Translate Smart, not Hard: Cascaded Translation Systems with Quality-Aware Deferral

António Farinhas 1* , Nuno M. Guerreiro 1† , Sweta Agrawal 2* Ricardo Rei 1† , André F.T. Martins 3,4,5,6†

¹Sword Health, ²Google, ³Instituto de Telecomunicações ⁴Instituto Superior Técnico, Universidade de Lisboa, ⁵TransPerfect, ⁶ELLIS Unit Lisbon antonio.farinhas@tecnico.ulisboa.pt

Abstract

Larger models often outperform smaller ones but come with high computational costs. Cascading offers a potential solution. By default, it uses smaller models and defers only some instances to larger, more powerful models. However, designing effective deferral rules remains a challenge. In this paper, we propose a simple yet effective approach for machine translation, using existing quality estimation (QE) metrics as deferral rules. We show that QE-based deferral allows a cascaded system to match the performance of a larger model while invoking it for a small fraction (30% to 50%) of the examples, significantly reducing computational costs. We validate this approach through both automatic and human evaluation.

1 Introduction

Larger models consistently outperform smaller ones in NLP tasks, but the trade-off is the increased computational cost (Brown et al., 2020; Kaplan et al., 2020; Chowdhery et al., 2023). This raises the question:

How can we maintain high performance while reducing computational load?

A promising solution is **model cascading**, where smaller models handle examples by default, and only a subset of hard instances is deferred to a larger model. However, this approach requires a robust deferral system that reliably determines when to defer. Common approaches often involve designing and training specialized deferral models, which determine when a large model is needed—*e.g.*, based on reliability or uncertainty estimates (Chen et al., 2023b; Gupta et al., 2024). But do we really need to train new models for every task, or can existing resources speed up this process?

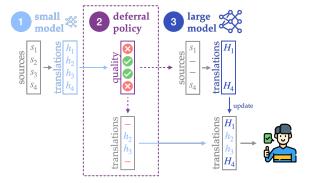


Figure 1: Cascaded translation system with quality-aware deferral. A small model translates a batch of source sentences, and a relatively lightweight QE model scores the hypotheses. Sources with the lowest-scoring translations are deferred to a larger model. By default, the deferral rate is set by a predefined compute budget, but the threshold can also be computed dynamically without requiring full batches (§5.5).

For machine translation (MT), extensive research on reference-free automatic evaluation offers an appealing alternative (Zerva et al., 2022, 2024; Blain et al., 2023). In this paper, we leverage recent quality estimation (QE) metrics to create straightforward and relatively lightweight deferral rules. This approach draws inspiration from professional translation workflows, where QE metrics help identify translations that should be deferred to expert post-editing (Castilho and O'Brien, 2017; Béchara et al., 2021). Our main contributions are:

- We introduce a **cascaded translation system** that uses pretrained QE metrics to determine whether to defer examples from a smaller model to a larger one, balancing efficiency and quality (§3). See Fig. 1 for an illustration.
- We confirm that the benefits of QE-based model cascading hold across different combinations of translation and QE models (§5).
- We perform human evaluation, further val-

^{*}Work done while at Instituto de Telecomunicações.

[†]Work done while at Unbabel.

idating our approach on two language pairs (en-es and en-ja) in the WMT24 test set (§6).

2 Adaptive Inference in NLP

Adaptive inference techniques are increasingly being adopted in natural language processing tasks (Mamou et al., 2022; Varshney and Baral, 2022; Chen et al., 2023b; Ong et al., 2024). These methods typically use models of different sizes and predictive power (often two, though most frameworks can easily accommodate more), with the primary goal of reducing the computational load by using the larger, more computationally expensive model only when necessary (e.g., for more difficult examples or when a model is highly uncertain about its prediction). Current strategies include **routing**, where a decision rule determines which model to use, ensuring only one model is used to handle each input, and cascading, which starts with a smaller model and may invoke a larger one afterward based on the small model's output and a deferral rule. In this paper, we focus on the second approach.

The computational efficiency of model cascading comes at the cost of designing a **robust deferral system** that can reliably identify when to defer to the larger model. This is often handled using simple decision rules, such as nonparametric methods or other approaches based on uncertainty measures (Ramírez et al., 2024; Gupta et al., 2024). A recent alternative involves training external models specifically to predict when deferral is needed—for a given example, these models can be trained, *e.g.*, to assess if a given candidate is correct (Chen et al., 2023b). Here, we propose a simple and effective deferral rule for MT that is conceptually similar to this approach while offering a particularly straightforward solution for this task.

Complementary to routing and cascading is the line of research on speculative decoding (Stern et al., 2018; Chen et al., 2023a; Leviathan et al., 2023; Xia et al., 2024), where a small model drafts multiple tokens in advance, and a large model verifies them in parallel. Only the tokens that pass this verification are accepted as final outputs. While speculative decoding is typically applied to accelerate autoregressive generation, its principles can be combined with cascading to further improve efficiency (Narasimhan et al., 2025).

3 Quality-Aware Deferral for MT

Although human evaluations and reference-based metrics remain the standard for evaluating machine translations, reference-free/quality estimation (QE) metrics have shown strong correlations with human judgments (Zerva et al., 2024), holding promise in distinguishing between the quality of translations for the same source (Agrawal et al., 2024). Since QE models are typically much smaller than current translation models (Kocmi et al., 2024a), we propose to leverage them for an efficient deferral rule. Rather than training new bespoke decision models (§2), existing QE models can evaluate translations from a lightweight model and determine when to accept them or defer to a larger one.

How to choose which examples to defer? Setting a fixed threshold on QE scores is challenging too high a threshold wastes computational resources, while too low a threshold risks compromising quality (Jitkrittum et al., 2023; Gupta et al., 2024). Throughout this paper, we use a **budget**constrained computation approach: we first translate all examples in a batch with the smaller model, then rank them based on QE scores, deferring only the lowest-scoring subset according to a predefined compute budget (the fraction of examples deferred to the larger model). This assumes parallel processing of entire batches rather than processing individual instances sequentially (see Fig. 1 for an illustration with 50% of deferral). However, realworld systems may process inputs sequentially or in a streaming fashion. To mimic such settings, we also explore a dynamic thresholding approach in §5.5. Inspired by Ramírez et al. (2024), this approach sets and updates the QE threshold on the fly, using the distribution of scores observed so far—allowing adaptive, approximate control over deferral without access to full batches upfront.

Computational efficiency. The standard approximation for the number of floating point operations (FLOPs) required for inference with a transformer model is 2ND, where N represents the number of model parameters and D is the number of tokens generated at inference time (Sardana et al., 2024; Snell et al., 2024). For a cascaded approach with superscripts S and L denoting the smaller and larger models, respectively, this becomes:

$$2BD_S(N_S + N_{QE}) + 2\eta BD_L N_L, \qquad (1)$$

¹Likewise, routing typically involves training external models to (*i*) predict the performance of the small model (Šakota et al., 2024), or (*ii*) determine if the small model is likely to outperform the large one (Ding et al., 2024).

where B is the batch size, N_{QE} is the number of parameters of the QE model, and η is the proportion of instances the larger model processes. Assuming $D_S \approx D_L$, this approach achieves computational parity with the larger model (i.e., $2BDN_L$) when:

$$\eta^{\star} = 1 - \frac{N_S + N_{QE}}{N_L}.\tag{2}$$

This expression provides a simple rule of thumb: to maintain computational efficiency, the larger model should handle at most η^{\star} of the examples. For instance, if it is $10\times$ larger than the smaller model and the QE model is negligible ($N_{QE}\ll N_S$), then $\eta^{\star}\approx 0.9$. This means the cascading is more efficient than always using the larger model as long as fewer than 90% of the examples are deferred.

4 Experimental Setup

4.1 Generation models

Through the paper, we experiment with generation models of different size and predictive power:

- Tower-v2 70B (Rei et al., 2024): With 70B parameters, this is an improved iteration of Tower (Alves et al., 2024), obtained by continued pertaining Llama-3 (AI@Meta, 2024) on a multilingual dataset with 25 billions of tokens, followed by supervised finentuning for translation-related tasks. Compared to the first iteration of Tower, this model is better at paragraph and document-level translation and supports more languages (15, instead of 10), including all the languages in the WMT24 test sets. Combined with quality-aware decoding (Fernandes et al., 2022), this is the winning submission of the WMT24 general translation shared task (Kocmi et al., 2024a).
- Tower-v2 7B (Rei et al., 2024): A more lightweight version of Tower-v2 70B using Mistral instead (Jiang et al., 2023).
- Tower-v2 7B (L): We follow the recipe described above to train a smaller version of Tower-v2 70B based on Llama-3. This model slightly underperforms its Mistral counterpart.
- EuroLLM Instruct (9B and 1.7B) (Martins et al., 2024): EuroLLM models are openweight multilingual models trained on 4 trillion tokens covering all European Union and many other relevant languages across several

	M ↑	C ↑	Win rate
Tower-v2 7B	-3.01	83.94	43% 32%
Tower-v2 7B (L)	-3.07	83.73	45% 32%
EuroLLM 9B	-4.01	80.56	52% 28%
EuroLLM 1.7B	-4.60	77.42	66% 20%
Tower-v2 70B	-2.79	84.71	NA

Table 1: Translation quality measured with METRICX (M) and COMET (C) on the WMT24 test set. Win rates against Tower-v2 70B, according to M. The bars represent the proportions of losses, ties, and wins.

data sources (web data, parallel data, and highquality datasets). The instruction-tuned models are obtained after finetuning the base models on the EuroBlocks dataset, which includes general instruction-following and MT tasks.

We generate all translations with greedy decoding using vLLM (Kwon et al., 2023) for faster inference. Table 1 shows the performance of these models on the WMT24 test sets (Kocmi et al., 2024a),² according to METRICX and COMET (results are averaged across all language pairs), along with win rates against Tower-v2 70B. Following Kocmi et al. (2024b), translations with differences in METRICX below 0.122 are considered ties when comparing two systems (90% human accuracy). We use this threshold for detecting ties at the segment level.

4.2 Deferral strategy and baselines

We use two versions of COMETKIWI: wmt22-cometkiwi-da (Rei et al., 2022b), which with only 0.5B parameters achieves a strong correlation with human judgments (Zerva et al., 2022); and wmt23-cometkiwi-da-xxl (Rei et al., 2023), a scaled version with 10.5B parameters. As baselines, we compare against several simple, yet often effective, heuristics:

- random selection, which uniformly defers a random subset of examples;
- source length computed using Tower-v2's tokenizer, *i.e.*, deferring either the shortest (length) or the longest (-length) sources—source length is a common proxy for translation difficulty (Kocmi and Bojar, 2017; Wan et al., 2022; Wang et al., 2023), as longer texts are typically harder to translate because

²Publicly available for research purposes at https://www2.statmt.org/wmt24/translation-task.html.
Our use of datasets and models aligns with their intended purposes as defined by the licenses.

they require maintaining consistency and adequacy for larger contexts, and most systems are trained with sentence-level data;

• logprobs, which uses the smaller model's normalized log-probability (directly obtained as a by-product of the translation), *i.e.*, deferring texts with the lowest likelihoods, which is typically a helpful quality proxy, as shown by Fomicheva et al. (2020) and Guerreiro et al. (2023) for hallucination detection.

These heuristics are not only computationally inexpensive but have also been shown in prior work to be effective, often outperforming more sophisticated methods. We also report results with oracle deferral, which maximizes translation quality according to humans in §6, and compare our approach with quality-aware decoding (Fernandes et al., 2022), which selects among multiple hypotheses produced by the same model in §7.

4.3 Evaluation

We use the WMT24 test sets (Kocmi et al., 2024a), which span multiple domains (news, social, speech, and literary) and 11 language pairs (en-cs, en-de, en-es, en-hi, en-is, en-ja, en-ru, en-uk, en-zh, cs-uk, and ja-zh). For each language pair, we treat the full test set as a single batch for computing QE thresholds (§3), unless otherwise stated. Results are then averaged across language pairs for better visualization. We evaluate systems with METRICX (Juraska et al., 2023) to reduce the risk of "reward hacking" (Fernandes et al., 2022) and better reflect real quality improvements (we also provide COMET (Rei et al., 2022a) scores in App. A). Since biases may still exist when using a different evaluation metric than the reward model (Kovacs et al., 2024), we also conduct a human study (§6) using DA+SQM (direct assessment + scalar quality metric) source contrastive evaluation (Kocmi et al., 2022).

5 Results and Analysis

In this section, we present a series of analyses to evaluate the effectiveness of quality-aware deferral. We first show that larger models are not always better than smaller ones (§5.1), motivating the use of cascaded approaches. We then demonstrate that quality-aware deferral effectively balances performance and efficiency (§5.2), and that these gains hold across different QE models (§5.3) and small model backbones (§5.4). Finally, we explore how

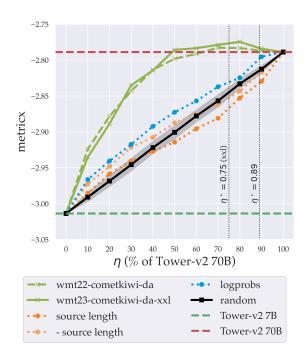


Figure 2: Translation quality of cascading combining Tower-v2 7B and Tower-v2 70B according to METRICX, as the inference computation budget varies. Horizontal lines show the performance of each model alone.

our method performs under streaming constraints where deferral must be made dynamically (§5.5).

5.1 Larger is not necessarily better

Although Tower-v2 70B outperforms Tower-v2 7B across all language pairs (Table 1 shows aggregated results), a closer look at its win rates shows it only outperforms the smaller model in 43% of individual examples. Even when compared to the significantly smaller EuroLLM 1.7B, Tower-v2 70B wins only 66% of the time. This is not surprising, as EuroLLM was trained explicitly with large amounts of translation data to perform well on translation tasks, despite its smaller size (Martins et al., 2024). This demonstrates that larger models are not consistently superior at the instance level (Table 2 provides some examples), opening the possibility of strategically using smaller models for a subset of examples. This approach can improve efficiency without compromising translation quality.

5.2 QE is an effective deferral rule

Fig. 2 shows the quality of a cascaded system combining Tower-v2 7B and Tower-v2 70B according to METRICX under varying inference budgets (results are averaged across language pairs). Each curve represents a different deferral rule. The dotted vertical lines show the values of η after which

	Example 1: Tower-v2 70B is better	Notes
Source	"It's your birthday soon, isn't it?" Ivory asked, remembering that the princess' birthday was in a few days. Kari nodded, a sad glint in her light blue eyes.	
Tower-v2 7B	—¿No es pronto tu cumpleaños, verdad?—preguntó Ivory, acordándose de que el cumpleaños de la princesa era en unos días. Kari asintió, con un brillo triste en sus ojos azul claro.	33/100, awkward phrasing, less natural structure
Tower-v2 70B	«Pronto es tu cumpleaños, ¿no es así?», preguntó Ivory, recordando que el cumpleaños de la princesa era en unos pocos días. Kari asintió, con un brillo triste en sus ojos azul claro.	100/100, more fluid and literary
	Example 2: Tower-v2 7B is better	Notes
Source Tower-v2 7B Tower-v2 70B	A quarter of the way through the year update well ahead of pace! ¡Ya llevamos una cuarta parte de la actualización anual ¡muy por encima del ritmo previsto! ¡Ya hemos avanzado un cuarto del camino en la actualización del año ¡y estamos muy por delante del ritmo!	100/100, idiomatic phrasing 67/100, slightly literal, less natural
	Example 3: Tower-v2 7B is better	Notes
Source Tower-v2 7B Tower-v2 70B	I'll keep posting my bakes to my lovely masto-peeps, as they gobble it up. 美味しそうに食べてくれるマストフォロワーの皆さんに、これからもお菓子の投稿を続けていきます。 皆さんに喜んでいただけるので、美味しそうな焼き菓子をマストドンに投稿し続けます。	100/100, conversational 33/100, more formal

Table 2: Illustrative examples where Tower-v2 70B is either better or worse than Tower-v2 7B, based on human quality assessments. These examples show that the larger model does not always produce the preferred output. Key differences are <u>underlined</u>.

our method would become more expensive than always running the larger model only, depending on the size of the QE model.

As expected, the random baseline fails to identify examples that benefit from larger models, resulting in suboptimal performance. Source length-based decision rules or using the small model's logprobs perform slightly better or worse than random, suggesting that simple heuristics cannot capture finegrained differences in translation quality and are inefficient for deferral. In contrast, QE-based deferral (our proposal) achieves the best overall performance, enabling the cascaded system to match the performance of the large model while invoking it for only 50% to 60% of the examples. From Eq. (2), computational parity is reached at $\eta^{\star} = 89\%$ when using wmt22-cometkiwi-da $(N_Q=0.5B)$ and $\eta^{\star} = 75\%$ with wmt23-cometkiwi-da-xxl $(N_Q = 10.5B)$. Matching Tower 70B's performance at such a small η shows that our approach effectively balances efficiency and quality.

5.3 What if we use another QE model?

We have seen that QE-based cascading works well with COMETKIWI models of different sizes (Fig. 2). Here, we show that this is also the case when using two reference-free versions of METRICX (Juraska et al., 2024): metricx-24-hybrid-large-v2p6, with 1.2B

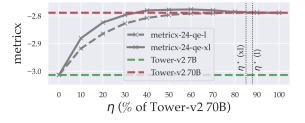


Figure 3: Translation quality of cascaded systems with deferral based on metricx-24-hybrid-large-v2p6 and metricx-24-hybrid-xl-v2p6.

parameters, and metricx-24-hybrid-x1-v2p6, with 3.7B parameters (Fig. 3, gray curves). Later in §6, we confirm this with human evaluation.

5.4 How does the quality of the small model impact performance?

We have shown that QE-based cascading works well across QE models of different sizes (Fig. 2). Here, we study whether it still provides gains when the smaller model is relatively weaker. We train another version of Tower 7B using Llama-3 instead of Mistral, referred to as **Tower 7B** (L), and use two versions of **EuroLLM** (Martins et al., 2024) with 1.7B ($\eta^* = 0.97$) and 9B parameters ($\eta^* = 0.86$). Fig. 4 shows that while these models underperform Tower-v2 7B, cascading with Tower 70B remains competitive. This indicates that QE-based cascad-

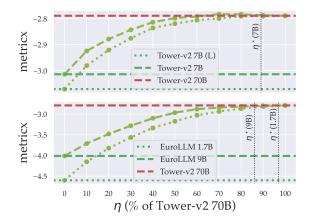


Figure 4: Translation quality of cascaded systems with deferral based on wmt22-cometkiwi-da. Large model: Tower-v2 70B. Small models: Tower-v2 7B (L), Tower-v2 7B (top); EuroLLM 1.7B, EuroLLM 9B (bottom).

ing is robust across different generation models, even when both belong to the same family (top) or when the small model is much smaller (bottom). However, when the win rates of the larger model against the smaller model increase, we can expect η to also increase.

5.5 What if deferral must be made on the fly?

So far, we have presented results by varying the percentage of deferred sentences, which is equivalent to thresholding QE scores with different cutoff values. However, in some real-world cases (e.g., real-time streaming scenarios), we may not have immediate access to an entire batch of sentences. To address this, we use a dynamic thresholding method similar to that of Ramírez et al. (2024), where the threshold is updated incrementally during inference. Specifically: i) for the first B examples, we route all inputs to the smaller model and compute their QE scores; ii) we then set an initial threshold based on the score that would defer a target percentage of examples (e.g., 40%) from that block; iii) after each subsequent block of B examples, we update the threshold using all previously computed QE scores to maintain the desired deferral rate. This mirrors a real-time scenario where the system continuously adapts its deferral decisions over time, aiming to approximate the target deferred percentage based on the sentences processed so far. Unlike our earlier budget-constrained computation setup, this approach does not guarantee the exact deferred percentage.

Fig. 5 shows that similar curves can be obtained using a dynamic thresholding setup with B=1

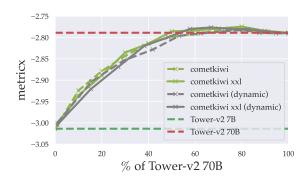


Figure 5: Translation quality of cascaded systems with deferral based on wmt22-cometkiwi-da and wmt23-cometkiwi-da-xxl, using a budget-constrained or a dynamic procedure.

(gray curves), without requiring access to the full batch of source sentences. For the rest of the paper, we use the budget-constrained setup.

6 Human Evaluation

Since using QE metrics during inference can bias automatic evaluations, we conduct a human study to obtain reference-quality translation judgments and validate our approach. We recruited professional translators who were native speakers of the target language on the freelancing site Upwork.³ We followed a DA+SQM (direct assessment + scalar quality metric) source contrastive evaluation (Kocmi et al., 2022) using Appraise (Federmann, 2018). We randomly sampled 500 source instances from the WMT24 test set for en-ja and en-es and asked one translator per language pair to read two alternative translations for each source and evaluate them on a continuous scale from 0 (no overlap in meaning) to 100 (perfect translation). The scale featured seven labeled tick marks (from 0 to 6) indicating different quality labels combining accuracy and grammatical correctness. Translators could further adjust their scores to reflect preferences or assign the same score to translations of similar quality. They were paid a market rate of around 20 USD per hour, and completing the task took approximately 12 to 14 hours for each language pair. Further details are in App. B.

6.1 QE remains an effective deferral rule

Fig. 6 shows the performance of cascaded systems using QE-based deferral. We use a paired-permutation test (Good, 2000; Zmigrod et al., 2022) to compare the performance of Tower-v2 70B with

³https://upwork.com.

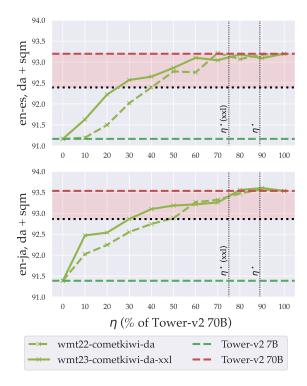


Figure 6: Translation quality of a cascaded system combining Tower-v2 7B and Tower-v2 70B according to human scores (in a scale from 0 to 100), as the inference computation budget varies. Systems in the shaded area are not significantly different from Tower-v2 70B according to the paired-permutation test with p=0.01.

each our cascaded systems under varying budgets (represented by the green lines and marked with an "x"). Using p=0.01, we found that all the systems in the shaded area are not significantly different from Tower-V2 70B, whereas systems below the shaded area (*i.e.*, below 92.4 and 92.9 for enes and en-ja, respectively) are significantly worse. This shows that our approach achieves performance comparable to Tower-v2 70B while invoking it for only 30% to 50% of the examples, confirming that it substantially reduces computational costs without compromising translation quality.

6.2 What if we use another QE model? What if we had a perfect QE model?

The effectiveness of our framework depends on the quality of existing QE models, and improving them can further strengthen our approach. First, we confirm with human evaluation that the benefits of quality-aware deferral go beyond COMETKIWI models, as explained in §5.2 (check the orange

curves in Fig. 7 for results using two reference-free versions of METRICX). Second, to access the performance ceiling of cascading, we report results with oracle deferral, *i.e.*, a deferral strategy that maximizes translation quality according to humans (Fig. 7, black curves).⁵ The high oracle values indicate significant potential for improvement, suggesting that having better QE models could directly boost the effectiveness of our cascaded approach.

7 Comparison with Quality-Aware Decoding

There is a large body of work on **reranking for lan**guage generation, where we start by generating multiple hypotheses with a language model, and then use a reranker to select the best one (Farinhas et al., 2024). For MT, an example is quality-aware decoding (Fernandes et al., 2022). The cheapest approach is QE reranking, where we first generate multiple translation hypotheses and then rerank them using a QE model. This strategy is often used to reduce the propensity of language models to hallucinate (Guerreiro et al., 2023; Farinhas et al., 2023). While our approach is conceptually different—designed with efficiency in mind, whereas QE reranking is often computationally expensive—there is also a key structural distinction: QE reranking selects among multiple hypotheses produced by the same model, whereas our deferral strategy uses a single hypothesis from the smaller model and invokes the larger model only when the predicted quality falls below a certain threshold.

Computational efficiency. Following the discussion in §3, the number of FLOPS required for inference with a large model on a batch of B examples is $2BDN_L$, where N_L represents the number of model parameters and D is the number of generated tokens. In this section, we assume that our goal is to **reduce the computational cost by** (1-X)%, meaning that we operate under a computational budget of $X \cdot 2BDN_L$. The number of FLOPs required to run inference with our cascaded approach is given by:

$$2BD(N_S + N_{QE} + \eta N_L), \tag{3}$$

which leads to the following expression for X:

$$X = \eta + \frac{N_S + N_{QE}}{N_L}. (4)$$

⁴We run the same statistical test to compare the small model (Tower-v2 7B) with our cascaded systems, concluding that all systems within the shaded region (including the large model itself) are also significantly better than Tower-v2 7B.

⁵Oracle performance goes down after reaching a *plateau* due to our budget-constrained approach, which enforces deferral for a fixed percentage of examples.

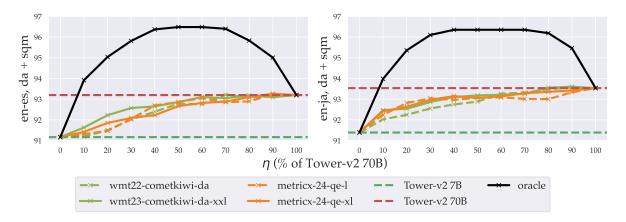


Figure 7: Translation quality of a cascaded system combining Tower-v2 7B and Tower-v2 70B according to human scores (in a scale from 0 to 100), as the inference computation budget varies. Deferral is based on different QE models (green and orange curves). The black curve shows the oracle selection.

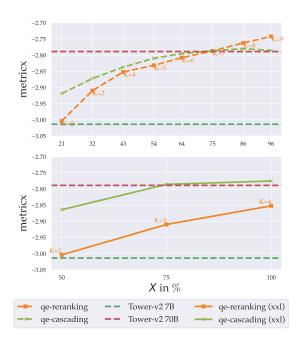


Figure 8: Translation quality of a cascaded system combining Tower-v2 7B and Tower-v2 70B (in green) v.s. QE reranking with hypotheses generated by Tower-v2 7B (in orange), measured with METRICX, as X varies.

For QE reranking, the computational cost is:

$$2BDK(N_S + N_{QE}), (5)$$

where K is the number of generated hypotheses. This yields:

$$X = K \cdot \left(\frac{N_S + N_{QE}}{N_L}\right). \tag{6}$$

These expressions allow us to obtain the values of η for which our approach incurs the same computational cost as QE reranking with K hypotheses:

$$\eta = (K - 1) \cdot \left(\frac{N_S + N_{QE}}{N_L}\right). \tag{7}$$

Experiments and discussion. We generate up to 9 hypotheses with Tower-v2 7B using ϵ -sampling with $\epsilon = 0.02$ (Freitag et al., 2023). cording to Eq. (6), the number of FLOPs required for QE reranking with more than 9 hypotheses already exceeds the budget of $2BDN_L$ if we use wmt22-cometkiwi-da. When using wmt23-cometkiwi-da-xxl, computational parity is achieved with K=4. Fig. 8 illustrates the tradeoff between computational efficiency and translation quality (measured with METRICX) for a cascaded approach with QE-based deferral against QE reranking. As expected, quality improves as the computational budget increases for both methods. While QE reranking is also an effective way to improve translation quality when generating multiple hypotheses is feasible, our cascaded approach achieves better quality at lower computation costs, making it a more efficient alternative when computational efficiency is a priority.

8 Conclusions and Future Work

We propose a simple yet effective approach to model cascading for MT using QE metrics for deferral. Our method matches the quality of larger models while requiring them to handle only a subset of examples, significantly reducing computational costs. This is shown through automatic and human evaluations. The effectiveness of our framework depends on the quality of existing QE models, and improving them can further strengthen our approach.

Limitations

We highlight the main limitations of our work. First, we focus on a two-stage cascade, where examples are handled by a small model or deferred to a larger one. Extending this to a multistage setup with more than two models could further improve efficiency but also add complexity. Second, our study is limited to machine translation. QE-based deferral works particularly well in MT due to the availability of high-quality human-labeled data for training QE models. Extending this approach to other tasks where such data is scarce is not straightforward. Third, our method assumes the smaller model is reasonably competitive with the larger one, which is a fair assumption for MT, as shown in our experiments. If the gap in win rates is too large, cascading offers little benefit, as most examples would require deferral. Finally, from the standpoint of deploying in a real-world application, any cascading solution may introduce latency if not engineered properly. Our analysis about the computational efficiency of our method in §3 does not account for that.

Acknowledgments

We thank Duarte Alves, Patrick Fernandes, Emmanouil Zaranis, and the SARDINE lab team for helpful discussions. This work was supported by EU's Horizon Europe Research and Innovation Actions (UTTER, contract 101070631), by the project DECOLLAGE (ERC-2022-CoG 101088763), by the Portuguese Recovery and Resilience Plan through project C645008882-00000055 (Center for Responsible AI), and by FCT/MECI through national funds and when applicable co-funded EU funds under UID/50008: Instituto de Telecomunicações.

References

Sweta Agrawal, António Farinhas, Ricardo Rei, and Andre Martins. 2024. Can automatic metrics assess high-quality translations? In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 14491–14502, Miami, Florida, USA. Association for Computational Linguistics.

AI@Meta. 2024. Llama 3 model card.

Duarte Miguel Alves, José Pombal, Nuno M Guerreiro, Pedro Henrique Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, José G. C. de Souza, and Andre Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. In *First Conference on Language Modeling*.

Frederic Blain, Chrysoula Zerva, Ricardo Rei, Nuno M. Guerreiro, Diptesh Kanojia, José G. C. de Souza, Beatriz Silva, Tânia Vaz, Yan Jingxuan, Fatemeh Azadi, Constantin Orasan, and André Martins. 2023. Findings of the WMT 2023 shared task on quality estimation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 629–653, Singapore. Association for Computational Linguistics.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, and 12 others. 2020. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Hannah Béchara, Constantin Orăsan, Carla Parra Escartín, Marcos Zampieri, and William Lowe. 2021. The role of machine translation quality estimation in the post-editing workflow. *Informatics*, 8(3).

Sheila Castilho and Sharon O'Brien. 2017. Acceptability of machine-translated content: A multi-language evaluation by translators and end-users. *Linguistica Antverpiensia, New Series–Themes in Translation Studies*, 16.

Charlie Chen, Sebastian Borgeaud, Geoffrey Irving, Jean-Baptiste Lespiau, Laurent Sifre, and John Jumper. 2023a. Accelerating large language model decoding with speculative sampling. *Preprint*, arXiv:2302.01318.

Lingjiao Chen, Matei Zaharia, and James Zou. 2023b. Frugalgpt: How to use large language models while reducing cost and improving performance. *Preprint*, arXiv:2305.05176.

Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, Parker Schuh, Kensen Shi, Sasha Tsvyashchenko, Joshua Maynez, Abhishek Rao, Parker Barnes, Yi Tay, Noam Shazeer, Vinodkumar Prabhakaran, and 48 others. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.

Dujian Ding, Ankur Mallick, Chi Wang, Robert Sim, Subhabrata Mukherjee, Victor Rühle, Laks V. S. Lakshmanan, and Ahmed Hassan Awadallah. 2024. Hybrid LLM: Cost-efficient and quality-aware query routing. In *The Twelfth International Conference on Learning Representations*.

António Farinhas, José de Souza, and Andre Martins. 2023. An empirical study of translation hypothesis

- ensembling with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 11956–11970, Singapore. Association for Computational Linguistics.
- António Farinhas, Haau-Sing Li, and Andre Martins. 2024. Reranking laws for language generation: A communication-theoretic perspective. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Christian Federmann. 2018. Appraise evaluation framework for machine translation. In *Proceedings of the 27th International Conference on Computational Linguistics: System Demonstrations*, pages 86–88, Santa Fe, New Mexico. Association for Computational Linguistics.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.
- Marina Fomicheva, Shuo Sun, Lisa Yankovskaya, Frédéric Blain, Francisco Guzmán, Mark Fishel, Nikolaos Aletras, Vishrav Chaudhary, and Lucia Specia. 2020. Unsupervised quality estimation for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:539–555.
- Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation. In *Findings of the Association for Computational Linguistics: EMNLP* 2023, pages 9198–9209, Singapore. Association for Computational Linguistics.
- Phillip Good. 2000. Permutation tests: a practical guide to resampling methods for testing hypotheses. Springer New York, NY.
- Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.
- Neha Gupta, Harikrishna Narasimhan, Wittawat Jitkrittum, Ankit Singh Rawat, Aditya Krishna Menon, and Sanjiv Kumar. 2024. Language model cascades: Token-level uncertainty and beyond. In *The Twelfth International Conference on Learning Representations*.
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego

- de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b. *Preprint*, arXiv:2310.06825.
- Wittawat Jitkrittum, Neha Gupta, Aditya Krishna Menon, Harikrishna Narasimhan, Ankit Singh Rawat, and Sanjiv Kumar. 2023. When does confidence-based cascade deferral suffice? In *Thirty-seventh Conference on Neural Information Processing Systems*.
- Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. MetricX-24: The Google submission to the WMT 2024 metrics shared task.
 In Proceedings of the Ninth Conference on Machine Translation, pages 492–504, Miami, Florida, USA. Association for Computational Linguistics.
- Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings* of the Eighth Conference on Machine Translation, pages 756–767, Singapore. Association for Computational Linguistics.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *Preprint*, arXiv:2001.08361.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, and 3 others. 2024a. Findings of the WMT24 general machine translation shared task: The LLM era is here but MT is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1–46, Miami, Florida, USA. Association for Computational Linguistics.
- Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. Findings of the 2022 conference on machine translation (WMT22). In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Tom Kocmi and Ondřej Bojar. 2017. Curriculum learning and minibatch bucketing in neural machine translation. In *Proceedings of the International Conference Recent Advances in Natural Language Process*

- ing, RANLP 2017, pages 379–386, Varna, Bulgaria. INCOMA Ltd.
- Tom Kocmi, Vilém Zouhar, Christian Federmann, and Matt Post. 2024b. Navigating the metrics maze: Reconciling score magnitudes and accuracies. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1999–2014, Bangkok, Thailand. Association for Computational Linguistics.
- Geza Kovacs, Daniel Deutsch, and Markus Freitag. 2024. Mitigating metric bias in minimum Bayes risk decoding. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1063–1094, Miami, Florida, USA. Association for Computational Linguistics.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. 2023. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles*.
- Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 19274–19286. PMLR.
- Jonathan Mamou, Oren Pereg, Moshe Wasserblat, and Roy Schwartz. 2022. Tangobert: Reducing inference cost by using cascaded architecture. *Preprint*, arXiv:2204.06271.
- Pedro Henrique Martins, Patrick Fernandes, João Alves, Nuno M Guerreiro, Ricardo Rei, Duarte M Alves, José Pombal, Amin Farajian, Manuel Faysse, Mateusz Klimaszewski, and 1 others. 2024. Eurollm: Multilingual language models for europe. *arXiv* preprint arXiv:2409.16235.
- Harikrishna Narasimhan, Wittawat Jitkrittum, Ankit Singh Rawat, Seungyeon Kim, Neha Gupta, Aditya Krishna Menon, and Sanjiv Kumar. 2025. Faster cascades via speculative decoding. In *The Thirteenth International Conference on Learning Representations*.
- Isaac Ong, Amjad Almahairi, Vincent Wu, Wei-Lin Chiang, Tianhao Wu, Joseph E. Gonzalez, M Waleed Kadous, and Ion Stoica. 2024. Routellm: Learning to route llms with preference data. *Preprint*, arXiv:2406.18665.
- Guillem Ramírez, Alexandra Birch, and Ivan Titov. 2024. Optimising calls to large language models with uncertainty-based two-tier selection. *Preprint*, arXiv:2405.02134.
- Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins.

- 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 841–848, Singapore. Association for Computational Linguistics.
- Ricardo Rei, Jose Pombal, Nuno M. Guerreiro, João Alves, Pedro Henrique Martins, Patrick Fernandes, Helena Wu, Tania Vaz, Duarte Alves, Amin Farajian, Sweta Agrawal, Antonio Farinhas, José G. C. De Souza, and André Martins. 2024. Tower v2: Unbabel-IST 2024 submission for the general MT shared task. In *Proceedings of the Ninth Conference on Machine Translation*, pages 185–204, Miami, Florida, USA. Association for Computational Linguistics.
- Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 634–645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Nikhil Sardana, Jacob Portes, Sasha Doubov, and Jonathan Frankle. 2024. Beyond chinchilla-optimal: Accounting for inference in language model scaling laws. In *Forty-first International Conference on Machine Learning*.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. 2024. Scaling Ilm test-time compute optimally can be more effective than scaling model parameters. *Preprint*, arXiv:2408.03314.
- Mitchell Stern, Noam Shazeer, and Jakob Uszkoreit. 2018. Blockwise parallel decoding for deep autoregressive models. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, NIPS'18, page 10107–10116, Red Hook, NY, USA. Curran Associates Inc.
- Neeraj Varshney and Chitta Baral. 2022. Model cascading: Towards jointly improving efficiency and accuracy of NLP systems. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 11007–11021, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Yu Wan, Baosong Yang, Derek Fai Wong, Lidia Sam Chao, Liang Yao, Haibo Zhang, and Boxing Chen.

2022. Challenges of neural machine translation for short texts. *Computational Linguistics*, 48(2):321–342.

Longyue Wang, Chenyang Lyu, Tianbo Ji, Zhirui Zhang, Dian Yu, Shuming Shi, and Zhaopeng Tu. 2023. Document-level machine translation with large language models. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 16646–16661, Singapore. Association for Computational Linguistics.

Heming Xia, Zhe Yang, Qingxiu Dong, Peiyi Wang, Yongqi Li, Tao Ge, Tianyu Liu, Wenjie Li, and Zhifang Sui. 2024. Unlocking efficiency in large language model inference: A comprehensive survey of speculative decoding. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7655–7671, Bangkok, Thailand. Association for Computational Linguistics.

Chrysoula Zerva, Frederic Blain, José G. C. De Souza, Diptesh Kanojia, Sourabh Deoghare, Nuno M. Guerreiro, Giuseppe Attanasio, Ricardo Rei, Constantin Orasan, Matteo Negri, Marco Turchi, Rajen Chatterjee, Pushpak Bhattacharyya, Markus Freitag, and André Martins. 2024. Findings of the quality estimation shared task at WMT 2024: Are LLMs closing the gap in QE? In *Proceedings of the Ninth Conference on Machine Translation*, pages 82–109, Miami, Florida, USA. Association for Computational Linguistics.

Chrysoula Zerva, Frédéric Blain, Ricardo Rei, Piyawat Lertvittayakumjorn, José G. C. de Souza, Steffen Eger, Diptesh Kanojia, Duarte Alves, Constantin Orăsan, Marina Fomicheva, André F. T. Martins, and Lucia Specia. 2022. Findings of the WMT 2022 shared task on quality estimation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 69–99, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Ran Zmigrod, Tim Vieira, and Ryan Cotterell. 2022. Exact paired-permutation testing for structured test statistics. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4894–4902, Seattle, United States. Association for Computational Linguistics.

Marija Šakota, Maxime Peyrard, and Robert West. 2024. Fly-swat or cannon? cost-effective language model choice via meta-modeling. In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining*, WSDM '24, page 606–615, New York, NY, USA. Association for Computing Machinery.

A Automatic Evaluation

Here, we include versions of Figs. 2 to 4 but using COMET (Rei et al., 2022a) scores instead of METRICX (Juraska et al., 2023). Please check

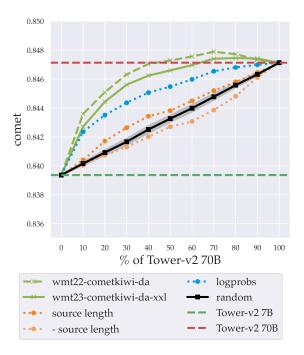


Figure 9: Translation quality of cascading combining Tower-v2 7B and Tower-v2 70B according to COMET, as the inference computation budget varies.

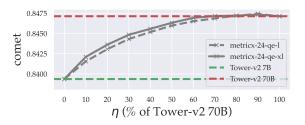


Figure 10: Translation quality of cascaded systems with deferral based on metricx-24-hybrid-large-v2p6 and metricx-24-hybrid-xl-v2p6 according to COMET.

Figs. 9 to 11. Crucially, the main conclusions do not change.

B Human Evaluation

B.1 Annotation guidelines

We share below the annotation guidelines shared with the freelancers.

Task overview. This task involves evaluating two alternative translations of a source text and assigning a rating to each translation based on its overall quality and adherence to the source content. You should consider accuracy, fluency, and overall quality when assessing the different translations.

Annotation scale. Each translation should be evaluated on a continuous scale from 0 to 6 with

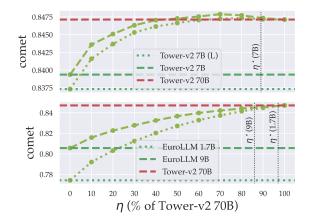


Figure 11: Translation quality of cascaded systems with deferral based on wmt22-cometkiwi-da, according to COMET. Large model: Tower-v2 70B. Small models: Tower-v2 7B (L), Tower-v2 7B (top); EuroLLM 1.7B, EuroLLM 9B (bottom).

the quality levels described below:

- 6 (perfect meaning and grammar): The meaning of the translation is completely consistent with the source and the surrounding context, if applicable. The grammar is also correct.
- 4 (most meaning preserved and few grammar mistakes): The translation retains most of the meaning of the source. It may have some grammar mistakes or minor contextual inconsistencies.
- 2 (some meaning preserved): The translation preserves some of the meaning of the source but misses significant parts. The narrative is hard to follow due to fundamental errors. Grammar may be poor.
- **0** (nonsense/no meaning preserved): Nearly all information is lost between the translation and source. Grammar is irrelevant.

Annotation interface. Figs. 12 and 13 show the annotation interface. If two candidates were the same or of the same quality, the annotators were asked to use "**match sliders**" to give them the exact same score. And, they could also use the absolute scale range to show preference between the translations.

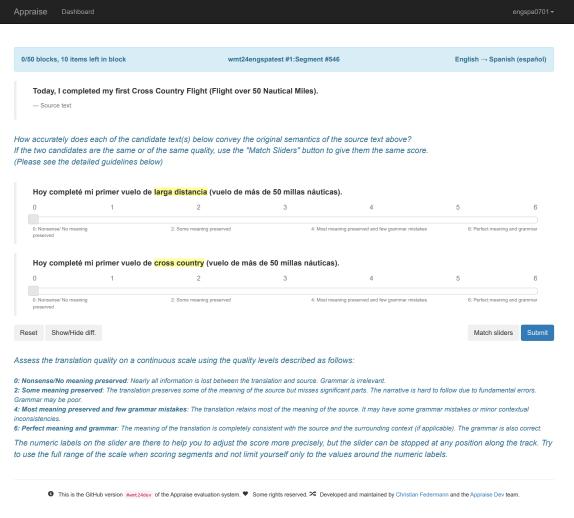


Figure 12: Annotation interface for en-es.

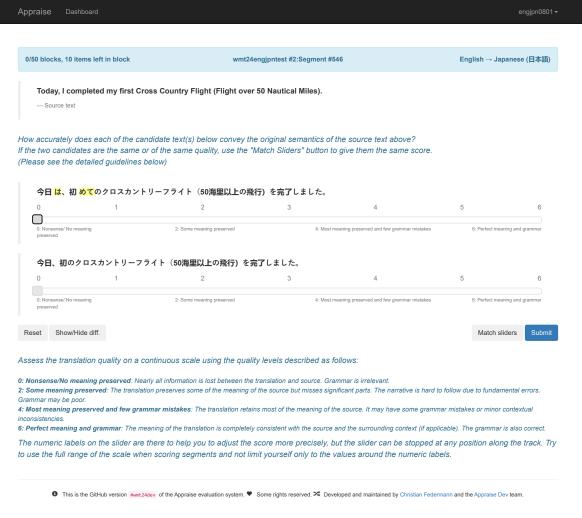


Figure 13: Annotation interface for en-ja.