

# Improving Minimum Bayes Risk Decoding with Multi-Prompt

Anonymous ACL submission

## Abstract

While instruction fine-tuned LLMs are effective text generators, sensitivity to prompt construction makes performance unstable and sub-optimal in practice. Relying on a single ‘best’ prompt cannot capture all differing approaches to a generation problem. Using this observation, we propose *multi-prompt decoding*, where many candidate generations are decoded from a prompt bank at inference-time. To ensemble candidates, we use Minimum Bayes Risk (MBR) decoding, which selects a final output using a trained value metric. We show multi-prompt improves MBR across a comprehensive set of conditional generation tasks (Figure 1), and show this is a result of estimating a more diverse and higher quality candidate space than that of a single prompt. Our experiments confirm multi-prompt improves generation across tasks, models and metrics.<sup>1</sup>

## 1 Introduction

Minimum Bayes Risk (MBR) decoding (Bickel and Doksum, 1977) has been shown to improve generation quality of large language models (LLMs) compared to typical single-output decoding methods, such as beam search and sampling, across NLP tasks (Shi et al., 2022; Suzgun et al., 2023). A special case of MBR, self-consistency (Wang et al., 2023), has been widely-used to improve LLM reasoning capabilities by ensembling reasoning paths. MBR leverages a set of candidates and selects the one with the highest expected utility, using all other hypotheses as references (see Fig. 2, left), following a simple intuition that a desirable output should be highly probable and consistent with others.

A central question to improve MBR is how to balance between diversity and adequacy within the candidate set. Prior work has found success using sampling-based decoding to generate hypotheses

<sup>1</sup>Our experiment code, data and prompts are available at [https://anonymized\\_url](https://anonymized_url)

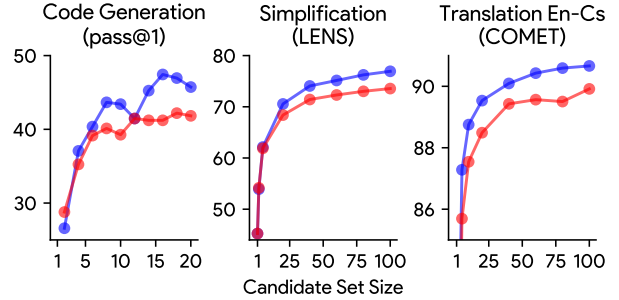


Figure 1: Multi-prompt and single prompt MBR results for code generation on HUMANEVAL, text simplification on SIMPEVAL, and translation on WMT ‘22 EN-Cs generated with open-source 7B LLMs (details in §4).

from a given input (Eikema and Aziz, 2020; Freitag et al., 2022a, 2023). However, naively increasing the sampling temperature eventually degrades the quality of the candidates. Recently, instruction fine-tuned LLMs (Ouyang et al., 2022; Chung et al., 2022) has opened up the possibility of writing the “prompts” in various formats to elicit more diverse and high quality outputs, as these models are observed to be sensitive to prompt design, where a slight change in phrasing or the inclusion of more relevant example can significantly impact model outputs (Srivastava et al., 2023; White et al., 2023).

Taking advantage of the prompt sensitivity of LLMs, we introduce multi-prompt MBR decoding, which samples candidates using a bank of human- or model-written prompts (see Figure 2, right). Intuitively, exploring a variety of prompts enables the generation of diverse, high quality hypotheses that provide a closer representation of the true output distribution. By guiding the model towards different modes or regions of the output space, each prompt captures unique sequences that are coherent and relevant to the input.

We experiment with three distinct generation tasks: text simplification (Maddela et al., 2023), machine translation (Kocmi et al., 2022), and code generation (Chen et al., 2021). Each task assess the impact of different prompt components on multi-

prompt MBR, such as instance-level prompts for code, task descriptions for simplification, and in-context examples for translation. To account for the relative quality between prompts, we develop different strategies for selecting prompts that significantly improve over random choice. These strategies include *sampling* prompts from a large prompt bank based on their usage on a training set and *selecting* prompts using embedding-based heuristics when a training set is unavailable.

We evaluate multi-prompt MBR on a broad range of LLMs including both open-source models like Llama 2 (Touvron et al., 2023) and state-of-the-art closed-source models such as GPT-4 (Achiam et al., 2023). The results show that multi-prompt MBR consistently improves single-prompt MBR across all three tasks and model scales, with gains of up to 14% on HumanEval (Chen et al., 2021) and 8 points of LENS on SIMPEVAL (Maddela et al., 2023). Figure 1 displays the results for models at the 7B scale. Additionally, we study the dynamics between different utility and evaluation metrics, revealing that multi-prompt MBR with one metric improves performance universally across metrics.

## 2 Preliminaries

Instruction fine-tuned LLMs are trained to follow arbitrary natural language task descriptions (Wei et al., 2022). Given an input  $x$  and prompt  $\rho$ , an autoregressive language model  $\pi_\theta$  parameterized by  $\theta$  estimates an output sequence  $y \sim \pi_\theta(x, \rho)$  using an decoding algorithm by sampling the next token conditioned on the input  $\pi_\theta(y_i | y_{<i}, x, \rho)$ . The decoding algorithm aims to generate  $y$  by maximizing the sequence likelihood over the language model distribution  $\pi_\theta(y|x, \rho) = \prod_{i=1}^T \pi_\theta(y_i | y_{<i}, x, \rho)$ .

**Minimum Bayes Risk Decoding.** As often observed in practice (Freitag et al., 2022a), unfortunately, the highest likelihood generation is not necessarily the highest quality (Jaeger and Levy, 2006). Building on this observation, MBR decoding (Bickel and Doksum, 1977; Eikema and Aziz, 2020) first samples a set of hypotheses  $\mathcal{H}$  from the model  $\pi_\theta$ , approximating the true distribution of output space  $\mathcal{Y}$ , then selects the output  $\hat{y}_{MBR}$  that maximizes the expected utility (or minimizes the expected loss in traditional formulation) with respect to a set of references  $\mathcal{R}$ :

$$\hat{y}_{MBR} = \arg \max_{y \in \mathcal{H}} (\mathbb{E}_{\mathcal{H} \sim \pi_\theta} [U(y, \mathcal{R})]), \quad (1)$$

where  $U(y, \mathcal{R}) = \mathbb{E}_{y' \sim \mathcal{R}} [u(y, y')]$  and  $u(y, y')$  is a

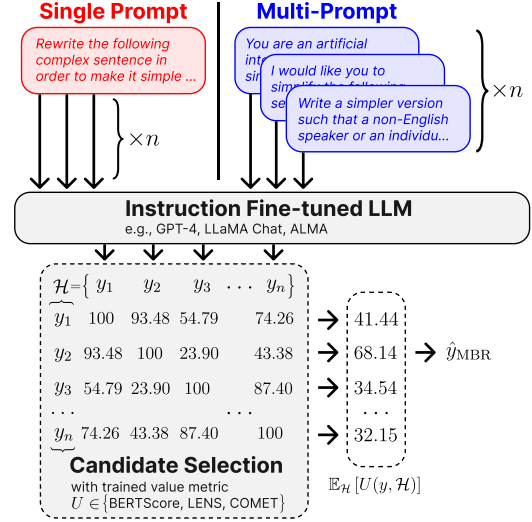


Figure 2: Multi-prompt MBR generates candidates using a human- or model-written prompt bank and selects the highest pairwise score with a trained value metric.

utility function that evaluates hypothesis  $y$  against a reference  $y'$ . In practice,  $\mathcal{R}$  is also sampled from the same model  $\pi_\theta$  under the assumption that the model produces reliable outputs in expectation, and is usually set as identical to hypothesis set  $\mathcal{H}$ .

Bertsch et al. (2023) show that some successful techniques that improve LLMs' performance such as self-consistency (Wang et al., 2023) and output ensemble (Kobayashi, 2018) are special cases of MBR. For example, self-consistency, which takes the majority vote among answers extracted from multiple sampled reasoning chains, can be viewed as MBR with utility function as  $u(y, y') = \mathbb{1}[\text{ans}(y) = \text{ans}(y')]$ , where  $\text{ans}(y)$  is the answer extracted from the reasoning path  $y$ .

## 3 Multi-Prompt MBR Decoding

Prior work on MBR decoding explores models trained for specific tasks, where the hypothesis set is generated given a single input  $x$  (Freitag et al., 2022a; Fernandes et al., 2022). With instruction fine-tuned LLMs, the input  $x$  is contained within a structured prompt  $\rho$ , consisting of task instruction and/or in-context examples. Earlier studies have extensively documented that the design of the prompt has a dramatic impact on overall performance (Mishra et al., 2022; Khashabi et al., 2022; Lu et al., 2022; Sclar et al., 2023).

To investigate these phenomena, we show in Figure 3a (bottom) the likelihoods and quality of samples from 10 prompts of varying performance for a text simplification task, measuring quality

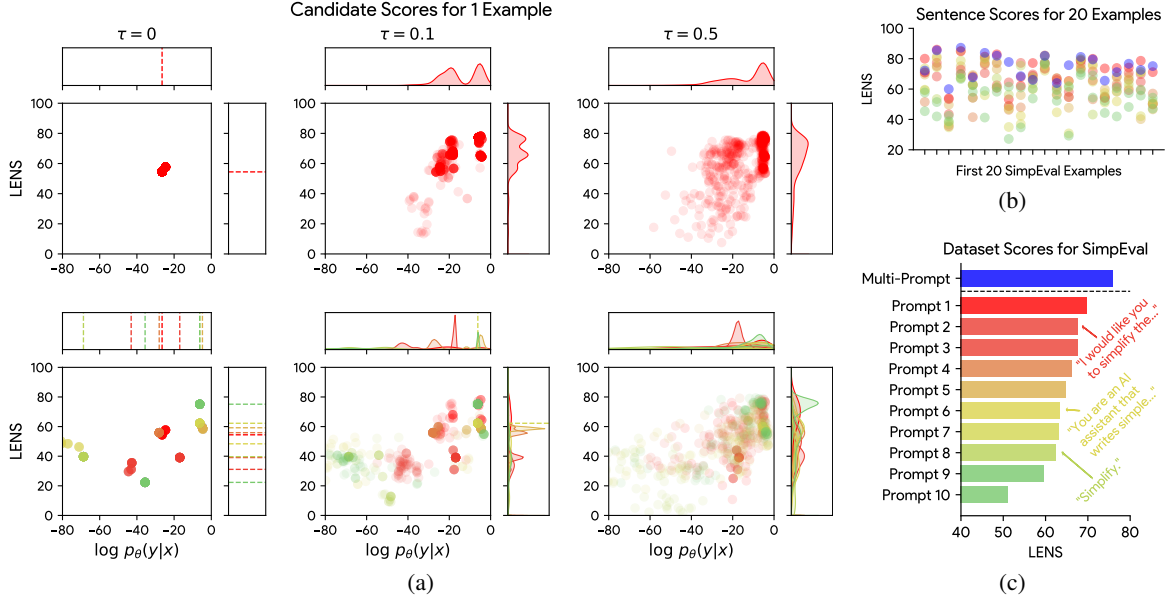


Figure 3: (a) LENS score and sequence probability for 1000 generations on a single text simplification example decoded from Llama 2 7B Chat with temperatures  $\tau = [0, 0.1, 0.5]$  using a single prompt (top) and multiple prompts (bottom). As the temperature increases, we find each prompt estimates candidate sequences centered at different modes. (b) LENS scores of the best generation per-prompt for the first 20 sentences in SIMPEVAL, showing no single prompt produces the best overall output. (c) Dataset-level LENS performance of each prompt when performing single prompt MBR vs. multi-prompt MBR.

as the LENS metric score against a set of gold references. Greedy sampling ( $\tau = 0$ ) estimates different sequences for each instruction, with single prompt (Figure 3a, top) generating a single sequence. As we increase temperature  $\tau$ , generations from a single prompt simply exhibit noise centered around the mode of the highest likelihood sequence, while multi-prompt estimates a generations around modes uniquely defined by each prompt. For instance, one of the prompts (i.e., Prompt 9 highlighted in green) produces the highest quality generation for this one input sentence, despite having a low performance over the entire dataset. In fact, no prompt consistently produces the highest quality sequences, as illustrated in Figure 3b, rather prompts are most effective at different inputs.

Building upon these insights, we propose Multi-Prompt MBR decoding, depicted in Figure 2, where the MBR hypothesis set  $\mathcal{H}$  consists of outputs sampled from  $n$  distinct prompts  $\rho$ :

$$\mathcal{H} = \bigcup_{i=1}^n \mathcal{H}_i, \text{ where } \mathcal{H}_i = \{y|y \sim \pi_\theta(x, \rho_i)\}. \quad (2)$$

Bertsch et al. (2023) show that MBR seeks the mode of some distribution  $q$  over a quality feature  $\phi(y)$  applied to the output space rather than the mode of the model’s distribution:

$$\hat{y}_{\text{MBR}} \approx \arg \max_{y \in \mathcal{H}} q(\phi(y)|x). \quad (3)$$

We hypothesize, in expectation, the mode of  $\phi(y)$  across outputs from multiple prompts has higher downstream performance compared to that derived from a single prompt. This is empirically supported by our example, where Figure 3c shows that multi-prompt MBR outperforms individual single-prompt MBR across the full task dataset.

Although multi-prompt ensembles hypothesis spaces between prompts, some notion of objective quality still exists when constructing the prompt bank. As shown in Figure 3c, the majority of the 10 human-written prompts fall within a 10-point range of LENS scores when evaluated on the task dataset but a few prompts consistently produce low-quality generation. Therefore, to account for the hierarchy in prompt quality, we propose two methods for choosing the prompts used at generation time from a prompt bank  $\mathcal{P}$ : sampling from a learned distribution of prompts, based on a small unlabeled train set (§3.1); and selecting a subset of prompts based on heuristics in the absence of a train set (§3.2).

### 3.1 Prompt Sampling

In this approach, we first calculate the probability of each prompt  $p(\rho)$  as the proportion of times that prompt generates the highest scoring output on a separate training set. At inference time, prompts are sampled with replacements from this learned probability distribution, and candidate outputs are

then generated given these prompts.

**Top- $p$  Prompt Sampling.** Inspired by the principle of nucleus sampling (Holtzman et al., 2020), our goal is to keep the prompts with high probability and truncate the least used prompts by setting their probabilities to zero. We define the top- $p$  prompt set as the minimal set  $\mathcal{P}_{\text{top-}p} \subseteq \mathcal{P}$  such that:

$$\sum_{i=0}^{|\mathcal{P}_{\text{top-}p}|} p(\rho_i) \geq p. \quad (4)$$

We then re-normalize the distribution of  $\mathcal{P}_{\text{top-}p}$  and sample prompts from the new distribution:

$$p'(\rho) = \begin{cases} \frac{p(\rho)}{\sum_{\rho \in \mathcal{P}_{\text{top-}p}} p(\rho)} & \text{if } \rho \in \mathcal{P}_{\text{top-}p} \\ 0 & \text{otherwise.} \end{cases} \quad (5)$$

### 3.2 Prompt Selection

Prompt selection chooses a fixed subset  $\mathcal{P}_{\text{best}} \subset \mathcal{P}$  of  $|\mathcal{P}_{\text{best}}| = k$  prompts based on heuristics. Compared to sampling, this does not require an additional training set to evaluate prompt efficacy. We consider the following heuristics for selecting  $\mathcal{P}_{\text{best}}$ : prompts that have the closest similarity and greatest dissimilarity with others, and prompts that are randomly selected from each  $k$ -NN cluster, which is also useful when a training set is presented, allowing the selection of high-performing prompts within each cluster. In our experiments, we calculate the semantic (dis)similarity of prompts based on their SentenceBERT (Reimers and Gurevych, 2019) embeddings.

## 4 Experiment Setup

In this section, we describe the experimental details for evaluating the efficacy of multi-prompt MBR decoding across tasks, prompt setups, models, and utility metrics, with results and analyses in §5.

### 4.1 Tasks & Datasets

Unlike previous work applying MBR to a single generation task (Shi et al., 2022; Eikema and Aziz, 2022), we deliberately select three unique tasks to demonstrate the universality of multi-prompt: text simplification with task-level instructions, code generation with example-level instructions, and machine translation with in-context examples.

**Code Generation.** We use HumanEval (Chen et al., 2021) benchmark, where models are tasked with generating a Python program given a description with unit tests. Since each example is a unique

coding task, we generate a unique prompt bank for each input. Following Zhang et al. (2023), we reject empty, degenerate (e.g., pass, return None), or non-compiling programs before applying MBR.

**Text Simplification.** We use the SIMPEVAL<sub>2022</sub> test set (Maddela et al., 2023), containing complex sentences from Wikipedia, paired with human-written simplifications. The prompt bank is generated based on author-written examples (Table 4) and are used for the entire dataset.

**Machine Translation.** We purposely choose the EN  $\rightarrow$  CS language pair from the WMT 22 (Kocmi et al., 2022) newstest corpus, ensuring its exclusion from the training data of recent translation LLMs or metrics (Xu et al., 2024). Results on additional language pairs are in Appendix C.2.

### 4.2 Constructing the Prompt Bank

Following existing work studying prompt sensitivity (Mizrahi et al., 2023; Gonen et al., 2023), our experiments rely on a small set of manually written seed prompts, and use an LLM to generate diverse paraphrases of prompts. Model-written prompts are generated using GPT-4 Turbo. For seed prompts, the authors manually write 10 for text simplification (Table 4) and use the original HUMAN-EVAL instruction from each example as the seed prompt for code generation. The translation prompts consist of randomly sampled in-context examples from previous WMT shared tasks.

For experiments, we select from the prompt bank with top- $p$  prompt sampling (§5.2) using  $p = 0.6$ , where the prompt usage  $p(\rho)$  is calculated using a held-out 20% split of each dataset. Human-written prompts and prompt generation instructions are included in Appendix A.

### 4.3 Models

Our main experiments are performed with Llama 2-7B Chat (Touvron et al., 2023) for simplification, ALMA-7B-R (Xu et al., 2024) for translation and CodeLLaMA-13B Instruct (Roziere et al., 2023) for code generation, all fine-tuned to follow instructions. In §5.3 we further explore a wide range of model architectures and sizes, including state-of-the-art and task-specific fine-tuned models. Unless otherwise specified, we generate the hypothesis set using nucleus sampling (Holtzman et al., 2020) with  $\tau = 0.9$ ,  $p = 0.95$ . We include a detailed review of all models in this work in Appendix B.2.



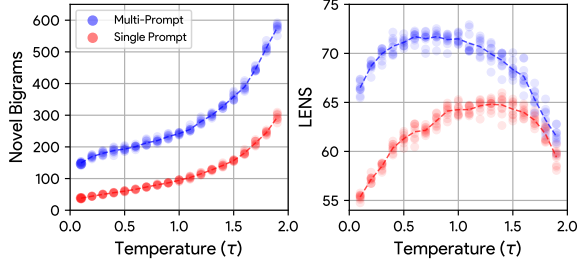


Figure 4: Candidate set diversity and LENS scores across temperatures for simplification task. At low temperatures, the increased candidate diversity from multi-prompt directly translates to improved performance.

#### 4.4 Utility Metrics & Evaluation

Our core experiments use the trained LENS (Madela et al., 2023) for simplification and COMET (Rei et al., 2020) for translation as the candidate selection metric. For code generation, we use MBR-EXEC (Shi et al., 2022), which executes each candidate program against a set of test cases, selecting the program with the highest agreement over all test cases’ outputs. As in Zhang et al. (2023), we use the docstring examples as test cases for MBR-EXEC and evaluate with pass@1. Given the growing body of work on metric development, we verify our multi-prompt results across a broad range of utility and evaluation metrics in §5.4.

### 5 Experiment Results

We compare multi-prompt decoding to traditional MBR (§5.1), ablate the prompt sampling mechanism (§5.2), vary model architectures (§5.3), evaluate across utility metrics (§5.4) and finally evaluate multi-prompt on efficient MBR alternatives (§5.5).

#### 5.1 How does multi-prompt MBR perform?

**Multi-prompt Improves MBR.** We report our main results in Figure 1 and Table 2, comparing single prompt and multi-prompt performance as the number of generated candidates increases, with detailed results in Figure 7 in Appendix. Multi-prompt MBR consistently outperforms traditional MBR for all tasks.

**Candidate Diversity  $\Rightarrow$  Quality.** To measure the impact of temperature on the candidate set quality, we report performance and diversity, as measured by novel bi-grams, across temperatures in Figure 4. For low temperatures, we find that multi-prompt generates a consistently more diverse candidate space, which directly translates to higher-quality generation. While single prompt MBR performance improves with temperature  $\tau > 1$ , despite

	pass@1	LENS	COMET
<i>Single Prompt</i> ( $ \mathcal{H}  = 100$ )	48.78	69.45	90.14
<i>Multi-Prompt + Prompt Sampling</i> ( $ \mathcal{P}  = 100$ )			
Random Selection	–	74.91	89.98
Prompt Sampling	–	78.29	90.33
Top- $p$ Prompt Random	–	78.61	90.11
Top- $p$ Prompt Sampling	–	<b>79.08</b>	<b>90.36</b>
<i>Single Prompt</i> ( $ \mathcal{H}  = 10$ )	41.55	51.64	87.54
<i>Multi-Prompt + Prompt Selection</i> ( $\mathcal{P}_{\text{best}} \subset \mathcal{P}$ , $ \mathcal{P}_{\text{best}}  = 10$ )			
Random Selection	39.63	60.00	87.81
$k$ -NN Cluster Random	40.24	58.73	87.80
Farthest Similarity	<b>44.51</b>	58.32	<b>88.14</b>
Closest Similarity	37.80	61.53	87.73
Highest Performance	–	62.43	87.65
$k$ -NN Cluster Performance	–	<b>66.12</b>	87.73

Table 1: Results for prompt sampling using 100 prompts (top) and subset selection using 10 of 100 prompts (bottom). Sampling from a weighted, truncated distribution improves multi-prompt across candidate set sizes.

generating an equal or greater diversity set than multi-prompt, multi-prompt MBR still produces higher quality candidates. As  $\tau \rightarrow 2$ , the quality of single and multi-prompt MBR begins to degrade as their candidate sets become too noisy to generate high-quality sequences. Framing the decoding process as each prompt estimating a unique distribution of candidate generations (§3), the ability of multi-prompt to achieve higher quality generation as a result of candidate set diversity is intuitively the byproduct of combining multiple candidate distributions defined by each instruction.

#### 5.2 What is the impact of the prompt bank?

##### Sampling Prompts Improves Candidate Quality.

Table 1 (top) reports results for multi-prompt across different prompt sampling methods for text simplification and translation. Note that, code generation is excluded as a unique set of prompts is generated for each HumanEval example, rather than the same prompts used across the entire dataset. We find sampling prompts by usage and truncating the top- $p$  prompts improves multi-prompt over a random selection baseline, with top- $p$  prompt sampling performing the best on both tasks.

##### Multi-prompt is Sensitive to the Prompt Bank.

Table 1 (bottom) reports results for different prompt subset selection methods, which use heuristics to select a smaller set of prompts for multi-prompt to maximize performance. This includes the 10 closest and furthest prompt embeddings, the 10 highest performing prompts, and a  $k$ -NN cluster of prompt embeddings where a single prompt is selected from

	Single Prompt	Multi-prompt
<i>Text Simplification</i> ( $n = 100$ ) – SIMPEVAL (LENS)		
Ctrl T5 3B	72.6	—
Ctrl T5 11B	74.4	—
GPT-3.5	75.37	80.09 (+4.72)
GPT-4	73.27	80.60 (+7.33)
LLaMA 2 7B Chat	70.51	76.29 (+5.78)
LLaMA 2 13B Chat	71.29	77.93 (+6.64)
LLaMA 2 70B Chat	75.09	80.53 (+5.44)
<i>Translation</i> ( $n = 100$ ) – WMT '22 EN-Cs (COMET)		
WMT '22 Winners	91.9	—
MS Translate API	90.6	—
GPT-3.5	91.89	92.39 (+0.50)
GPT-4	91.57	91.92 (+0.35)
ALMA 7B R	90.14	90.36 (+0.22)
ALMA 13B R	90.56	90.97 (+0.41)
<i>Code Generation</i> ( $n = 20$ ) – HUMANEval (pass@1)		
StarCoder 2 15B	44.51	45.73 (+1.22)
GPT-3.5	66.46	73.17 (+6.71)
GPT-4	71.34	85.36 (+14.0)
CodeLLaMA 7B	35.97	39.68 (+3.71)
CodeLLaMA 13B	43.29	48.17 (+4.88)
CodeLLaMA 34B	47.56	53.65 (+6.09)
CodeLLaMA 70B	60.97	68.29 (+7.32)

Table 2: Metric scores for state-of-the-art systems compared to multi-prompt LLMs using  $n$  candidates. Translation and simplification baselines are as reported in Hendy et al. (2023) and Maddela et al. (2023).

each cluster. Each selection method had a significant impact on performance when compared to a random selection of 10 prompts (+0.03 pass@1, +14 LENS and +0.6 COMET). For text simplification, decoding with the 10 highest performing prompts is further improved by selecting prompts from a  $k$ -NN clustering of prompt embeddings, which enforces a dis-similarity between prompts. Translation does not benefit from clustering, and instead both translation and code generation benefit from simply generating with farthest similarity, or semantically distant prompts. These results highlight multi-prompt’s sensitivity to the prompt construction, and shows that enforcing both diversity via multi-prompt and performance via prompt selection improves candidate generation.

### 5.3 Does multi-prompt MBR apply across model scale and architecture?

**Increasing Returns as Models Scale.** To argue multi-prompt improves generation across instruction fine-tuned models and at scale, we experiment with widely used LLMs. Figure 5 reports improvement of multi-prompt over single prompt as a  $\Delta$  change in score, with analysis of per-model results in Appendix C.3. On text simplification, instruction fine-tuned models appear to converge to a +5

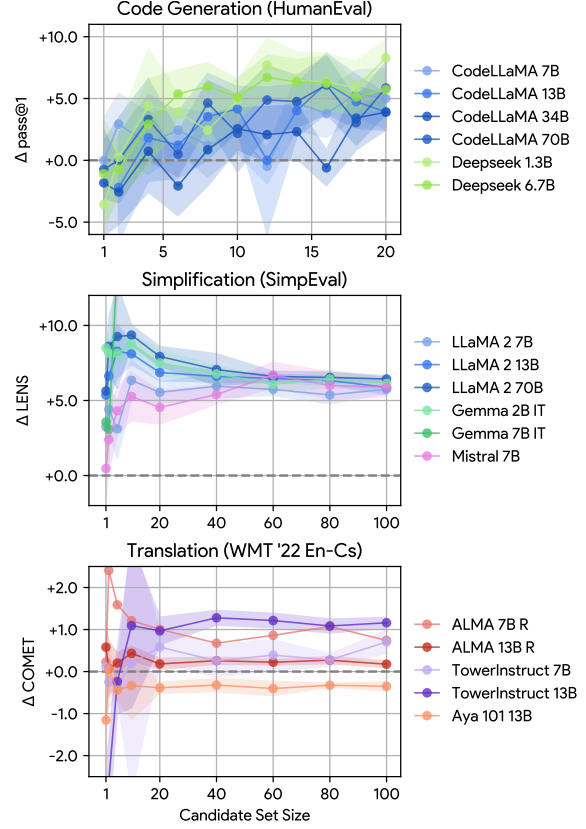


Figure 5:  $\Delta$  metric improvement from single prompt to multi-prompt across model sizes and architectures, reported with a 95% CI bootstrapped over 5 iterations. For absolute performance, see Figure 9.

improvement in LENS score as candidate set size increases, consistent across model sizes and types, while code generation models saw increasing returns using multi-prompt as candidate set size increased. We find the same trend of convergence to a score improvement for translation, but saw inconsistent results, which may be a result of the vast difference in training data for translation LLMs.

**LLMs with Multi-prompt Outperform Fine-tuned Models.** Whether general instruction fine-tuned LLMs can outperform an LLM trained or fine-tuned on a conditional generation task is still an active question (Chung et al., 2022), so we compare state-of-the-art models in each task to instruction fine-tuned LLMs using multi-prompt. In Table 2, we report previous SOTA results for each task: an 11B T5-based text simplification model trained using control tokens corresponding to simplification operations (Sheang and Saggion, 2021), the EN-Cs results for the WMT '22 winning submission (Kocmi et al., 2022) and StarCoder 15B, a code infilling and generation LLM (Li et al., 2023), not explicitly trained to follow natural language instructions. For text simplification model of com-

MBR Utility Metric	Evaluation Metric					
	Text Simplification (LLaMA 7B Chat)					
	SARI	BERTSCORE	LENS	LENS-SALSA <sup>RF</sup>	SLE <sup>RF</sup>	
SARI		+1.08*	+1.06*	+7.24*	+4.33*	+0.38*
BERTSCORE	+1.44*		+1.09*	+6.18*	+3.11*	+0.45*
LENS	-0.67	-0.05		+5.78*	+4.69*	+0.82*
LENS-SALSA <sup>RF</sup>	-0.83	+0.35*	+8.10*		+4.65*	+0.97*
SLE <sup>RF</sup>	-5.25	-4.71	+2.39*	-4.51		+1.05*
MBR Utility Metric	Translation (ALMA 7B)					
	BLEU	BERTSCORE	COMET-22	COMETKIWI <sup>RF</sup>	XCOMET	METRICX-QE <sup>RF</sup>
BLEU		+0.34*	+0.47*	+0.67*	-0.14	+0.04 +0.11*
BERTSCORE	+0.51*		+1.59*	+1.68*	+2.48*	+0.22* +0.29*
COMET-22	+0.71*	+0.89*		+1.72*	+3.29*	+0.13* +0.18*
COMETKIWI <sup>RF</sup>	+0.80*	+1.03*	+1.06*		+2.87*	+0.07* +0.08*
XCOMET	+0.14	+0.85*	+0.84*	+3.34*		+0.09* +0.04*
METRICX	+0.36*	+0.81*	+0.36	+3.93*	+0.07*	-0.04
METRICX-QE <sup>RF</sup>	+0.60*	+1.68*	+2.11*	+5.31*	+0.08*	+0.03*

Table 3:  $\Delta$  metric improvement from single prompt to multi-prompt across metrics. RF = Reference-free reranker. \* = Statistically significant improvement with  $p < 0.05$ . For absolute performance, see Table 6.

parable size only surpass fine-tuned performance when using multi-prompt, with LLaMA 13B showing a +5 LENS over fine-tuned T5 11B.

#### 5.4 Is multi-prompt MBR over-fitting to the utility metric?

An inherent challenge of evaluating MBR is that the utility metric used to select candidates is typically also used for the final evaluation, in such cases it is difficult to attribute the metric improvement to higher quality generation (Bertsch et al., 2023). Given growing attention to metric development, we leverage various trained metrics to test whether multi-prompt using one utility metric improves performance cross all other utility metrics. We experiment with traditional overlap-based metrics, (BLEU, SARI), embedding similarity (BERTSCORE), small ( $\sim 100M$  parameter) trained metrics with references (LENS, COMET-22) and without references (COMETKIWI, LENS-SALSA, SLE), and large (3B+ parameter) trained metrics (XCOMET, METRICX, METRICX-QE). These metrics represent diverse text evaluation approaches and encompass the full state of evaluation in both tasks. We include a full description of metric architectures in Appendix B.1.

**Multi-prompt MBR Improves Across Metrics.** Table 3 reports results for cross-metric evaluation,

with the diagonal reflecting the traditional MBR evaluation setup (i.e., calculate MBR and evaluate using the same metric) and other cells indicate generalization from one metric to all others. We also perform a hypothesis test for the statistical significance of multi-prompt outperforming single prompt using bootstrap sampling (Berg-Kirkpatrick et al., 2012) with  $b = 10^3$ . Multi-prompt improves performance on most evaluation setups, with a few notable exceptions such as disagreement between trained and overlap-based metrics for simplification and COMET-based metrics for translation. For simplification, trained metrics’ failure when evaluated by SARI and BERTSCORE may be a byproduct of the test set size, as these metrics typically require a substantial number of references for stable evaluation (Alva-Manchego et al., 2020), more than what are provided in SIMPEVAL. Interestingly, the magnitude of performance improvement is highly variable to the specific utility metric, with no clear relationship between the metric architecture and improvement of multi-prompt, but typically a lower baseline performance indicates multi-prompt performs better (Table 6 in Appendix for more details).

#### 5.5 How does the metric type impact multi-prompt MBR?

As discussed by Fernandes et al. (2022), the MBR operation requires each candidate evaluate against every other candidate (i.e.,  $\mathcal{O}(n^2)$  comparisons), this becomes inefficient in practice for a large  $n$ , especially when using a trained utility metric. Therefore, we explore multi-prompt MBR alternatives using reference-free utility metrics:

- **Reranker.** Re-ranking directly estimates the quality of each candidate using a reference-free metric:  $\hat{y}_{\text{MBR}} = \arg \max_{y \in \mathcal{H}} [U(y)]$ . We use the trained LENS-SALSA for simplification (Heineman et al., 2023) and COMET-MQM (Rei et al., 2021) for translation. For code generation, we use Code Reviewer (Shi et al., 2022), which calculates agreement between the per-token probability of the generation given the docstring and the original docstring given the generation. Reference-free re-ranking simply requires  $\mathcal{O}(n)$  metric calculations as it directly estimates generation quality.
- **Reranker + MBR.** We use a two-stage MBR selection where we first rerank all candidates and select the top  $m$  to use for MBR, where the reranker can distill the candidate set and the expensive MBR metric can perform the final selection.



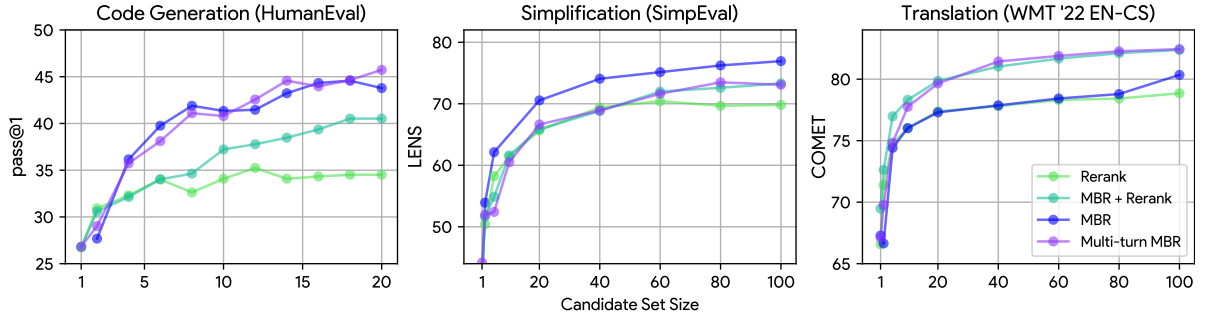


Figure 6: Alternative MBR formulations for multi-prompt across candidate set sizes for code generation, text simplification and translation. Efficient MBR methods show inconsistent results, dependent on task and metric.

- **Multi-turn MBR.** Similar to the previous approach, we select the top  $m$  MBR candidates and re-compute MBR using the smaller candidate set.

**Results.** We report results across candidate selection methods in Figure 6, finding the multi-prompt achieves performance improvement across reference-based and reference-free metrics, yet the relative performance of methods varies between tasks. With text simplification, we find the more expensive MBR performs better than the reference-free alternatives. For translation, both using a re-ranker first to narrow the candidate set (MBR + Rerank) and iteratively performing MBR (Multi-turn MBR) outperform vanilla MBR, despite these methods being more computationally efficient. We speculate the first pass may prune the lowest quality generations such that the second pass only considers a distilled candidate set, which better informs the MBR calculation. For code generation, we find the re-ranker performs relatively poorly compared to MBR, which may be reflective of the performance of Code Reviewer compared to MBR-EXEC, as the latter has access to multiple test cases.

## 6 Related Work

**Prompt Selection.** Current work on prompting for text generation has instead focused on optimization, such as in-context example selection (Min et al., 2022), example ordering (Lu et al., 2022) and prompt selection (Gonen et al., 2023). Notably, Agrawal et al. (2023) show selecting in-context examples for MT by maximizing  $n$ -gram overlap between the source and examples improves few-shot performance. Zhou et al. (2023) experiment with LLMs as prompt generators, and Yang et al. (2023) show using LLMs to iteratively rewrite prompts on a development set can distill a single, high-performant prompt. Our work uses LLM-written prompts and basic heuristics to distill the prompt bank, further improving multi-prompt.

**Output Selection.** Ensembling outputs under a candidate space has become a popular technique for improving LLM performance in classification tasks, such as majority vote over prompt chains (Wang et al., 2023), or merging outputs from multiple models (Kobayashi, 2018; Martínez Lorenzo et al., 2023). To our knowledge this work is the first to apply a multi-prompt approach to text generation.

**MBR Decoding.** Automatic evaluators have been incorporated into the training signal for task-specific models (Shen et al., 2016), used to improve the decoding process (Shen et al., 2004) and even evaluate the metrics themselves (Amrhein and Sennrich, 2022). MBR decoding has been explored extensively in improving translation quality (Kumar and Byrne, 2004; Eikema and Aziz, 2020; Müller and Sennrich, 2021) and has been proposed for text simplification (Maddela et al., 2023), summarization and style transfer (Suzgun et al., 2023). While our work is the first to propose generating the MBR hypothesis space using a prompt bank, Farinhas et al. (2023) perform preliminary experiments with paraphrases of a single sentence prompt, but found no difference in performance. Recent work argues sampling strategies like nucleus (Eikema and Aziz, 2022) or epsilon (Freitag et al., 2023) offer slightly better performance over beam search for MBR, with this work extending their findings by attributing candidate set quality to sampling diversity.

## 7 Conclusion

In this work, we propose multi-prompt, a generalized case of MBR for conditional text generation. Multi-prompt successfully ensembles outputs of instruction fine-tuned language models across prompt constructions and in-context examples. We highlight the importance of prompt selection and sampling when constructing the prompt bank with top- $p$  prompt sampling and further verify our results across tasks, models and utility metrics.



## Limitations

We limit our study of the prompt bank to a basic set of seed prompts and GPT-written paraphrases for each task. Notably, we do not study the impact of prompt formats (e.g., `passage:{ }\n answer{ }` vs. `Passage: :{ } Answer: :{ }`, Sclar et al., 2023), in-context example ordering (Lu et al., 2022) or example selection (Agrawal et al., 2023) on multi-prompt performance, although multi-prompt may extend to such methods. We leave the question of exhaustively constructing a prompt bank to future work, perhaps by extending work in prefix tuning (Li and Liang, 2021).

An inherent limitation of MBR is the increase in inference time, where we generate up to 100 samples in our experiments, and use a neural utility metric with either linear or quadratic comparisons between candidates. While recent work has lowered the number of metric comparisons (Cheng and Vlachos, 2023), MBR remains prohibitively expensive for use in compute-limited scenarios.

## References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Sweta Agrawal, Chunting Zhou, Mike Lewis, Luke Zettlemoyer, and Marjan Ghazvininejad. 2023. *In-context examples selection for machine translation*. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 8857–8873, Toronto, Canada. Association for Computational Linguistics.
- Fernando Alva-Manchego, Louis Martin, Antoine Bordes, Carolina Scarton, Benoît Sagot, and Lucia Specia. 2020. *ASSET: A dataset for tuning and evaluation of sentence simplification models with multiple rewriting transformations*. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4668–4679, Online. Association for Computational Linguistics.
- Fernando Alva-Manchego, Carolina Scarton, and Lucia Specia. 2021. The (un) suitability of automatic evaluation metrics for text simplification. *Computational Linguistics*, 47(4):861–889.
- Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, et al. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.
- Chantal Amrhein and Rico Sennrich. 2022. *Identifying weaknesses in machine translation metrics through minimum Bayes risk decoding: A case study for COMET*. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1125–1141, Online only. Association for Computational Linguistics.
- 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.
- Taylor Berg-Kirkpatrick, David Burkett, and Dan Klein. 2012. *An empirical investigation of statistical significance in NLP*. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 995–1005, Jeju Island, Korea. Association for Computational Linguistics.
- Amanda Bertsch, Alex Xie, Graham Neubig, and Matthew Gormley. 2023. *It’s MBR all the way down: Modern generation techniques through the lens of minimum Bayes risk*. In *Proceedings of the Big Picture Workshop*, pages 108–122, Singapore. Association for Computational Linguistics.
- Peter J Bickel and Kjell A Doksum. 1977. *Mathematical statistics: Basic ideas and selected topics, volumes I-II package*. Chapman and Hall/CRC.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Julius Cheng and Andreas Vlachos. 2023. *Faster minimum Bayes risk decoding with confidence-based pruning*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12473–12480, Singapore. Association for Computational Linguistics.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2022. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*.
- Liam Cripwell, Joël Legrand, and Claire Gardent. 2023. *Simplicity level estimate (SLE): A learned referenceless metric for sentence simplification*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12053–12059,

675	Singapore. Association for Computational Linguistics.	
676		
677	Bryan Eikema and Wilker Aziz. 2020. <a href="#">Is MAP decoding all you need? the inadequacy of the mode in neural machine translation</a> . In <i>Proceedings of the 28th International Conference on Computational Linguistics</i> , pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.	
678		
679		
680		
681		
682		
683	Bryan Eikema and Wilker Aziz. 2022. <a href="#">Sampling-based approximations to minimum Bayes risk decoding for neural machine translation</a> . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	
684		
685		
686		
687		
688		
689		
690	António Farinhas, José G. C. de Souza, and André F. T. Martins. 2023. <a href="#">An empirical study of translation hypothesis ensembling with large language models</a> . <i>Preprint</i> , arXiv:2310.11430.	
691		
692		
693		
694	Christian Federmann, Tom Kocmi, and Ying Xin. 2022. <a href="#">NTREX-128 – news test references for MT evaluation of 128 languages</a> . In <i>Proceedings of the First Workshop on Scaling Up Multilingual Evaluation</i> , pages 21–24, Online. Association for Computational Linguistics.	
695		
696		
697		
698		
699		
700	Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. <a href="#">Quality-aware decoding for neural machine translation</a> . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1396–1412, Seattle, United States. Association for Computational Linguistics.	
701		
702		
703		
704		
705		
706		
707		
708		
709	Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. <a href="#">Experts, errors, and context: A large-scale study of human evaluation for machine translation</a> . <i>Transactions of the Association for Computational Linguistics</i> , 9:1460–1474.	
710		
711		
712		
713		
714		
715	Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. <a href="#">Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 9198–9209, Singapore. Association for Computational Linguistics.	
716		
717		
718		
719		
720		
721		
722	Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022a. <a href="#">High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics</a> . <i>Transactions of the Association for Computational Linguistics</i> , 10:811–825.	
723		
724		
725		
726		
727	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022b. <a href="#">Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more</a>	
728		
729		
730		
731		
	<a href="#">robust</a> . In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 46–68, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	732
		733
		734
		735
	Hila Gonen, Srini Iyer, Terra Blevins, Noah Smith, and Luke Zettlemoyer. 2023. <a href="#">Demystifying prompts in language models via perplexity estimation</a> . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 10136–10148, Singapore. Association for Computational Linguistics.	736
		737
		738
		739
		740
		741
	Naman Goyal, Jingfei Du, Myle Ott, Giri Anantharaman, and Alexis Conneau. 2021. <a href="#">Larger-scale transformers for multilingual masked language modeling</a> . In <i>Proceedings of the 6th Workshop on Representation Learning for NLP (RepL4NLP-2021)</i> , pages 29–33, Online. Association for Computational Linguistics.	742
		743
		744
		745
		746
		747
		748
	Nuno M Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André FT Martins. 2023. <a href="#">xCOMET: Transparent machine translation evaluation through fine-grained error detection</a> . <i>arXiv preprint arXiv:2310.10482</i> .	749
		750
		751
		752
		753
	David Heineman, Yao Dou, Mounica Maddela, and Wei Xu. 2023. <a href="#">Dancing between success and failure: Edit-level simplification evaluation using SALSA</a> . In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 3466–3495, Singapore. Association for Computational Linguistics.	754
		755
		756
		757
		758
		759
		760
	Amr Hendy, Mohamed Abdelrehim, Amr Sharaf, Vikas Raunak, Mohamed Gabr, Hitokazu Matsushita, Young Jin Kim, Mohamed Afify, and Hany Hassan Awadalla. 2023. <a href="#">How good are GPT models at machine translation? A comprehensive evaluation</a> . <i>arXiv preprint arXiv:2302.09210</i> .	761
		762
		763
		764
		765
		766
	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. 2020. <a href="#">The curious case of neural text de-generation</a> . In <i>International Conference on Learning Representations</i> .	767
		768
		769
		770
	T Jaeger and Roger Levy. 2006. <a href="#">Speakers optimize information density through syntactic reduction</a> . <i>Advances in neural information processing systems</i> , 19.	771
		772
		773
	Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al. 2023. <a href="#">Mistral 7B</a> . <i>arXiv preprint arXiv:2310.06825</i> .	774
		775
		776
		777
		778
	Chao Jiang, Mounica Maddela, Wuwei Lan, Yang Zhong, and Wei Xu. 2020. <a href="#">Neural CRF model for sentence alignment in text simplification</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7943–7960, Online. Association for Computational Linguistics.	779
		780
		781
		782
		783
		784
	Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. <a href="#">MetricX-23: The Google submission to the</a>	785
		786
		787

788	WMT 2023 metrics shared task. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 756–767, Singapore. Association for Computational Linguistics.	846
789		847
790		848
791		
792	Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022. DEMETR: Diagnosing evaluation metrics for translation. In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	849
793		850
794		851
795		852
796		853
797		854
798		855
799		856
800	Daniel Khashabi, Xinxi Lyu, Sewon Min, Lianhui Qin, Kyle Richardson, Sean Welleck, Hannaneh Hajishirzi, Tushar Khot, Ashish Sabharwal, Sameer Singh, and Yejin Choi. 2022. Prompt waywardness: The curious case of discretized interpretation of continuous prompts. In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 3631–3643, Seattle, United States. Association for Computational Linguistics.	857
801		858
802		859
803		860
804		861
805		862
806		863
807		
808		
809		
810	Hayato Kobayashi. 2018. Frustratingly easy model ensemble for abstractive summarization. In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 4165–4176, Brussels, Belgium. Association for Computational Linguistics.	864
811		865
812		866
813		867
814		868
815		869
816	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 <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.	870
817		871
818		
819		
820		
821		
822		
823		
824		
825		
826		
827		
828	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.	872
829		873
830		874
831		875
832		876
833		877
834		878
835		879
836	Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, et al. 2023. Starcoder: May the source be with you! <i>arXiv preprint arXiv:2305.06161</i> .	880
837		881
838		882
839		883
840		884
841	Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. 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 4582–4597, Online. Association for Computational Linguistics.	885
842		886
843		887
844		888
845		889
	Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and Pontus Stenetorp. 2022. Fantastically ordered prompts and where to find them: Overcoming few-shot prompt order sensitivity. In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 8086–8098, Dublin, Ireland. Association for Computational Linguistics.	890
		891
		892
		893
		894
		895
		896
		897
	Mounica Maddela, Yao Dou, David Heineman, and Wei Xu. 2023. LENS: A learnable evaluation metric for text simplification. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 16383–16408, Toronto, Canada. Association for Computational Linguistics.	898
		899
		900
		901
		902
		903
	Abelardo Carlos Martínez Lorenzo, Pere Lluís Huguet Cabot, and Roberto Navigli. 2023. AMRs assemble! learning to ensemble with autoregressive models for AMR parsing. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1595–1605, Toronto, Canada. Association for Computational Linguistics.	904
		905
		906
		907
		908
		909
		910
		911
		912
		913
		914
		915
		916
		917
		918
		919
		920
		921
		922
		923
		924
		925
		926
		927
		928
		929
		930
		931
		932
		933
		934
		935
		936
		937
		938
		939
		940
		941
		942
		943
		944
		945
		946
		947
		948
		949
		950
		951
		952
		953
		954
		955
		956
		957
		958
		959
		960
		961
		962
		963
		964
		965
		966
		967
		968
		969
		970
		971
		972
		973
		974
		975
		976
		977
		978
		979
		980
		981
		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992
		993
		994
		995
		996
		997
		998
		999
		1000



904	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane	961
905	Jing Zhu. 2002. <a href="#">Bleu: a method for automatic evalu-</a>	Suhr. 2023. Quantifying language models’ sensitiv-	962
906	<a href="#">ation of machine translation</a> . In <i>Proceedings of the</i>	ity to spurious features in prompt design or: How I	963
907	<i>40th Annual Meeting of the Association for Comput-</i>	learned to start worrying about prompt formatting.	964
908	<i>ational Linguistics</i> , pages 311–318, Philadelphia,	<i>arXiv preprint arXiv:2310.11324</i> .	965
909	Pennsylvania, USA. Association for Computational		
910	Linguistics.		
911	Matt Post. 2018. <a href="#">A call for clarity in reporting BLEU</a>	Kim Cheng Sheang and Horacio Saggion. 2021. <a href="#">Con-</a>	966
912	<a href="#">scores</a> . In <i>Proceedings of the Third Conference on</i>	<a href="#">trollable sentence simplification with a unified text-</a>	967
913	<i>Machine Translation: Research Papers</i> , pages 186–	<a href="#">to-text transfer transformer</a> . In <i>Proceedings of the</i>	968
914	191, Brussels, Belgium. Association for Computa-	<i>14th International Conference on Natural Language</i>	969
915	tional Linguistics.	<i>Generation</i> , pages 341–352, Aberdeen, Scotland, UK.	970
		Association for Computational Linguistics.	971
916	Ricardo Rei, Ana C Farinha, Chrysoula Zerva, Daan	Libin Shen, Anoop Sarkar, and Franz Josef Och. 2004.	972
917	van Stigt, Craig Stewart, Pedro Ramos, Taisiya	<a href="#">Discriminative reranking for machine translation</a> .	973
918	Glushkova, André F. T. Martins, and Alon Lavie.	In <i>Proceedings of the Human Language Technol-</i>	974
919	2021. <a href="#">Are references really needed? unbabel-IST</a>	<i>ogy Conference of the North American Chapter</i>	975
920	<a href="#">2021 submission for the metrics shared task</a> . In <i>Pro-</i>	<i>of the Association for Computational Linguistics:</i>	976
921	<i>ceedings of the Sixth Conference on Machine Trans-</i>	<i>HLT-NAACL 2004</i> , pages 177–184, Boston, Mas-	977
922	<i>lation</i> , pages 1030–1040, Online. Association for	sachusetts, USA. Association for Computational Lin-	978
923	Computational Linguistics.	guistics.	979
924	Ricardo Rei, Nuno M. Guerreiro, Josã© Pombal, Daan	Shiqi Shen, Yong Cheng, Zhongjun He, Wei He, Hua	980
925	van Stigt, Marcos Treviso, Luisa Coheur, José G.	Wu, Maosong Sun, and Yang Liu. 2016. <a href="#">Minimum</a>	981
926	C. de Souza, and André Martins. 2023. <a href="#">Scaling up</a>	<a href="#">risk training for neural machine translation</a> . In <i>Pro-</i>	982
927	<a href="#">CometKiwi: Unbabel-IST 2023 submission for the</a>	<i>ceedings of the 54th Annual Meeting of the Associa-</i>	983
928	<a href="#">quality estimation shared task</a> . In <i>Proceedings of the</i>	<i>tion for Computational Linguistics (Volume 1: Long</i>	984
929	<i>Eighth Conference on Machine Translation</i> , pages	<i>Papers)</i> , pages 1683–1692, Berlin, Germany. Associ-	985
930	841–848, Singapore. Association for Computational	ation for Computational Linguistics.	986
931	Linguistics.		
932	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon	Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke	987
933	Lavie. 2020. <a href="#">COMET: A neural framework for MT</a>	Zettlemoyer, and Sida I. Wang. 2022. <a href="#">Natural lan-</a>	988
934	<a href="#">evaluation</a> . In <i>Proceedings of the 2020 Conference</i>	<a href="#">guage to code translation with execution</a> . In <i>Proce-</i>	989
935	<i>on Empirical Methods in Natural Language Process-</i>	<i>edings of the 2022 Conference on Empirical Methods</i>	990
936	<i>ing (EMNLP)</i> , pages 2685–2702, Online. Association	<i>in Natural Language Processing</i> , pages 3533–3546,	991
937	for Computational Linguistics.	Abu Dhabi, United Arab Emirates. Association for	992
		Computational Linguistics.	993
938	Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro,	Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, and	994
939	Chrysoula Zerva, Ana C Farinha, Christine Maroti,	et al. 2023. <a href="#">Beyond the imitation game: Quantifying</a>	995
940	José G. C. de Souza, Taisiya Glushkova, Duarte	<a href="#">and extrapolating the capabilities of language models</a> .	996
941	Alves, Luisa Coheur, Alon Lavie, and André F. T.	<i>Transactions on Machine Learning Research</i> .	997
942	Martins. 2022. <a href="#">CometKiwi: IST-unbabel 2022 sub-</a>		
943	<a href="#">mission for the quality estimation shared task</a> . In	Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky.	998
944	<i>Proceedings of the Seventh Conference on Machine</i>	2023. <a href="#">Follow the wisdom of the crowd: Effective</a>	999
945	<i>Translation (WMT)</i> , pages 634–645, Abu Dhabi,	<a href="#">text generation via minimum Bayes risk decoding</a> .	1000
946	United Arab Emirates (Hybrid). Association for Com-	In <i>Findings of the Association for Computational</i>	1001
947	putational Linguistics.	<i>Linguistics: ACL 2023</i> , pages 4265–4293, Toronto,	1002
		Canada. Association for Computational Linguistics.	1003
948	Nils Reimers and Iryna Gurevych. 2019. <a href="#">Sentence-</a>	Gemma Team, Thomas Mesnard, Cassidy Hardin,	1004
949	<a href="#">BERT: Sentence embeddings using Siamese BERT-</a>	Robert Dadashi, Surya Bhupatiraju, Shreya Pathak,	1005
950	<a href="#">networks</a> . In <i>Proceedings of the 2019 Conference on</i>	Laurent Sifre, Morgane Rivi�re, Mihir Sanjay Kale,	1006
951	<i>Empirical Methods in Natural Language Processing</i>	Juliette Love, et al. 2024. Gemma: Open models	1007
952	<i>and the 9th International Joint Conference on Natu-</i>	based on gemini research and technology. <i>arXiv</i>	1008
953	<i>ral Language Processing (EMNLP-IJCNLP)</i> , pages	<i>preprint arXiv:2403.08295</i> .	1009
954	3982–3992, Hong Kong, China. Association for Com-		
955	putational Linguistics.	Hugo Touvron, Louis Martin, Kevin Stone, Peter Al-	1010
956	Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten	bert, Amjad Almahairi, Yasmine Babaei, Nikolay	1011
957	Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi,	Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti	1012
958	Jingyu Liu, Tal Remez, J�r�my Rapin, et al. 2023.	Bhosale, et al. 2023. Llama 2: Open founda-	1013
959	Code LLaMA: Open foundation models for code.	<i>tion and fine-tuned chat models. arXiv preprint</i>	1014
960	<i>arXiv preprint arXiv:2308.12950</i> .	<i>arXiv:2307.09288</i> .	1015
		Ahmet �st�n, Viraat Aryabumi, Zheng-Xin Yong, Wei-	1016
		Yin Ko, Daniel D’souza, Gbemileke Onilude, Neel	1017



1018	Bhandari, Shivalika Singh, Hui-Lee Ooi, Amr Kayid,	Tianyi Zhang, Tao Yu, Tatsunori Hashimoto, Mike	1074
1019	et al. 2024. Aya model: An instruction finetuned	Lewis, Wen-tau Yih, Daniel Fried, and Sida Wang.	1075
1020	open-access multilingual language model. <i>arXiv</i>	2023. <i>Coder reviewer reranking for code generation.</i>	1076
1021	<i>preprint arXiv:2402.07827.</i>	In <i>International Conference on Machine Learning</i> ,	1077
		pages 41832–41846. PMLR.	1078
1022	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le,	Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han,	1079
1023	Ed H. Chi, Sharan Narang, Aakanksha Chowdhery,	Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy	1080
1024	and Denny Zhou. 2023. <a href="#">Self-consistency improves</a>	Ba. 2023. <a href="#">Large language models are human-level</a>	1081
1025	<a href="#">chain of thought reasoning in language models.</a> In	<a href="#">prompt engineers.</a> In <i>The Eleventh International</i>	1082
1026	<i>The Eleventh International Conference on Learning</i>	<i>Conference on Learning Representations.</i>	1083
1027	<i>Representations.</i>		
1028	Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu,		
1029	Adams Wei Yu, Brian Lester, Nan Du, Andrew M.		
1030	Dai, and Quoc V Le. 2022. <a href="#">Finetuned language mod-</a>		
1031	<a href="#">els are zero-shot learners.</a> In <i>International Confer-</i>		
1032	<i>ence on Learning Representations.</i>		
1033	Jules White, Quchen Fu, Sam Hays, Michael Sandborn,		
1034	Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse		
1035	Spencer-Smith, and Douglas C Schmidt. 2023. A		
1036	prompt pattern catalog to enhance prompt engineer-		
1037	ing with ChatGPT. <i>arXiv preprint arXiv:2302.11382.</i>		
1038	Haoran Xu, Young Jin Kim, Amr Sharaf, and		
1039	Hany Hassan Awadalla. 2023. A paradigm shift		
1040	in machine translation: Boosting translation perfor-		
1041	mance of large language models. <i>arXiv preprint</i>		
1042	<i>arXiv:2309.11674.</i>		
1043	Haoran Xu, Amr Sharaf, Yunmo Chen, Weiting Tan,		
1044	Lingfeng Shen, Benjamin Van Durme, Kenton Mur-		
1045	ray, and Young Jin Kim. 2024. Contrastive prefer-		
1046	ence optimization: Pushing the boundaries of LLM		
1047	performance in machine translation. <i>arXiv preprint</i>		
1048	<i>arXiv:2401.08417.</i>		
1049	Wei Xu, Chris Callison-Burch, and Courtney Napoles.		
1050	2015. <a href="#">Problems in current text simplification re-</a>		
1051	<a href="#">search: New data can help.</a> <i>Transactions of the Asso-</i>		
1052	<i>ciation for Computational Linguistics</i> , 3:283–297.		
1053	Wei Xu, Courtney Napoles, Ellie Pavlick, Quanze Chen,		
1054	and Chris Callison-Burch. 2016. <a href="#">Optimizing sta-</a>		
1055	<a href="#">tistical machine translation for text simplification.</a>		
1056	<i>Transactions of the Association for Computational</i>		
1057	<i>Linguistics</i> , 4:401–415.		
1058	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,		
1059	Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and		
1060	Colin Raffel. 2021. <a href="#">mT5: A massively multilingual</a>		
1061	<a href="#">pre-trained text-to-text transformer.</a> In <i>Proceedings</i>		
1062	<i>of the 2021 Conference of the North American Chap-</i>		
1063	<i>ter of the Association for Computational Linguistics:</i>		
1064	<i>Human Language Technologies</i> , pages 483–498, On-		
1065	line. Association for Computational Linguistics.		
1066	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu,		
1067	Quoc V Le, Denny Zhou, and Xinyun Chen. 2023.		
1068	Large language models as optimizers. <i>arXiv preprint</i>		
1069	<i>arXiv:2309.03409.</i>		
1070	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q.		
1071	Weinberger, and Yoav Artzi. 2020. <a href="#">BERTScore:</a>		
1072	<a href="#">Evaluating text generation with BERT.</a> In <i>Internat-</i>		
1073	<i>ional Conference on Learning Representations.</i>		

Human-Written Text Simplification Prompt
Rewrite the following complex sentence in order to make it easier to understand by non-native speakers of English. You can do so by replacing complex words with simpler synonyms (i.e. paraphrasing), deleting unimportant information (i.e. compression), and/or splitting a long complex sentence into several simpler ones. The final simplified sentence needs to be grammatical, fluent, and retain the main ideas of its original counterpart without altering its meaning.
Simplify the sentence please.
You are an artificial intelligence designed to simplify human written text. The text you are given will contain complex ideas, phrases or concepts and your job is to rewrite that text in a simple and easy to understand way. Your simplification should be completely fluent and retain the ideas of the simplification.
I would like you to simplify the following sentence such that the text is as concise and easy to read as possible.
Text simplification is an operation used in natural language processing to change, enhance, classify, or otherwise process an existing body of human-readable text so its grammar and structure is greatly simplified while the underlying meaning and information remain the same. Text simplification is an important area of research because of communication needs in an increasingly complex and interconnected world more dominated by science, technology, and new media. But natural human languages pose huge problems because they ordinarily contain large vocabularies and complex constructions that machines, no matter how fast and well-programmed, cannot easily process. However, researchers have discovered that, to reduce linguistic diversity, they can use methods of semantic compression to limit and simplify a set of words used in given texts. Please simplify the following sentence.
Please simplify the below sentence by using a combination of these three operations. Elaboration. An addition of meaningful, relevant and correct information, such as clarifying vague terminology, providing background information on an entity or subject, or explicating general world knowledge unknown to the audience. Generalization. A deletion of unnecessary, irrelevant or complicated concepts. Paraphrase. Swapping complex spans with equivalent, simpler alternatives. The final sentence should be grammatical, concise and easier to read compared to the original sentence.
You are an AI assistant that writes text simplification. Text simplification can be defined as any process that reduces the syntactic or lexical complexity of a text while attempting to preserve its meaning and information content. The aim of text simplification is to make text easier to comprehend for a human user, or process by a program. Please simplify the following sentence.
Simplify.
You are to act as a text simplification bot. As a text simplification bot, you will simplify the following sentence such that it is syntactically easier to read and semantically easier to understand. Please do not make the text more complex, longer or difficult for a reader.
I am writing a sentence, please take a look at this sentence and write a simpler version such that a non-English speaker or an individual with disabilities could better understand the sentence.

Table 4: Text simplification prompts used for the decoding experiment in Figure 3 and used as examples to write GPT-4 prompts for experiments in §5.

## A Prompt Bank Construction

Table 4 contains the human-written prompts for text simplification. These human-written prompts are provided as examples to GPT-4 when automatically generating prompts for large-scale experiments in §5. For code generation, we extract the docstring in the original HUMANEVAL examples as the human-written prompt, and provide it as an example prompt to GPT-4. For machine translation, our few-shot examples were sampled randomly from the WMT newstest19 test corpus (Barrault et al., 2019).

## B Detailed System Descriptions

In this section, we include a full description of the generation models and utility metrics used in experiments throughout §5.3 and §5.4. All experiments were inference-based and were run on up to 4xN-

Prompt-Generation Instruction
Please write a variation of the following instruction for a coding task. You may be creative in proposing potential solutions, or explaining the nature of the task. Please do not write any examples. Example: {example_prompt} Prompt:
Create a prompt for a language model to simplify a sentence, this prompt will explain the text simplification task and instructions for how to perform the task. The prompt should be diverse, include a description of simplification and clearly state what is expected of the language model. Example: {example_prompt_1} Example: {example_prompt_2} Prompt:

Table 5: Instruction templates provided to GPT-4 when generating task instructions for code generation (top) and text simplification (bottom).

VIDIA A40 GPUs, depending on the requirements of the specific model or utility metric. The use of models, metrics and datasets in this project follows their respective licenses and intended use.

## B.1 Utility Metrics

### B.1.1 Simplification

SARI (Xu et al., 2016) is an  $n$ -gram overlap based metric that compares edits on inputs, outputs and a bank of references.

BERTSCORE (Zhang et al., 2020) calculates a word-level cosine similarity of BERT embeddings. Alva-Manchego et al. (2021) find BERTSCORE is an adequate measure of quality generation, but that it does not correlate with simplicity.

LENS (Maddela et al., 2023) is a RoBERTa-based metric trained using human ratings of text simplification model outputs. The authors train on an adaptive loss to allow a high score for generations was close to *any* references, encouraging the metric to consider different simplification types.

LENS-SALSA (Heineman et al., 2023) extends the LENS architecture by fine-tuning on a dual sentence- and word-level quality objective. The authors show LENS-SALSA is more sensitive to specific edit operations, while not requiring any reference simplifications.

SLE (Cripwell et al., 2023) is a RoBERTa-based metric trained to estimate the simplicity of text, with the simplicity score defined as the difference in simplicity between the complex and simplified sentences. SLE was trained on 0-4 readability scores of news articles in the Newsela corpus (Xu et al., 2015), with an additional label softening for individual sentences in the corpus.

### B.1.2 Translation

**BLEU** (Papineni et al., 2002) is an  $n$ -gram overlap based metric comparing a translation to a bank of references. BLEU remains a widely-used standard for automatic evaluation, despite lower correlation to human judgement compared to learned metrics (Freitag et al., 2022b). We use the ScaBLEU implementation (Post, 2018).

**COMET** (Rei et al., 2020) is a widely used RoBERTa-based metric, trained on direct assessments of simplification quality. For reference-free evaluation, we use the CometKiwi-XXL variant (Rei et al., 2022, 2023), trained to predict sentence- and word-level scores simultaneously.

**XCOMET** (Guerreiro et al., 2023) is a fine-tuned XLM-R model (Goyal et al., 2021) based on the CometKiwi architecture, but scaling the model size and training data, including with synthetic data created by randomly swapping  $n$ -grams or entire sentences with unrelated translations. We use the 11B XCOMET-XXL in our experiments.

**METRICX** (Juraska et al., 2023) is a recent fine-tuned 11B mT5-XXL (Xue et al., 2021) trained on DA data from 2015-20, MQM data from 2020-21 (Freitag et al., 2021) and synthetic data based on the MQM and DEMETR (Karpinska et al., 2022) taxonomies of translation errors. Notably, the MetricX architecture encodes both candidates and references together, while COMET encodes both separately and combines the outputs to calculate the final score. We also use the QE variant METRICX-QE trained without references. The WMT '22 test data used in this work is not included in the training data of any translation metrics we considered.

### B.1.3 Code Generation

**MBR-EXEC** (Shi et al., 2022) executes candidate generations on a series of test cases, and selects the candidate with the highest agreement on its output with all other candidates. While the authors do not evaluate on HUMANEVAL, we replicate the setup in Zhang et al. (2023) by using the test cases in the docstring to calculate the agreement. We use a soft loss over all test cases, as many HUMANEVAL docstring examples are trivial or edge cases. If two candidates have the same MBR score, we break ties using the candidate with higher probability under the language model.

**Code Reviewer** (Zhang et al., 2023) attempts to find a consensus between the likelihood of the generated program  $p(y|x)$  and the original docstring

using a minified version of the generation  $p(x|y)$ . We use their implementation for rejecting degenerate samples, minifying code and calculating the reviewer score. We use the same models for generation and re-ranking.

## B.2 Model Architectures

### B.2.1 Simplification

**Instruction Fine-tuned Models.** We experiment with widely used instruction fine-tuned LLMs, aiming for a broad coverage of current models: Llama 2 Chat (Touvron et al., 2023), Gemma (Team et al., 2024) and Mistral (Jiang et al., 2023).

**Fine-tuned Control T5** (Sheang and Saggion, 2021) is a T5-based text simplification model fine-tuned on the Wiki-Auto (Jiang et al., 2020) dataset of aligned English-Simple English Wikipedia articles. We use their same control token setup: `<NC_0.95> <LS_0.75> <DR_0.75> <WR_0.75>`.

### B.2.2 Translation

**ALMA-R** (Xu et al., 2024) is a class of translation LLMs. The base ALMA (Xu et al., 2023) is a fine-tuned LLaMA model with text in each target language and then parallel translation data. ALMA-R is an extension trained on a contrastive preference loss to incorporate ratings of translation quality.

**TowerInstruct** (Alves et al., 2024) is a fine-tuned Llama 2 model on multi-lingual instructions, aiming to incorporate tasks beyond translation, such as paraphrasing, post editing and grammar error correction.

**Aya 101** (Üstün et al., 2024) is an mT5-based model fine-tuned on multi-lingual data in 101 languages. While mT5 is instruction-following model, Aya is not fine-tuned on instruction data.

Additionally, we provide results from the WMT '22 winning submission, and the Microsoft Translate API, as reported in Hendy et al. (2023).

### B.2.3 Code Generation

**StarCoder 2** (Li et al., 2023) is trained from-scratch on 4T tokens from 600+ programming languages. Although the model is not instruction fine-tuned, we see a slight performance improvement with multi-prompt, likely because comments and code descriptions are included in its pre-training.

**CodeLLaMA** (Roziere et al., 2023) is a fine-tuned Llama 2 model on 500B-1T tokens of code-related datasets, including Python, substantially outperforming the base Llama 2 model on HumanEval.

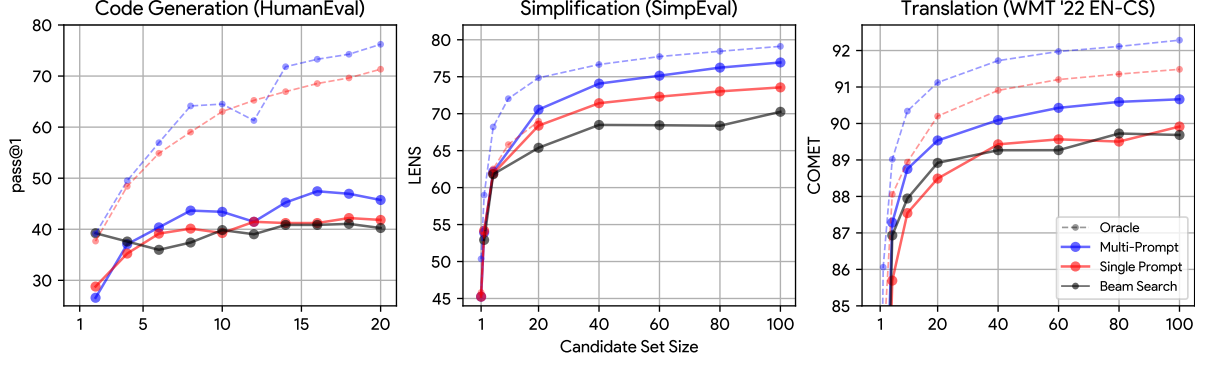


Figure 7: Multi-prompt, single prompt and beam search MBR decoding performance across candidate set sizes for code generation, text simplification and translation. Results averaged over 5 bootstrap iterations.

## C Further Results

### C.1 Beam Search & Oracle Performance

Following related work in MBR, we report upper-bound ‘oracle’ results (similar to Shi et al., 2022) and a lower-bound beam search baseline (similar to Freitag et al., 2023) in comparison to our main results (Figure 1) in Figure 7.

**Beam Search.** The MBR candidate set historically has consisted of the top beam search candidates, but as language models have become better generators recent work has argued sampling leads to a better estimation of the hypothesis space (Freitag et al., 2023). For this reason, we exclusively use nucleus sampling in §5, but we report beam search as a baseline in Figure 7, with a ‘candidate set size’ of  $n$  corresponding to the top  $n$  beam candidates, or  $n$  candidates with nucleus sampling for other results.

**Oracle.** As the final MBR performance can be impacted both by the quality of the candidate set and the choice of utility metric, we report an upper-bound performance by deliberately selecting the best candidate generations. Given a test set with gold-standard references  $\mathcal{R}$ , we define the oracle performance as the set of the highest scoring possible selection of candidates:

$$\text{Oracle}(\mathcal{R}^*) = \sum_{r \in \mathcal{R}^*} \max_{y \in \mathcal{H}} [U(y, r)] \quad (6)$$

Since code generation is evaluated using pass@1, its oracle uses expected pass@k (Shi et al., 2022), which measures whether at least one candidate within the candidate set passes all unit tests  $\mathcal{T}$ :

$$\text{ExPass}@K = \mathbb{E}_{|\mathcal{H}|=K} \left[ \max_{y \in \mathcal{H}} \min_{t \in \mathcal{T}} \mathbb{1}[t(y)] \right] \quad (7)$$

**Results.** As oracle performance measures candidate set quality independent of the utility metric,

we find an increase in oracle performance coincides with an improvement when using multi-prompt, indicating that a utility metric can naturally select candidates when the candidate set is higher quality. This suggests improving utility metrics may be a promising direction to bridge the gap between candidate quality and candidate selection. Beam search was a particularly strong baseline for small candidate set sizes, particularly for code generation, but beam search is not as sensitive to improvement as the candidate set size increases. Additionally, as code generation is evaluated using the binary pass@1 metric, rather than a scalar quality metric as used by translation and simplification, there is a large gap between MBR and oracle performance, also observed by Shi et al. (2022).

### C.2 En-XX Translation Results

For brevity, we limit our multi-prompt experiments to only the English-Czech language pair, but report results across the full ALMA test set, including WMT '22 test data and a subset of NTREX (Fedorov et al., 2022), in Figure 8, where we observe improvement with multi-prompt is dependent on the language pair. Generally, high resource languages (such as French, German, Russian) do not have a substantial difference, which may be a result of the low prompt sensitivity for such pairs.

### C.3 Detailed Multi-Model Results

See Figure 9 contains separated results for multi-prompt and single prompt for each model, as reported in Figure 5 and discussed in §5.3.



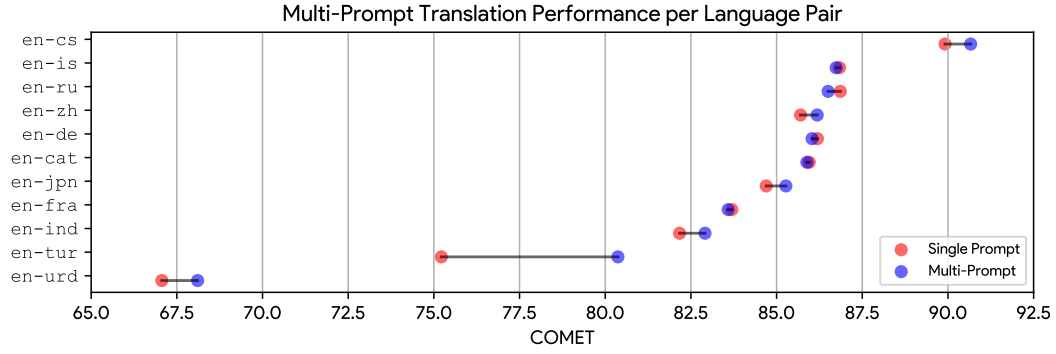


Figure 8: **Multi-prompt** and **single prompt** performance of ALMA 7B R across En-XX translation pairs. For low resource language pairs (e.g., Urdu, Turkish, Czech) we observe significant performance improvements, but not for most high resource pairs (e.g., French, German, Russian).

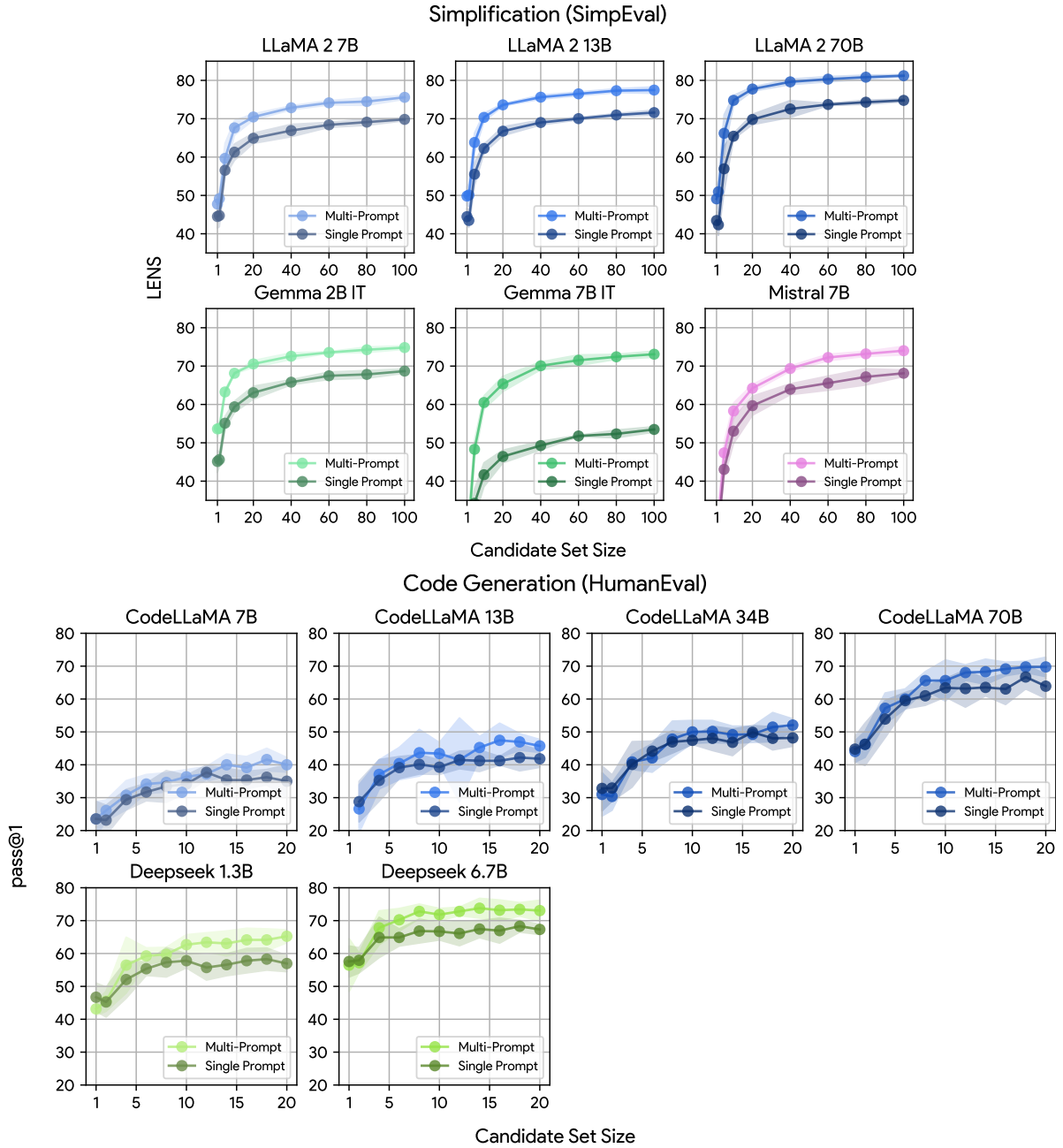


Figure 9: Results of multi-prompt across model sizes and architectures bootstrapped over 5 iterations with a 95% CI. Multi-prompt consistently improves performance across architectures and as models scale.

#### C.4 Detailed Cross Metric Evaluation

Table 6 contains the full results for the MBR experiments across metrics as discussed in §5.4. While evaluating on the same metric used for MBR clearly improves performance the most (see entries on the diagonal), we find multi-prompt performed on any metric universally improves performance when evaluated on any other metric. Recent neural metrics, which achieve higher correlation with human judgements, also have a higher overall performance. Note, METRICX scores translations on a  $[0, 25]$  scale corresponding to an MQM rating, where lower is better and SLE scores simplifications on a  $[0, 4]$  corresponding to a Newsela simplification rating, where higher is better. For clarity, we negate the METRICX results in Table 3 such that all the green cells indicate a metric improvement.

MBR Utility Metric	Evaluation Metric						Evaluation Metric						
	Text Simplification (LLaMA 7B Chat)						Text Simplification (LLaMA 7B Chat)						
	BERTSCORE LENS LENS-SALSA <sup>RF</sup> SLE <sup>RF</sup> SARI						BERTSCORE LENS LENS-SALSA <sup>RF</sup> SLE <sup>RF</sup> SARI						
	SARI	44.33	92.64	58.73	72.31	1.42	SARI	43.25	91.58	51.49	67.97	1.04	
	BERTSCORE	45.46	93.71	60.86	71.47	1.37	BERTSCORE	44.02	92.62	54.68	68.36	0.92	
	LENS	39.98	92.18	76.29	79.55	2.30	LENS	40.64	92.24	70.51	74.86	1.49	
	LENS-SALSA <sup>RF</sup>	38.55	91.29	73.31	84.59	2.47	LENS-SALSA <sup>RF</sup>	39.38	90.94	65.21	79.93	1.51	
	SLE <sup>RF</sup>	33.57	85.36	52.33	64.74	3.84	SLE <sup>RF</sup>	38.82	90.07	49.94	69.26	2.79	
	Translation (ALMA 7B)						Translation (ALMA 7B)						
	BERTSCORE COMET-22 COMETKIWI <sup>RF</sup> xCOMET METRICX METRICX-QE <sup>RF</sup>						BERTSCORE COMET-22 COMETKIWI <sup>RF</sup> xCOMET METRICX METRICX-QE <sup>RF</sup>						
BLEU	90.91	87.12	81.16	72.43	1.15	1.24	BLEU	90.57	86.65	80.49	72.57	1.20	1.35
BERTSCORE	91.41	88.11	82.15	73.59	1.10	1.15	BERTSCORE	90.90	86.52	80.48	71.10	1.31	1.44
COMET-22	90.45	91.18	86.17	76.71	0.61	0.63	COMET-22	89.74	90.28	84.44	73.42	0.74	0.81
COMETKIWI <sup>RF</sup>	90.67	90.56	85.64	81.16	0.51	0.57	COMETKIWI <sup>RF</sup>	89.87	89.53	84.58	78.29	0.58	0.65
xCOMET	90.15	90.03	83.19	86.73	0.70	0.79	xCOMET	90.01	89.18	82.35	83.39	0.79	0.83
METRICX	89.35	89.07	82.00	69.26	0.47	0.69	METRICX	88.99	88.26	81.63	65.32	0.54	0.66
METRICX-QE <sup>RF</sup>	89.58	89.29	83.93	68.78	0.43	0.25	METRICX-QE <sup>RF</sup>	88.98	87.61	81.82	63.47	0.50	0.27

Table 6: Multi-prompt and single prompt performance across metrics. RF = Reference-free reranker.