -best reranker using as features *all* the reference-free metrics in §3.2, along with the model log-likelihood $\log p_\theta(y|x)$. The weights in Eq. 3 are optimized to maximize a given reference-based metric \mathcal{M}^{ref} using MERT (Och, 2003), a coordinate-ascent optimization algorithm widely used in previous work. The decoding cost is $\mathcal{O}(N \times \sum_i C_{\mathcal{M}_i^{\text{QE}}})$ for all metrics $\{\mathcal{M}_i^{\text{QE}}\}$.

MBR Decoding. Choosing as the utility function a reference-based metric \mathcal{M}^{ref} (§3.1), we estimate the utility using a simple Monte Carlo sum, as shown in Eq. 4. The estimation requires computing pairwise comparisons and thus the cost of running MBR decoding is $\mathcal{O}(N^2 \times C_{\mathcal{M}^{\text{ref}}})$.

N -best Reranker \rightarrow MBR. Using a large number of samples in MBR decoding is expensive due to its quadratic cost. To circumvent this issue, we explore a *two-stage* ranking approach: we first rank all the candidates using a tuned N -best reranker, followed by MBR decoding using the top M candidates. The computational cost becomes $\mathcal{O}(N \times \sum_i C_{\mathcal{M}_i} + M^2 \times C_{\mathcal{M}^{\text{ref}}})$. The first ranking stage *prunes* the candidate list to a smaller, higher quality subset, making possible a more accurate estimation of the utility with less samples, and potentially allowing a better reranker than *plain* MBR for almost the same computational budget.

3.1 Reference-based Metrics

Reference-based metrics are the standard way to evaluate MT systems; the most used ones rely on the lexical overlap between hypotheses and reference translations (Papineni et al., 2002; Lavie

and Denkowski, 2009; Popović, 2015). However, lexical-based approaches have important limitations: they have difficulties recognizing correct translations that are paraphrases of the reference(s); they ignore the source sentence, an important indicator of meaning for the translation; and they do not always correlate well with human judgments, particularly at segment-level (Freitag et al., 2021b).

In this work, apart from BLEU and chrF, we use the following state-of-the-art trainable reference-based metrics for both ranking and performance evaluation of MT systems:

- BLEURT (Sellam et al., 2020; Pu et al., 2021b), trained to regress on human direct assessments (DA; Graham et al. 2013). We use the largest multilingual version, *BLEURT-20*, based on the RemBERT model (Chung et al., 2021).
- COMET (Rei et al., 2020a), based on XLM-R (Conneau et al., 2020), trained to regress on quality assessments such as DA using both the reference and the source to assess the quality of a given translation. We use the publicly available model developed for the WMT20 metrics shared task (*wmt20-comet-da*).

These metrics have shown much better correlation at segment-level than previous lexical metrics in WMT metrics shared tasks (Mathur et al., 2020b; Freitag et al., 2021b). Hence, as discussed in §2.2, they are good candidates to be used either *indirectly* as an optimization objective for learning the tuned reranker’s feature weights, or *directly* as a utility function in MBR decoding. In the former, the higher the metric correlation with human judgment, the better the translation picked by the tuned reranker. In the latter, we approximate the expected utility in Eq. 4 by letting a candidate generated by the model be a reference translation – a suitable premise *if* the model is good in expectation.

3.2 Reference-free Metrics

MT evaluation metrics have also been developed for the case where references are not available – they are called *reference-free* or *quality estimation* (QE) metrics. In the last years, considerable improvements have been made to such metrics, with state-of-the-art models having increasing correlations with human annotators (Freitag et al., 2021b; Specia et al., 2021). These improvements enable the use of such models for ranking translation hypotheses in a more reliable way than before.

In this work, we explore four recently proposed reference-free metrics as features for N -best reranking, all at the sentence-level:

- COMET-QE (Rei et al., 2020b), a reference-free version of COMET (§3.1). It was the winning submission for the QE-as-a-metric subtask of the WMT20 shared task (Mathur et al., 2020b).
- TransQuest (Ranasinghe et al., 2020), the winning submission for the sentence-level DA prediction subtask of the WMT20 QE shared task (Specia et al., 2020). Similarly to COMET-QE this metric predicts a DA score.
- MBART-QE (Zerva et al., 2021), based on the mBART (Liu et al., 2020) model, trained to predict both the *mean* and the *variance* of DA scores. It was a top performer in the WMT21 QE shared task (Specia et al., 2021).
- OpenKiwi-MQM (Kepler et al., 2019; Rei et al., 2021), based on XLM-R, trained to predict the *multidimensional quality metric* (MQM; Lommel et al. 2014).² This reference-free metric was ranked second on the QE-as-a-metric subtask from the WMT 2021 metrics shared task.

4 Experiments

4.1 Setup

We study the benefits of quality-aware decoding over MAP-based decoding in two regimes:

- A high-resource, unconstrained, setting with *large* transformer models (6 layers, 16 attention heads, 1024 embedding dimensions, and 8192 hidden dimensions) trained by Ng et al. (2019) for the WMT19 news translation task (Barrault et al., 2019), using English to German (EN → DE) and English to Russian (EN → RU) language pairs. These models were trained on over 20 million parallel and 100 million back-translated sentences, being the winning submissions of that year’s shared task. We consider the non-ensembled version of the model and use *newstest19* for validation and *newstest20* for testing.
- A more constrained scenario with a *small* transformer model (6 layers, 4 attention heads, 512 embedding dimensions, and 1024 hidden dimensions) trained from scratch in *Fairseq* (Ott et al.,

²MQM annotations are expert-level type of annotations more fine-grained than DA, with individual errors annotated.

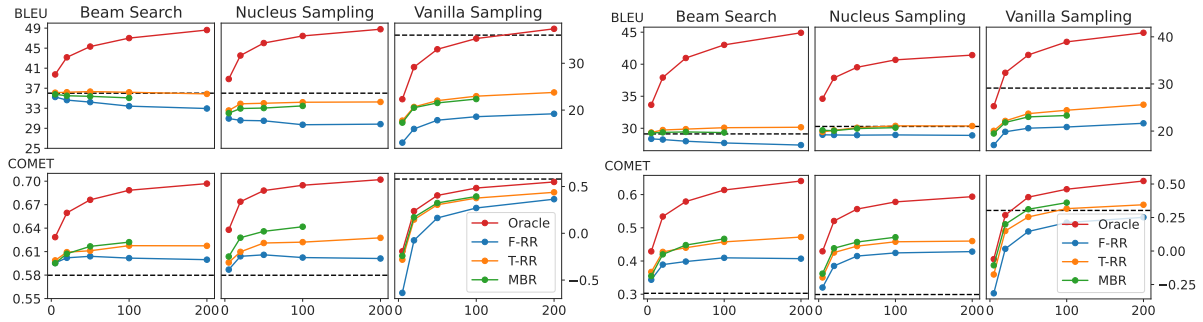


Figure 2: Values for BLEU (top) and COMET (bottom) for EN \rightarrow DE as we increase the number of candidates for different generation and ranking procedures, as well as oracles with the respective metrics, for the *large* (left) and *small* (right) models. Baseline values are marked with a dashed horizontal line.

2019) on the smaller IWSLT17 datasets (Cettolo et al., 2012) for English to German (EN \rightarrow DE) and English to French (EN \rightarrow FR), each with a little over 200k training examples. We chose these datasets because they have been extensively used in previous work (Bhattacharyya et al., 2021) and smaller model allows us to answer questions about how the training methodology affects ranking performance (see § 4.2.2). Further training details can be found in Appendix A.

We use beam search with a beam size of 5 as our decoding baseline. For tuned N -best reranking, we use Travatar’s (Neubig, 2013) implementation of MERT (Och, 2003) to optimize the weight of each feature, as described in §3.2. Finally, we evaluate each system using the metrics discussed in §3.1, along with BLEU and chrF (Popović, 2015).

4.2 Results

Overall, given all the metrics, candidate generation, and ranking procedures, we evaluate over 150 systems per dataset. We report subsets of this data separately to answer specific research questions, and defer to Appendix B for additional results.

4.2.1 Impact of Candidate Generation

First, we explore the impact of the candidate generation procedure and the number of candidates.

Which candidate generation method works best, beam search or sampling? We generate candidates with beam search, vanilla sampling, and nucleus sampling. For the latter, we use $p = 0.6$ based on early results showing improved performance for all metrics.³ For N -best reranking, we

³We picked nucleus sampling over top- k sampling because it allows varying support size and has outperformed top- k in text generation tasks (Holtzman et al., 2020).

use up to 200 samples; for MBR decoding, due to the quadratic computational cost, we use up to 100.

Figure 2 shows BLEU and COMET for different candidate generation and ranking methods for the EN \rightarrow DE WMT20 and IWSLT17 datasets, with increasing number of candidates. To assess the performance *ceiling* of the rankers, we also report results with an *oracle* ranker for the reported metrics, picking the candidate that maximizes it. For the *fixed* N -best reranker, we use COMET-QE as a metric, albeit the results for other reference-free metrics are similar. Performance seems to scale well with the number of candidates, particularly for vanilla sampling and for the *tuned* N -best reranker and MBR decoder – this is in line with the findings of previous work (Lee et al., 2021; Müller and Senrich, 2021). However, all the rankers using vanilla sampling severely under-perform the baseline in most cases (we will come back to this in §4.2.2). In contrast, the rankers using beam search or nucleus sampling are competitive or outperform the baseline in terms of BLEU, and greatly outperform it in terms of COMET. For the larger models trained on WMT20, we see that the performance according to the lexical metrics degrades with more candidates. In this scenario, rankers using candidates generated by nucleus sampling seem to have an edge over the ones that use beam search for COMET.

Based on the findings above, and due to generally better performance of COMET over BLEU for MT evaluation (Kocmi et al., 2021), in following experiments we use nucleus sampling with the *large* model and beam search with the *small* model.

4.2.2 Impact of Label Smoothing

How does label smoothing affect candidate generation? Label smoothing (Szegedy et al., 2016) is a regularization technique that redistributes proba-

	Large (WMT20)				Small (IWSLT)			
	BLEU	chrF	BLEURT	COMET	BLEU	chrF	BLEURT	COMET
Baseline	36.01	63.88	0.7376	0.5795	29.12	56.23	0.6635	0.3028
F-RR w/ COMET-QE	29.83	59.91	<u>0.7457</u>	<u>0.6012</u>	<u>27.38</u>	54.89	<u>0.6848</u>	<u>0.4071</u>
F-RR w/ MBART-QE	<u>32.92</u>	<u>62.71</u>	0.7384	0.5831	27.30	<u>55.62</u>	0.6765	0.3533
F-RR w/ OpenKiwi	30.38	<u>59.56</u>	0.7401	0.5623	25.35	51.53	0.6524	0.2200
F-RR w/ Transquest	31.28	60.94	0.7368	0.5739	26.90	54.46	0.6613	0.2999
T-RR w/ BLEU	<u>35.34</u>	<u>63.82</u>	0.7407	0.5891	30.51	57.73	0.7077	0.4536
T-RR w/ BLEURT	33.39	62.56	<u>0.7552</u>	0.6217	30.16	57.40	<u>0.7127</u>	<u>0.4741</u>
T-RR w/ COMET	34.26	63.31	0.7546	<u>0.6276</u>	30.16	57.32	0.7124	0.4721
MBR w/ BLEU	<u>34.94</u>	<u>63.21</u>	0.7333	0.5680	29.25	56.36	0.6619	0.3017
MBR w/ BLEURT	32.90	62.34	<u>0.7649</u>	0.6047	28.69	56.28	<u>0.7051</u>	0.3799
MBR w/ COMET	33.04	62.65	0.7477	<u>0.6359</u>	<u>29.43</u>	<u>56.74</u>	0.6882	<u>0.4480</u>
T-RR+MBR w/ BLEU	<u>35.84</u>	63.96	0.7395	0.5888	<u>30.23</u>	<u>57.34</u>	0.6913	0.3969
T-RR+MBR w/ BLEURT	33.61	62.95	0.7658	0.6165	29.28	56.77	0.7225	0.4361
T-RR+MBR w/ COMET	34.20	63.35	0.7526	0.6418	29.46	57.13	0.7058	0.5005

Table 1: Evaluation metrics for EN \rightarrow DE for the *large* and *small* model settings, using a *fixed* N -best reranker (F-RR), a *tuned* N -best reranker (T-RR), MBR decoding, and a two-stage approach.

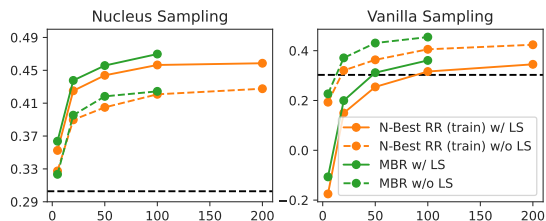


Figure 3: COMET scores for EN \rightarrow DE (IWSLT17) for models trained with and without label smoothing.

bility mass from the gold label to the other target labels, typically preventing the model from becoming overconfident (Müller et al., 2019). However, it has been found that label smoothing negatively impacts model fit, compromising the performance of MBR decoding (Eikema and Aziz, 2020, 2021). Thus, we train a small transformer model without label smoothing to verify its impact in the performance of N -best reranking and MBR decoding. Figure 3 shows that disabling label smoothing really helps when generating candidates using vanilla sampling. However, the performance *degrades* for candidates generated using nucleus sampling when we disable label smoothing, hinting that the pruning mechanism of nucleus sampling may help mitigate the negative impact of label smoothing in sampling based approaches. Even without label smoothing, vanilla sampling is not competitive with nucleus sampling or beam search with label smoothing, thus, we do not experiment further with it.

4.2.3 Impact of Ranking and Metrics

We now investigate the usefulness of the metrics presented in §3 as features and objectives for rank-

ing. For N -best reranking, we use all the available candidates (200) while, for MBR, due to the computational cost of using 100 candidates, we report results with 50 candidates only (we found that ranking with *tuned* N -best reranking with $N = 100$ and MBR with $N = 50$ takes about the same time). We report results in Table 1, and use them to answer some specific research questions.

Which QE metric works best in a fixed N -best reranker? We consider a *fixed* N -best reranker with a single reference-free metric as a feature (see Table 1, second group). While none of the metrics allows for improving the baseline results in terms of the lexical metrics (BLEU and chrF), rerankers using COMET-QE or MBART-QE outperform the baseline according to BLEURT and COMET, for both the *large* and *small* models. Due to the aforementioned better performance of these metrics for translation quality evaluation, we hypothesize that these rankers produce better translations than the baseline. However, since the sharp drop in the lexical metrics is concerning, we will verify this hypothesis in a human study, in §4.2.4.

How does the performance of a tuned N -best reranker vary when we change the optimization objective? We consider a *tuned* N -best reranker using as features *all* the reference-free metrics introduced in §3.2, and optimized using MERT. Table 1 (third group) shows the results for EN \rightarrow DE. For the *small* model, all the rankers show improved results over the baseline for all the metrics. In particular, optimizing for BLEU leads to the best results

469 in terms of the lexical metrics, while optimizing for
470 BLEURT leads to the best performance in terms of
471 the others. Finally, optimizing for COMET leads to
472 similar performance than optimizing for BLEURT.
473 For the *large* model, although none of the rerankers
474 is able to outperform the baseline in terms of the
475 lexical metrics, one can see similar trends as before
476 in terms of BLEURT and COMET.

477 **How does the performance of MBR decoding vary**
478 **when we change the utility function?** Table 1
479 (fourth group) shows the impact of the utility func-
480 tion (BLEU, BLEURT, or COMET). For the *small*
481 model, using COMET leads to the best perfor-
482 mance according to all the metrics except BLEURT
483 (for which the best result is attained when optimiz-
484 ing itself). For the *large* model, the best result
485 according to a given metric is obtained when using
486 that metric as the utility function.

487 **How do (tuned) *N*-best reranking and MBR com-**
488 **pare to each other?** Looking at the third and
489 fourth groups in Table 1 we see that, for the *small*
490 model, *N*-best reranking seems to perform better
491 than MBR decoding in terms of all the evaluation
492 metrics, including the one that was used as the
493 utility function in MBR decoding. The picture is
494 less clear for the *large* model, with MBR decoding
495 achieving best values for a given fine-tuned metric
496 when using it as the utility; this comes at the cost of
497 worse performance according to the others metrics,
498 hinting at a potential “overfitting” effect. Overall,
499 *N*-best reranking seems to have an edge over MBR
500 decoding. We will further clarify this question with
501 human evaluation in § 4.2.4.

502 **Can we improve performance by combining *N*-**
503 **best reranking with MBR decoding?** The results
504 in Table 1 show that, for both the *large* and the
505 *small* model, the two-stage ranking approach de-
506 scribed in §3 leads to the best performance accord-
507 ing to the fine-tuned metrics. In particular, the best
508 result is obtained when the utility function is the
509 same as the evaluation metric. These results sug-
510 gest that a promising research direction is to seek
511 more sophisticated pruning strategies as a preprocess-
512 ing step for MBR decoding.

513 4.2.4 Human Evaluation

514 **Which metric correlates more with human judg-**
515 **ments? How risky is it to optimize a metric and**
516 **evaluate on a related metric?** Our experiments
517 suggest that, overall, *quality-aware* decoding pro-

518 duces translations with better performance across
519 most metrics than *MAP-based* decoding. However,
520 for some cases (such as fixed *N*-best reranking and
521 most results with the *large* model), there is a con-
522 cerning “metric gap” between lexical-based and
523 fine-tuned metrics. While the latter have shown to
524 correlate better with human judgments, previous
525 work has not attempted to explicitly optimize these
526 metrics, and doing so could lead to ranking systems
527 that learn to exploit “pathologies” in these metrics
528 rather than improving translation quality.

529 To investigate this hypothesis, we perform a hu-
530 man study across all four datasets. We ask anno-
531 tators to rate, from 1 (no overlap in meaning) to
532 5 (perfect translation), the translations produced
533 by the 4 *ranking* systems mentioned §3, as well as
534 the baseline translation and the reference. Further
535 details can be found in Appendix C. We choose
536 COMET-QE as the feature for the fixed *N*-best
537 ranker and COMET as the optimization metric and
538 utility function for the tuned *N*-best reranker and
539 MBR decoding, respectively. The reasons for this
540 are two-fold: (1) they are currently the reference-
541 free and reference-based metrics with highest re-
542 ported correlation with human judgments (Kocmi
543 et al., 2021), (2) we saw the largest “metric gap” for
544 systems based on these metrics, hinting of a poten-
545 tial “overfitting” problem (specially since COMET-
546 QE and COMET are similar models).

547 Table 2 shows the results for the human evalu-
548 ation, as well as the automatic metrics. Overall,
549 we see that when fine-tuned metrics are explicitly
550 optimized for, their correlation with human judg-
551 ments decreases and they are no longer reliable in-
552 dicators of system-level ranking. This is especially
553 notable for the fixed *N*-best reranker with COMET-
554 QE, which outperforms the baseline in terms of
555 COMET in every single scenario, but results in
556 markedly lower quality translations. However, de-
557 spite the potential for overfitting these metrics, we
558 find that *tuned N*-best reranking, MBR, and their
559 combination consistently achieve better translation
560 quality than the baseline, especially with the small
561 model. In particular, *N*-best reranking seems to re-
562 sult in better translations than MBR, however their
563 combination is the best system in two of four LPs.

564 5 Related Work

565 **Reranking.** Inspired by the work of Shen et al.
566 (2004) on discriminative reranking for SMT, Lee
567 et al. (2021) trained a large transformer model us-

	EN-DE (WMT20)					EN-RU (WMT20)				
	BLEU	chrF	BLEURT	COMET	Human R.	BLEU	chrF	BLEURT	COMET	Human R.
Reference	-	-	-	-	4.51	-	-	-	-	4.07
Baseline	36.01	63.88	0.7376	0.5795	4.28	23.86	51.16	0.6953	0.5361	3.62
F-RR w/ COMET-QE	29.83	59.91	0.7457	0.6012	4.19	20.32	49.18	0.7130	0.6207	3.25
T-RR w/ COMET	34.26	63.31	0.7546	0.6276	4.33	22.42	50.91	0.7243	0.6441	3.65
MBR w/ COMET	33.04	62.65	0.7477	0.6359	4.27	23.67	51.18	0.7093	0.6242	3.66
F-RR + MBR w/ COMET	34.20	63.35	0.7526	0.6418	4.30	23.21	51.26	0.7238	0.6736	3.72 [†]

	EN-DE (IWSLT17)					EN-FR (IWSLT17)				
	BLEU	chrF	BLEURT	COMET	Human R.	BLEU	chrF	BLEURT	COMET	Human R.
Reference	-	-	-	-	4.38	-	-	-	-	4.00
Baseline	29.12	0.6635	56.23	0.3028	3.68	38.12	0.6532	63.20	0.4809	3.92
F-RR w/ COMET-QE	27.38	0.6848	54.89	0.4071	3.67	35.59	0.6628	60.90	0.5553	3.63
T-RR w/ COMET	30.16	0.7124	57.32	0.4721	3.90 [†]	38.60	0.7020	63.77	0.6392	4.05 [†]
MBR w/ COMET	29.43	0.6882	56.74	0.4480	3.79 [†]	37.77	0.6710	63.24	0.6127	4.05 [†]
F-RR + MBR w/ COMET	29.46	0.7058	57.13	0.5005	3.83 [†]	38.33	0.6883	63.53	0.6610	4.09 [†]

Table 2: Results for automatic and human evaluation. Top: WMT20 (large models); Bottom: IWSLT17 (small models). Methods with [†] are statistically significantly better than the baseline, with $p < 0.05$.

ing a reranking objective to optimize BLEU. Our work differs in which our rerankers are much simpler (a single feature or a linear combination of features) and therefore can simply be tuned on a validation set; and we use more powerful quality metrics instead of BLEU. Similarly, [Bhattacharyya et al. \(2021\)](#) learned an energy-based reranker to assign lower energy to the samples with higher BLEU scores. While the energy model plays a similar role to a QE system (the higher the quality, the lower the energy), our work differs in two ways: we use an existing, pretrained QE model instead of training a dedicated reranker, making our approach applicable to any MT system without requiring further training; and the QE model is trained to predict human assessments, rather than BLEU scores. [Leblond et al. \(2021\)](#) compare a reinforcement learning approach to some reranking approaches (but not MBR decoding, as we do). They investigate the use of reference-based metrics and, for the reward function, a reference-free metric based on a modified BERTScore ([Zhang et al., 2020](#)). This new multilingual BERTScore is not fine-tuned on human judgments as COMET and BLEURT and it is unclear what its level of agreement with human judgments is. Another line of work is *generative reranking*, where the reranker is not trained to optimize an evaluation metric directly, but rather as a generative noisy-channel model ([Yu et al., 2017](#); [Yee et al., 2019](#); [Ng et al., 2019](#)).

Minimum Bayes Risk Decoding. MBR decoding ([Kumar and Byrne, 2002, 2004](#)) has recently been revived for NMT using candidates generated with beam search ([Stahlberg et al., 2017](#); [Shu](#)

and [Nakayama, 2017b](#)) and sampling ([Eikema and Aziz, 2020, 2021](#); [Müller and Sennrich, 2021](#)). However, a comparison with N -best re-ranking was missing in these works, a gap our paper fills. A related line of work is *minimum risk training* (MRT; [Smith and Eisner 2006](#); [Shen et al. 2016](#)), which *trains* models to minimize risk, allowing arbitrary non-differentiable loss functions ([Edunov et al., 2018](#); [Wieting et al., 2019](#)) and avoiding exposure bias ([Wang and Sennrich, 2020](#); [Kiegl and Kreutzer, 2021](#)). However, MRT is considerably more expensive and difficult to train and the gains are often small. Incorporating our quality metrics in MRT is an exciting research direction.

6 Conclusions and Future Work

We leverage recent advances in MT quality estimation and evaluation and propose *quality-aware decoding* for NMT. We explore different candidate generation and ranking methods, with a comprehensive empirical analysis across four datasets and two model classes. We show that, compared to MAP-based decoding, quality-aware decoding leads to better translations, according to powerful automatic evaluation metrics and human judgments.

There are several directions for future work. Our ranking strategies, while leading to higher accuracies, are substantially more expensive, particularly when used with costly evaluation metrics such as BLEURT and COMET. While reranking-based pruning before MBR decoding was found helpful, additional strategies such as caching encoder representations and distillation of BERT-based metrics ([Pu et al., 2021a](#)) are promising directions.

References

- Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. 2019. [Findings of the 2019 conference on machine translation \(WMT19\)](#). In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 1–61, Florence, Italy. Association for Computational Linguistics.
- Sumanta Bhattacharyya, Amirmohammad Rooshenas, Subhajit Naskar, Simeng Sun, Mohit Iyyer, and Andrew McCallum. 2021. [Energy-based reranking: Improving neural machine translation using energy-based models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 4528–4537, Online. Association for Computational Linguistics.
- Mauro Cettolo, Christian Girardi, and Marcello Federico. 2012. [WIT3: Web inventory of transcribed and translated talks](#). In *Proceedings of the 16th Annual conference of the European Association for Machine Translation*, pages 261–268, Trento, Italy. European Association for Machine Translation.
- Hyung Won Chung, Thibault Févry, Henry Tsai, Melvin Johnson, and Sebastian Ruder. 2021. [Rethinking embedding coupling in pre-trained language models](#). *ArXiv*, abs/2010.12821.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Kevin Duh and Katrin Kirchhoff. 2008. [Beyond log-linear models: Boosted minimum error rate training for n-best re-ranking](#). In *Proceedings of ACL-08: HLT, Short Papers*, pages 37–40, Columbus, Ohio. Association for Computational Linguistics.
- Sergey Edunov, Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. [Classical structured prediction losses for sequence to sequence learning](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 355–364, New Orleans, Louisiana. Association for Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2020. [Is MAP decoding all you need? the inadequacy of the mode in neural machine translation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2021. [Sampling-based minimum bayes risk decoding for neural machine translation](#).
- Angela Fan, Mike Lewis, and Yann Dauphin. 2018. [Hierarchical neural story generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 889–898, Melbourne, Australia. Association for Computational Linguistics.
- Markus Freitag, George F. Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. [Experts, errors, and context: A large-scale study of human evaluation for machine translation](#). *CoRR*, abs/2104.14478.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondrej Bojar. 2021b. [Results of the WMT21 Metrics Shared Task: Evaluating Metrics with Expert-based Human Evaluations on TED and News Domain](#). In *Proceedings of the Sixth Conference on Machine Translation (WMT)*, pages 716–757. NRC.
- Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. 2013. [Continuous measurement scales in human evaluation of machine translation](#). In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 33–41, Sofia, Bulgaria. Association for Computational Linguistics.
- Alex Graves. 2012. [Sequence transduction with recurrent neural networks](#).
- Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. [The curious case of neural text de-generation](#). In *International Conference on Learning Representations*.
- Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. 2019. [OpenKiwi: An open source framework for quality estimation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 117–122, Florence, Italy. Association for Computational Linguistics.
- Samuel Kiegl and Julia Kreutzer. 2021. [Revisiting the weaknesses of reinforcement learning for neural machine translation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1673–1681, Online. Association for Computational Linguistics.
- Diederik P. Kingma and Jimmy Ba. 2015. [Adam: A method for stochastic optimization](#). In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

748	Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation.	805
749		806
750		807
751		808
752		809
753	Philipp Koehn and Rebecca Knowles. 2017. Six challenges for neural machine translation. In <i>Proceedings of the First Workshop on Neural Machine Translation</i> , pages 28–39, Vancouver. Association for Computational Linguistics.	810
754		811
755		812
756		813
757		814
758	Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 66–71, Brussels, Belgium. Association for Computational Linguistics.	815
759		816
760		
761		817
762		818
763		819
764		820
765	Shankar Kumar and William Byrne. 2002. Minimum bayes-risk word alignments of bilingual texts. In <i>Proceedings of the ACL-02 Conference on Empirical Methods in Natural Language Processing - Volume 10, EMNLP '02</i> , page 140–147, USA. Association for Computational Linguistics.	822
766		823
767		824
768		825
769		826
770		827
771	Shankar Kumar and William Byrne. 2004. Minimum Bayes-risk decoding for statistical machine translation. In <i>Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004</i> , pages 169–176, Boston, Massachusetts, USA. Association for Computational Linguistics.	828
772		829
773		830
774		831
775		
776		832
777		833
778		834
779		835
780		836
781	Alon Lavie and Michael J. Denkowski. 2009. The meteor metric for automatic evaluation of machine translation. <i>Machine Translation</i> , 23(2-3):105–115.	837
782		838
783		839
784	Rémi Leblond, Jean-Baptiste Alayrac, Laurent Sifre, Miruna Pislari, Lespiau Jean-Baptiste, Ioannis Antonoglou, Karen Simonyan, and Oriol Vinyals. 2021. Machine translation decoding beyond beam search. In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 8410–8434, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.	840
785		841
786		842
787		843
788		
789		844
790		845
791	Ann Lee, Michael Auli, and Marc’Aurelio Ranzato. 2021. Discriminative reranking for neural machine translation. In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 7250–7264, Online. Association for Computational Linguistics.	846
792		847
793		848
794		849
795		850
796		
797		851
798		852
799		853
800		854
801		855
802	Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.	856
803		857
804		858
		859
		860
		861
	Arle Lommel, Aljoscha Burchardt, and Hans Uszkoreit. 2014. Multidimensional quality metrics (MQM): A framework for declaring and describing translation quality metrics. <i>Tradumàtica: tecnologies de la traducció</i> , 0:455–463.	
	Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020a. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 4984–4997, Online. Association for Computational Linguistics.	
	Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020b. Results of the WMT20 metrics shared task. In <i>Proceedings of the Fifth Conference on Machine Translation</i> , pages 688–725, Online. Association for Computational Linguistics.	
	Clara Meister, Ryan Cotterell, and Tim Vieira. 2020. If beam search is the answer, what was the question? In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2173–2185, Online. Association for Computational Linguistics.	
	Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. 2019. When does label smoothing help? In <i>Advances in Neural Information Processing Systems</i> , volume 32. Curran Associates, Inc.	
	Kenton Murray and David Chiang. 2018. Correcting length bias in neural machine translation. In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 212–223, Brussels, Belgium. Association for Computational Linguistics.	
	Mathias Müller and Rico Sennrich. 2021. Understanding the properties of minimum bayes risk decoding in neural machine translation.	
	Graham Neubig. 2013. Travatar: A forest-to-string machine translation engine based on tree transducers. In <i>Proceedings of the ACL Demonstration Track</i> , Sofia, Bulgaria.	
	Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR’s WMT19 news translation task submission. In <i>Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)</i> , pages 314–319, Florence, Italy. Association for Computational Linguistics.	
	Franz Josef Och. 2003. Minimum error rate training in statistical machine translation. In <i>Proceedings of the 41st Annual Meeting of the Association for Computational Linguistics</i> , pages 160–167, Sapporo, Japan. Association for Computational Linguistics.	
	Myle Ott, Michael Auli, David Grangier, and Marc’Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In <i>Proceedings of the 35th International Conference on Machine Learning</i> , volume 80 of <i>Proceedings of Machine Learning Research</i> , pages 3956–3965. PMLR.	

862	Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan,	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020.	917
863	Sam Gross, Nathan Ng, David Grangier, and Michael	BLEURT: Learning robust metrics for text genera-	918
864	Auli. 2019. fairseq: A fast, extensible toolkit for	tion . In <i>Proceedings of the 58th Annual Meeting of</i>	919
865	sequence modeling . In <i>Proceedings of the 2019 Con-</i>	<i>ference of the North American Chapter of the Associa-</i>	920
866	<i>tion for Computational Linguistics (Demonstrations)</i> ,	pages 7881–7892, Online. Association for Computational	921
867	pages 48–53, Minneapolis, Minnesota. Association	Linguistics.	922
868	for Computational Linguistics.		
869			
870	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Uri Shaham and Omer Levy. 2021. What do you get	923
871	Jing Zhu. 2002. Bleu: a method for automatic evalua-	when you cross beam search with nucleus sampling?	924
872	tion of machine translation . In <i>Proceedings of the</i>		
873	<i>40th Annual Meeting of the Association for Computa-</i>	Libin Shen, Anoop Sarkar, and Franz Josef Och. 2004.	925
874	<i>tional Linguistics</i> , pages 311–318, Philadelphia,	Discriminative reranking for machine translation .	926
875	Pennsylvania, USA. Association for Computational	In <i>Proceedings of the Human Language Technol-</i>	927
876	Linguistics.	<i>ogy Conference of the North American Chapter</i>	928
		<i>of the Association for Computational Linguistics:</i>	929
		<i>HLT-NAACL 2004</i> , pages 177–184, Boston, Mas-	930
		sachusetts, USA. Association for Computational Lin-	931
		guistics.	932
877	Maja Popović. 2015. chrF: character n-gram F-score	Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua	933
878	for automatic MT evaluation . In <i>Proceedings of the</i>	Wu, Maosong Sun, and Yang Liu. 2016. Minimum	934
879	<i>Tenth Workshop on Statistical Machine Translation</i> ,	risk training for neural machine translation . In <i>Pro-</i>	935
880	pages 392–395, Lisbon, Portugal. Association for	<i>ceedings of the 54th Annual Meeting of the Associa-</i>	936
881	Computational Linguistics.	<i>tion for Computational Linguistics (Volume 1: Long</i>	937
		<i>Papers)</i> , pages 1683–1692, Berlin, Germany. Associ-	938
		ation for Computational Linguistics.	939
882	Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian	Raphael Shu and Hideki Nakayama. 2017a. Later-stage	940
883	Gehrmann, and Thibault Sellam. 2021a. Learning	minimum bayes-risk decoding for neural machine	941
884	compact metrics for MT . In <i>Proceedings of the 2021</i>	translation . <i>arXiv preprint arXiv:1704.03169</i> .	942
885	<i>Conference on Empirical Methods in Natural Lan-</i>		
886	<i>guage Processing</i> , pages 751–762, Online and Punta	Raphael Shu and Hideki Nakayama. 2017b. Later-stage	943
887	Cana, Dominican Republic. Association for Compu-	minimum bayes-risk decoding for neural machine	944
888	tational Linguistics.	translation .	945
889	Amy Pu, Hyung Won Chung, Ankur P. Parikh, Sebastian	David A. Smith and Jason Eisner. 2006. Minimum	946
890	Gehrmann, and Thibault Sellam. 2021b. Learning	risk annealing for training log-linear models . In <i>Pro-</i>	947
891	compact metrics for mt .	<i>ceedings of the COLING/ACL 2006 Main Conference</i>	948
		<i>Poster Sessions</i> , pages 787–794, Sydney, Australia.	949
		Association for Computational Linguistics.	950
892	Tharindu Ranasinghe, Constantin Orasan, and Ruslan	Lucia Specia, Frédéric Blain, Marina Fomicheva, Er-	951
893	Mitkov. 2020. TransQuest: Translation quality esti-	rick Fonseca, Vishrav Chaudhary, Francisco Guzmán,	952
894	mation with cross-lingual transformers . In <i>Proceed-</i>	and André F. T. Martins. 2020. Findings of the WMT	953
895	<i>ings of the 28th International Conference on Com-</i>	2020 shared task on quality estimation . In <i>Proceed-</i>	954
896	<i>putational Linguistics</i> , pages 5070–5081, Barcelona,	<i>ings of the Fifth Conference on Machine Translation</i> ,	955
897	Spain (Online). International Committee on Compu-	pages 743–764, Online. Association for Computa-	956
898	tational Linguistics.	tional Linguistics.	957
899	Ricardo Rei, Ana C Farinha, Chrysoula Zerva, Daan	Lucia Specia, Frédéric Blain, Marina Fomicheva,	958
900	van Stigt, Craig Stewart, Pedro G Ramos, Taisiya	Chrysoula Zerva, Zhenhao Li, Vishrav Chaudhary,	959
901	Glushkova, André Martins, and Alon Lavie. 2021.	and André F T Martins. 2021. Findings of the WMT	960
902	Are References Really Needed? Unbabel-IST 2021	2021 Shared Task on Quality Estimation . pages 667–	961
903	Submission for the Metrics Shared Task . In <i>Proceed-</i>	708.	962
904	<i>ings of the Sixth Conference on Machine Translation</i> ,		
905	Online. Association for Computational Linguistics.	Felix Stahlberg and Bill Byrne. 2019. On NMT search	963
		errors and model errors: Cat got your tongue? In	964
906	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	<i>Proceedings of the 2019 Conference on Empirical</i>	965
907	Lavie. 2020a. COMET: A neural framework for MT	<i>Methods in Natural Language Processing and the</i>	966
908	evaluation . In <i>Proceedings of the 2020 Conference</i>	<i>9th International Joint Conference on Natural Lan-</i>	967
909	<i>on Empirical Methods in Natural Language Process-</i>	<i>guage Processing (EMNLP-IJCNLP)</i> , pages 3356–	968
910	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Association	3362, Hong Kong, China. Association for Computa-	969
911	for Computational Linguistics.	tional Linguistics.	970
912	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon		
913	Lavie. 2020b. Unbabel’s participation in the WMT20		
914	metrics shared task . In <i>Proceedings of the Fifth Con-</i>		
915	<i>ference on Machine Translation</i> , pages 911–920, On-		
916	line. Association for Computational Linguistics.		

971	Felix Stahlberg, Adrià de Gispert, Eva Hasler, and Bill	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.	1029
972	Byrne. 2017. Neural machine translation by min-	Weinberger, and Yoav Artzi. 2020. Bertscore: Eval-	1030
973	imising the Bayes-risk with respect to syntactic trans-	uating text generation with bert. In <i>International</i>	1031
974	lation lattices. In <i>Proceedings of the 15th Confer-</i>	ence on Learning Representations.	1032
975	<i>ence of the European Chapter of the Association</i>		
976	<i>for Computational Linguistics: Volume 2, Short Pa-</i>		
977	<i>pers</i> , pages 362–368, Valencia, Spain. Association		
978	for Computational Linguistics.		
979	Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Se-		
980	quence to sequence learning with neural networks. In		
981	<i>Advances in Neural Information Processing Systems</i> ,		
982	volume 27. Curran Associates, Inc.		
983	Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe,		
984	Jon Shlens, and Zbigniew Wojna. 2016. Rethink-		
985	ing the inception architecture for computer vision.		
986	In <i>2016 IEEE Conference on Computer Vision and</i>		
987	<i>Pattern Recognition (CVPR)</i> , pages 2818–2826.		
988	Chaojun Wang and Rico Sennrich. 2020. On exposure		
989	bias, hallucination and domain shift in neural ma-		
990	chine translation. In <i>Proceedings of the 58th Annual</i>		
991	<i>Meeting of the Association for Computational Lin-</i>		
992	<i>guistics</i> , pages 3544–3552, Online. Association for		
993	Computational Linguistics.		
994	John Wieting, Taylor Berg-Kirkpatrick, Kevin Gimpel,		
995	and Graham Neubig. 2019. Beyond BLEU:training		
996	neural machine translation with semantic similarity.		
997	In <i>Proceedings of the 57th Annual Meeting of the As-</i>		
998	<i>sociation for Computational Linguistics</i> , pages 4344–		
999	4355, Florence, Italy. Association for Computational		
1000	Linguistics.		
1001	Yilin Yang, Liang Huang, and Mingbo Ma. 2018. Break-		
1002	ing the beam search curse: A study of (re-)scoring		
1003	methods and stopping criteria for neural machine		
1004	translation. In <i>Proceedings of the 2018 Conference</i>		
1005	<i>on Empirical Methods in Natural Language Process-</i>		
1006	<i>ing</i> , pages 3054–3059, Brussels, Belgium. Associa-		
1007	tion for Computational Linguistics.		
1008	Kyra Yee, Yann Dauphin, and Michael Auli. 2019. Simple		
1009	and effective noisy channel modeling for		
1010	neural machine translation. In <i>Proceedings of the</i>		
1011	<i>2019 Conference on Empirical Methods in Natu-</i>		
1012	<i>ral Language Processing and the 9th International</i>		
1013	<i>Joint Conference on Natural Language Processing</i>		
1014	<i>(EMNLP-IJCNLP)</i> , pages 5696–5701, Hong Kong,		
1015	China. Association for Computational Linguistics.		
1016	Lei Yu, Phil Blunsom, Chris Dyer, Edward Grefenstette,		
1017	and Tomáš Kociský. 2017. The neural noisy chan-		
1018	nel. In <i>5th International Conference on Learning</i>		
1019	<i>Representations, ICLR 2017, Toulon, France, April</i>		
1020	<i>24-26, 2017, Conference Track Proceedings.</i> Open-		
1021	Review.net.		
1022	Chrysoula Zerva, Daan van Stigt, Ricardo Rei, Ana		
1023	C Farinha, José G. C. de Souza, Taisiya Glushkova,		
1024	Miguel Vera, Fabio Kepler, and André Martins. 2021.		
1025	IST-Unbabel 2021 Submission for the Quality Es-		
1026	timation Shared Task. In <i>Proceedings of the Sixth</i>		
1027	<i>Conference on Machine Translation</i> , Online. Associa-		
1028	tion for Computational Linguistics.		

Supplemental Material

1033

A Training Details

1034

For the experiments using IWSLT17, we train a *small* transformer model (6 layers, 4 attention heads, 512 embedding dimensions, and 1024 hidden dimensions) from scratch, using *Fairseq* (Ott et al., 2019). We tokenize the data using SentencePiece (Kudo and Richardson, 2018), with a joint vocabulary with 20000 units. We train using the Adam optimizer (Kingma and Ba, 2015) with $\beta_1 = 0.9$ and $\beta_2 = 0.98$ and use an inverse square root learning rate scheduler, with an initial learning rate of 5×10^{-4} and with a linear warm-up in the first 4000 steps. For models trained with label smoothing, we use the default value of 0.1.

1035

1036

1037

1038

1039

1040

B Additional Results

1041

For completeness, we include in Table 3 results to evaluate the impact of the metrics presented in §3 as features and objectives for ranking using the other language pairs: EN \rightarrow RU (large model) and EN \rightarrow FR (small model).

1042

1043

1044

	Large (WMT20)				Small (IWSLT)			
	BLEU	chrF	BLEURT	COMET	BLEU	chrF	BLEURT	COMET
Baseline	23.86	51.16	0.6953	0.5361	38.12	63.20	0.6532	0.4809
F-RR w/ COMET-QE	20.32	49.18	<u>0.7130</u>	<u>0.6207</u>	35.59	60.90	0.6628	<u>0.5553</u>
F-RR w/ MBART-QE	<u>22.39</u>	<u>50.59</u>	0.6993	0.5481	<u>36.68</u>	<u>62.17</u>	0.6593	0.5091
F-RR w/ OpenKiwi	20.88	48.72	0.7040	0.5688	32.03	55.68	0.5996	0.2581
F-RR w/ Transquest	21.60	50.14	0.7060	0.5836	36.02	62.26	<u>0.6681</u>	0.5397
T-RR w/ BLEU	<u>23.87</u>	51.51	0.7042	0.5669	39.10	64.22	0.6968	0.6189
T-RR w/ BLEURT	22.84	51.25	<u>0.7265</u>	<u>0.6470</u>	38.60	63.76	<u>0.7042</u>	<u>0.6405</u>
F-RR w/ COMET	22.42	50.91	0.7243	<u>0.6441</u>	38.60	63.77	<u>0.7020</u>	0.6392
MBR w/ BLEU	<u>24.03</u>	51.12	0.6938	0.5393	<u>37.97</u>	63.13	0.6484	0.4764
MBR w/ BLEURT	23.01	50.87	<u>0.7314</u>	0.5984	37.29	62.82	<u>0.6886</u>	0.5361
MBR w/ COMET	23.67	<u>51.18</u>	0.7093	<u>0.6242</u>	37.77	<u>63.24</u>	<u>0.6710</u>	<u>0.6127</u>
T-RR+MBR w/ BLEU	24.11	<u>51.44</u>	0.6967	0.5482	38.96	<u>64.04</u>	0.6781	0.5636
T-RR+MBR w/ BLEURT	23.18	51.30	0.7344	<u>0.6277</u>	37.43	63.14	0.7092	0.5961
T-RR+MBR w/ COMET	23.21	51.26	0.7238	0.6736	38.33	63.53	0.6883	0.6610

Table 3: Evaluation metrics for EN \rightarrow RU for the *large* model setting and EN \rightarrow FR for *small* model settings, using a *fixed* N -best reranker (F-RR), a *tuned* N -best reranker (T-RR), MBR decoding, and a two-stage approach.

C Human Study

1045

In order to perform human evaluation, we recruited professional translators who were native speakers of the target language on the freelancing site Upwork.⁴ 300 sentences were evaluated for each language pair, sampled randomly from the test sets after a restriction that sentences were no longer than 30 words. All translation hypotheses for a single source sentence were first deduplicated, and then shown to the translator side-by-side in randomized order to avoid any ordering biases.

1046

1047

1048

1049

1050

Sentences were evaluated according to a 1-5 rubric slightly adapted from that of Wieting et al. (2019):

1051

1. There is no overlap in the meaning of the source sentence whatsoever.

1052

2. Some content is similar but the most important information in the sentence is different.

1053

3. The key information in the sentence is the same but the details differ.

1054

4. Meaning is essentially equal but some expressions are unnatural.

1055

5. Meaning is essentially equal and the sentence is natural.

1056

⁴<https://upwork.com>. Freelancers were paid a market rate of 18-20 US dollars per hour, and finished approximately 50 sentences in one hour.