# BETTER INSTRUCTION-FOLLOWING THROUGH MINI MUM BAYES RISK

Anonymous authors

Paper under double-blind review

### ABSTRACT

General-purpose LLM judges capable of human-level evaluation provide not only a scalable and accurate way of evaluating instruction-following LLMs but also new avenues for supervising and improving their performance. One promising way of leveraging LLM judges for supervision is through Minimum Bayes Risk (MBR) decoding, which uses a reference-based evaluator to select a high-quality output from amongst a set of candidate outputs. In the first part of this work, we explore using MBR decoding as a method for improving the test-time performance of instruction-following LLMs. We find that MBR decoding with reference-based LLM judges substantially improves over greedy decoding, best-of-N decoding with reference-free judges and MBR decoding with lexical and embedding-based metrics on AlpacaEval and MT-Bench. These gains are consistent across LLMs with up to 70B parameters, demonstrating that smaller LLM judges can be used to supervise much larger LLMs. Then, seeking to retain the improvements from MBR decoding while mitigating additional test-time costs, we explore iterative self-training on MBR-decoded outputs. We find that self-training using Direct Preference Optimisation leads to significant performance gains, such that the selftrained models with greedy decoding generally match and sometimes exceed the performance of their base models with MBR decoding.

027 028 029

030

025

026

004

010 011

012

013

014

015

016

017

018

019

021

### 1 INTRODUCTION

Instruction-following large language models (LLMs) (Chung et al., 2022; Wei et al., 2022) have
shown remarkable potential as generalist problem-solvers, prompting extensive efforts to improve
their performance. One task that has seen tremendous progress due to LLMs is the evaluation of
text generation itself. Recent works find that "LLM-as-a-Judge" frameworks (Zheng et al., 2023;
Li et al., 2023; Dubois et al., 2024a) demonstrate strong correlation with human evaluations and
significantly outperform lexical (Lin, 2004b; Papineni et al., 2002) and embedding-based methods
(Zhang et al.; Yuan et al., 2021; Qin et al., 2023) across a wide range of instruction-following tasks.

While the use of LLM judges has largely focused on the evaluation of outputs, LLM judges can also provide a way to supervise the generations of other LLMs. This generally involves using the judge as a *reference-free* evaluator to score candidate outputs produced by the LLM and then selecting the 040 highest-scoring candidate as the final output in what is known as best-of-N (BoN) decoding (Song 041 et al., 2024). However, prior works find that LLM judges, including powerful proprietary LLMs 042 such as GPT-4, significantly underperform when no human-curated reference answer is available 043 (Ye et al., 2024; Zheng et al., 2023). In contrast, reference-based evaluation, where a human-curated 044 reference answer is available, shows significantly higher correlation with human evaluations of out-045 puts. This poses a chicken-and-egg problem: how can we leverage reference-based LLM judges for 046 test time generation if no human references are available? 047

Minimum Bayes Risk (MBR) decoding (Bickel & Doksum, 1977) provides a way of overcoming this
 problem. In place of the inaccessible human references, MBR decoding leverages other candidate
 outputs as *pseudo-references*, and uses the evaluator, also known as the *utility metric*, to conduct
 reference-based evaluation of all candidate outputs against all pseudo-references. The final output
 is then chosen as the candidate output with the highest average score: see Figure 1.

In this work, we explore whether MBR decoding using LLM judges as utility metrics can be used to enhance instruction-following LLMs. We divide our work into two main parts. First, inspired



Figure 1: **Illustration of MBR decoding**. Multiple candidates are first sampled from an LLM. Then, each candidate is evaluated against all other candidates (pseudo-references) using a reference-based evaluator. The pseudo-references are marginalised to produce final scores, and the candidate with the highest score is selected.

067 by the effectiveness of scaling inference time compute (Welleck et al., 2024; Snell et al., 2024), we investigate whether MBR decoding with LLM judges can improve the performance of instruction-068 following LLMs during test time (denoted as "MBR inference") (Section 3). Second, following 069 recent works demonstrating that iterative self-training can improve LLM performance (Xu et al., 2023; Chen et al., 2024; Yuan et al., 2024a; Wu et al., 2024), we examine whether MBR decoding 071 with LLM judges can be used to select high-quality model outputs for use in subsequent iterations 072 of self-training (denoted as "MBR distillation") (Section 4). This both provides a way of training 073 models without access to external labels and also allows us to mitigate the inference-time costs 074 associated with MBR inference. 075

From our MBR inference experiments, we find that MBR decoding with LLM judge Prometheus-076 2-7B (Kim et al., 2024) improves performance by +3.6% on AlpacaEval and +0.28 on MT-Bench 077 on average across five LLMs relative to greedy decoding. Notably, Llama2-7b with Prometheus 078 MBR decoding outperforms Llama2-13b with greedy decoding on MT-Bench, while Prometheus 079 MBR decoding with Llama2-13b outperforms Llama2-70b with greedy decoding on AlpacaEval 2.0. Gains persist even for large 70B models, demonstrating that small LLM judges can supervise 081 larger LLMs through MBR decoding. We also compare MBR decoding against other methods that 082 use LLM judges for supervision. We show that Prometheus MBR decoding is far more effective than 083 MBR decoding with word match-based metrics (e.g. ROUGE) or semantic similarity-based metrics 084 (e.g. BERTScore). Comparing MBR to BoN decoding, we find that MBR decoding consistently 085 outperforms BoN decoding across multiple LLMs and LLM judges, and that the gains from MBR decoding increase as the supervising LLM judge increases in size and ability.

From our MBR distillation experiments, we find that self-training with Direct Preference Optimisation (DPO) (Rafailov et al., 2024) on preference pairs selected using MBR decoding (Yang et al., 2024) with Prometheus-2-7B substantially improves greedy decoding performance. For instance, MBR self-trained Llama2-13b improves by +7.1% on AlpacaEval 2.0 and +0.90 on MT-Bench relative to its baseline SFT counterpart when evaluated using only greedy decoding, far surpassing the corresponding gains from BoN self-training. We also find that MBR self-trained models evaluated with greedy decoding generally match and sometimes exceed the performance of their base models evaluated with MBR decoding, thereby demonstrating that MBR distillation is an effective way of mitigating the inference-time costs of MBR decoding while retaining improved performance.

096 097

098

### 2 BACKGROUND

Language models are *autoregressive* probabilistic models; i.e., the probability of a token  $y_i$  depends on prompt x and all previous tokens in the sequence:

101 102 103  $p(y|x) = \prod_{i=1}^{T} p(y_i|y_{i-1}, \dots, y_1, x).$ (1)

During inference, outputs are typically obtained either using maximum a-posteriori-based decoding methods that attempt to maximise probability, such as greedy decoding  $(y_i = \arg \max_{y_i} p(y_i | y_{< i}, x))$  and beam search (Graves, 2012), or by tokenwise sampling from the distribution  $(y_i \sim p(y_i | y_{< i}, x))$ . Both approaches rely on the model's distribution as indicative of output quality. Alternatively, we can first obtain a hypothesis set  $\mathcal{H}_{hyp}$  comprising  $N_{cand}$  candidate outputs from the model (for example, by sampling multiple times), and then select the final output from  $\mathcal{H}_{hyp}$ based on some external criteria. For example, given some reference-free evaluator u (e.g. an LLM judge), best-of-N (BoN) decoding selects the output  $\hat{y} \in \mathcal{H}_{hyp}$  such that

$$\hat{y} = \underset{y \in \mathcal{H}_{hyp}}{\arg\max u(y)}.$$
(2)

As reference-free estimation of output quality can be a difficult problem, MBR decoding replaces the reference-free evaluator with a reference-based evaluator  $u(y, y^*)$  (e.g. a reference-based LLM judge) that evaluates candidate y relative to a reference  $y^*$ .<sup>1</sup> In the MBR literature, this evaluator is known as a *utility metric* (Freitag et al., 2022; Fernandes et al., 2022; Finkelstein et al., 2024). MBR decoding selects the final output  $\hat{y}$  that maximises *expected utility* under the model distribution:

$$\hat{y} = \underset{y \in \mathcal{H}_{\text{hyp}}}{\arg \max} \qquad \underbrace{\mathbb{E}_{y^* \sim p(y|x)}[u(y, y^*)]}_{\approx \frac{1}{N_{\text{cand}}} \sum_{j=1}^{N_{\text{cand}}} u(y, y^{(j)})}, \tag{3}$$

123 where the expectation is approximated as a Monte Carlo sum using model samples  $y^{(1)}, \ldots, y^{(N_{\text{cand}})} \sim p(y|x)$ . In practice, this amounts to computing the utility of each candidate 124 in  $\mathcal{H}_{hyp}$  using all other candidates as (pseudo-)references, and then selecting the candidate with 125 the highest average utility as the final output<sup>2</sup> - see Appendix I.1 for an algorithmic description of 126 MBR decoding. The MBR-decoded output can therefore be interpreted as being the candidate with 127 the highest "consensus" utility as measured by the utility metric, as it achieves the highest aver-128 age utility when evaluated against all other candidate outputs. It is therefore crucial to choose a 129 reference-based metric that is a good proxy for human preferences as our utility function, as this 130 ensures that a "high-consensus" output corresponds to a "high-quality" output. 131

132

112 113 114

120 121 122

### 3 MBR INFERENCE

133 134 135

136 137

145

161

In this experiment, we investigate using MBR decoding with LLM judge utility metrics to improve instruction-following LLMs at test time.

### 138 3.1 EXPERIMENTAL SETUP

139 140 3.1.1 MODELS AND GENERATION PARAMETERS

We use the chat and instruct variants of the Llama2 (Touvron et al., 2023b) and Llama3 (Dubey et al., 2024) models in this experiment. All models have undergone prior SFT and demonstrate strong instruction-following and conversation abilities. We generate  $N_{\text{cand}} = 30$  candidates using temperature sampling with t = 0.3 for all MBR decoding experiments unless otherwise specified.

146 3.1.2 MBR UTILITY METRICS

LLM judge We choose Prometheus-2-7B (Kim et al., 2024) as our representative reference-based
LLM judge. Prometheus is a specialist judge model finetuned from Mistral-7b (Jiang et al., 2023)
that correlates strongly with human judges and GPT-4. It takes as inputs a task prompt, a scoring
rubric (see Appendix C), a candidate and a reference, and outputs an explanation of its judgement
followed by a score from 1 to 5, which we interpret as a utility score. Crucially, Prometheus can
also act as a reference-free judge by simply omitting the reference from its input. This allows us to
directly compare MBR with BoN decoding using the same LLM judge utility metric.

We compare using LLM judges for MBR decoding with three other classes of utility metrics:

ROUGE ROUGE (Lin, 2004a) is a word-level F-measure designed for measuring summarisation and machine translation performance (Lin, 2004a). We use ROUGE-1 in our main experiments but include results for other ROUGE variants in Appendix A.3.

<sup>159</sup> <sup>1</sup>Certain evaluators (e.g. LLM judges) are "task-aware", and take prompt x as an input when performing evaluation. Such utility metrics can then be written as u(y; x) and  $u(y, y^*; x)$ .

<sup>&</sup>lt;sup>2</sup>The expectation can also be computed over a separate set of model outputs known as the *evidence set* (Eikema & Aziz, 2020; Bertsch et al., 2023). We do not explore this setting in our work.

BERTScore BERTScore is a neural evaluation metric that computes the token-level contextual similarity between a candidate and a reference. Like ROUGE, BERTScore is not task-aware. As our model outputs may be longer than 512 tokens, we use Longformer-large (Beltagy et al., 2020) as our BERTScore model, which supports inputs up to 4094 tokens long.

**Dense embedders** Dense embedders generate contextual embeddings of text passages for use in downstream tasks. One such task is measuring the level of semantic similarity between two text passages (Agirre et al., 2012). This task is directly relevant to MBR decoding, as we can treat pairs of candidates as text passages and their similarity score as the utility. To the best of our knowledge, using dense embedders as a utility metric for MBR decoding has never been explored before. In our work, we use the instruction-following embedder **SFR-Embedder-2\_R** (Meng et al., 2024) as our representative dense embedder. We include results for two other dense embedders in Appendix A.3.

174 3.1.3 BASELINES

In addition to greedy decoding and beam search (BS) with k = 10 beams, we also experiment with LONGEST decoding, where we select the longest candidate from the hypothesis set (as measured in characters) as the final output, and best-of-N (BoN) decoding. We generate  $N_{cand} = 30$  candidates using temperature sampling with t = 0.3 for both longest and BoN decoding. See Appendices A.2 and A.3 for comparison to additional baselines.

180 181 a

182

173

3.1.4 EVALUATION

AlpacaEval 2.0 AlpacaEval (Li et al., 2023) is an LLM-based evaluation metric. It consists of 183 an 805-sample, highly diverse single-turn instruction-following conversational dataset and an as-184 sociated evaluation framework. In AlpacaEval 2.0 (Dubois et al., 2024b), evaluation is conducted 185 by performing head-to-head comparison of candidate answers against GPT-4-generated answers facilitated by a judge LLM. The judge model is prompted to output a single token representing its 187 choice of winner, with the log-probabilities of the token used to compute a weighted win rate. In 188 addition to standard win rates, AlpacaEval 2.0 also provides length-controlled (LC) win rates, which 189 are debiased versions of the standard win rates that control for the length of the outputs. Both the 190 AlpacaEval standard and LC evaluation demonstrate strong correlation with human judgements.

MT-Bench MT-Bench (Zheng et al., 2023) is an 80-sample, two-turn instruction-following conversational dataset. It can be evaluated using either head-to-head comparison or direct assessment with an LLM judge. In the direct assessment setting, the judge LLM is prompted to generate an explanation followed by a score between 1 and 10, with no reference answer used. MT-Bench with GPT-4-judge matches crowdsourced human preferences well, achieving over 80% agreement, which is the same level of agreement between human evaluators (Zheng et al., 2023).

We use GPT-40 (OpenAI et al., 2024) as the LLM judge for both AlpacaEval 2.0 and MT-Bench.
For AlpacaEval, we report LC win rates unless otherwise stated. For MT-Bench, we use direct assessment for all experiments. See Appendix B for further details on our evaluation strategy, and Appendix H for human study findings verifying the alignment of our automatic LLM evaluation results with human judgements.

202 203 204

### 3.2 EXPERIMENTAL RESULTS

	2-7B	2-13B	2-70B	3-8B	3-70B	Avg. $\Delta$
Greedy	14.4	19.0	22.8	34.4	42.7	0
BS	14.8	18.2	21.5	33.9	42.4	-0.50
Longest	10.5	15.2	19.8	29.8	40.4	-3.51
Prometheus BoN	<u>16.4</u>	<u>20.8</u>	<u>25.0</u>	35.5	<u>44.3</u>	<u>1.74</u>
ROUGE MBR	16.2	20.0	24.7	35.4	43.7	1.33
BERTScore MBR	16.2	20.5	24.4	35.7	44.0	1.50
SFR-Embedder MBR	12.1	16.6	22.2	32.5	42.8	-1.42
Prometheus MBR	17.7	23.4	26.2	37.9	46.0	3.62

212 213

Table 1: AlpacaEval 2.0 win rates (%) for various models and decoding strategies, along with the average win rate differences compared to greedy decoding across all models (denoted as **Avg.**  $\Delta$ ). MBR decoding with Prometheus consistently outperforms all baseline methods and other MBR decoding methods.

216		2-7B	2-13B	2-70B	3-8B	3-70B	Avg. $\Delta$
217	Greedy	5.72	5.90	6.50	7.54	8.29	0
218	BS	5.58	5.95	6.49	7.30	8.20	-0.09
219	Longest	5.67	6.03	6.59	7.22	8.22	-0.04
220	Prometheus BoN	5.77	6.08	6.65	<u>7.66</u>	<u>8.42</u>	<u>0.13</u>
220	ROUGE MBR	5.78	6.11	6.68	7.63	8.31	0.11
221	BERTScore MBR	5.68	6.02	6.72	7.52	8.42	0.08
222	SFR-Embedder MBR	5.73	6.04	6.54	7.45	8.33	0.03
223	Prometheus MBR	6.10	6.26	6.79	7.69	8.50	0.28

223 224 225

Table 2: MT-Bench scores for various models and decoding strategies, along with the average score differences compared to greedy decoding across all models (denoted as Avg.  $\Delta$ ). MBR decoding with Prometheus consistently outperforms all baseline methods and other MBR decoding methods.

226 227 228

<sup>229</sup> Our main experimental results are documented in Tables 1 and 2.

230 **Prometheus MBR decoding provides significant and consistent gains** The gains associated 231 with Prometheus MBR decoding are significantly larger than those associated with other utility 232 metrics, yielding an average improvement of +3.6% on AlpacaEval 2.0 and +0.28 on MT-Bench. 233 For comparison, the performance gap between Llama2-7b and Llama2-13b with greedy decoding is 234 +4.6% on AlpacaEval 2.0 and +0.18 on MT-Bench, while the corresponding gap between Llama2-235 13b and Llama2-70b is +3.8% and +0.60. Notably, Llama2-7b with Prometheus MBR decoding 236 outperforms Llama2-13b with greedy decoding on MT-Bench, while Prometheus MBR decoding with Llama2-13b outperforms Llama2-70b - a model over five times bigger - with greedy decoding 237 on AlpacaEval 2.0. We also find that Prometheus MBR decoding yields larger gains than Prometheus 238 BoN decoding; we explore this further in Section 3.3.1. 239

We also highlight that the performance gains associated with Prometheus MBR decoding are significant across models of all sizes, even for much larger models such as Llama3-70b. This scaling
property suggests that small judge models can still be used to supervise much larger models.

ROUGE and BERTScore MBR decoding provide small but consistent gains ROUGE and BERTScore MBR decoding improve average performance relative to greedy decoding by +1.3% and +1.5% on AlpacaEval 2.0 and by +0.11 and +0.08 on MT-Bench respectively. This benefit is present for all models. This improvement suggests that selecting outputs without awareness of the task and using only word- or token-level measures of consistency can still yield meaningful improvements even in the instruction-following setting.

249 SFR-Embedder MBR decoding fails to yield consistent gains SFR-Embedder MBR decoding 250 reduces performance relative to greedy decoding by -1.4% on AlpacaEval 2.0 while improving per-251 formance by +0.03 on MT-Bench on average. We hypothesise that embedder models, which are 252 trained to distinguish at a high level between text passages, cannot to detect nuanced differences 253 between semantically-similar outputs. We also note that embedder MBR decoding generally se-254 lects for longer outputs, which may explain the discrepancy between its performance on AlpacaEval 255 2.0 (which is length-controlled) and MT-Bench. See Appendix A.4 for analysis on the generation lengths of various decoding strategies. 256

Beam search and LONGEST decoding degrade performance Beam search and LONGEST decoding reduce performance relative to greedy decoding by -0.5% and -3.5% on AlpacaEval 2.0 and -0.09 and -0.04 on MT-Bench respectively. The poor performance of beam search further under-scores the idea that optimising for output probability alone is not enough to improve output quality.

261 262

263

3.3 ANALYSIS OF PROMETHEUS MBR DECODING

264 265 Given the promise shown by Prometheus MBR decoding, we conduct additional experiments to better understand its properties.

**Prometheus MBR decoding performance vs.** t and  $N_{cand}$  We plot AlpacaEval 2.0 win rates as a function of  $N_{cand}$  and t in Figure 2 for Llama2-70b with Prometheus MBR and BoN decoding. We find that performance initially increases with  $N_{cand}$  but plateaus at around  $N_{cand} = 20$ , suggesting that expanding the size of the hypothesis set beyond this yields little benefit. Performance



Figure 2: AlpacaEval 2.0 win rates (%) for Llama2-70b with varying hypothesis set size  $N_{\text{cand}}$  (left) and generation temperature t (right) values for Prometheus MBR and BoN decoding. Performance for both methods initially increases with  $N_{\text{cand}}$  and plateaus at around  $N_{\text{cand}} = 20$ . Performance also initially increases with t, but drops rapidly after t = 1.0.

also initially increases with t, highlighting the benefits of increased candidate diversity, although it rapidly degrades at high temperatures as the individual candidate outputs decline in quality.



Figure 3: Difference in AlpacaEval 2.0 win rates (%) between Prometheus MBR decoding and greedy decoding averaged over all five LLMs and broken down by question category. A positive value indicates that MBR decoding outperforms greedy decoding on the given category. Orange bars represent the standard error. We find that Prometheus MBR decoding improves performance across a wide range of question categories.

**Prometheus MBR decoding performance by question category** We classified questions from the AlpacaEval dataset into one of ten categories using GPT-40 (see Appendix D for details), and then computed the differences in win rates between Prometheus MBR decoding and greedy decoding by question category, averaged over all five LLMs. We find that MBR decoding improves output quality across most question categories. These include reasoning-based categories such as coding and mathematical reasoning, although the largest improvements are seen across writing-based categories such as technical, recommendations and creative writing. We hypothesise that this discrepancy arises due to (1) the higher bar for correctness associated with reasoning tasks which limits the number of good answers that can be found amongst candidate outputs; and (2) limitations of existing utility functions, which may struggle to handle difficult reasoning tasks.

# 316 3.3.1 FURTHER COMPARISONS WITH BEST-OF-N DECODING

As BoN decoding can also leverage LLM judges as a utility metric, we conduct additional experiments to compare its performance against MBR decoding. We compare BoN and MBR decoding for five different LLM judges on MT-Bench and report the results in Table 3. In addition to
Prometheus-2-7B, we also evaluate its larger sibling Prometheus-2-8x7B, as well as JudgeLM-7b
and JudgeLM-33b (Zhu et al., 2023c), which are two judge models finetuned from LLaMA models
(Touvron et al., 2023a). We also assess Llama3-8b-Instruct and Llama3-70b-Instruct as a zeroshot judge for MBR decoding (see Appendix E for our prompts). All of our chosen judges can act

324		2-7B	2-13B	2-70B	3-8B	3-70B	Avg. $\Delta$
325	Greedy	5.72	5.90	6.50	7.54	8.29	0
326	Prometheus-2-7B-BoN	5.77	6.08	6.65	7.66	8.42	0.13
327	Prometheus-2-7B-MBR	6.10	6.26	6.79	7.69	8.50	0.28
308	Prometheus-2-8x7B-BoN	6.01	6.17	6.80	7.75	8.41	0.24
320	Prometheus-2-8x7B-MBR	6.26	<u>6.32</u>	6.87	7.79	8.64	<u>0.39</u>
329	JudgeLM-7b-BoN	5.63	5.95	6.69	7.37	8.26	-0.01
330	JudgeLM-7b-MBR	6.00	6.11	6.79	7.69	8.44	0.22
331	JudgeLM-33b-BoN	5.68	6.03	6.58	7.37	8.35	0.01
332	JudgeLM-33b-MBR	5.94	6.27	<u>6.88</u>	7.92	8.50	0.31
333	Llama3-8b-Instruct-BoN	5.83	6.05	6.61	7.60	8.38	0.10
224	Llama3-8b-Instruct-MBR	5.96	6.28	6.84	7.80	8.47	0.28
334	Llama3-70b-Instruct-BoN	5.77	6.16	6.57	7.39	8.35	0.06
335	Llama3-70b-Instruct-MBR	6.22	6.43	6.94	7.87	8.52	0.41
000			1				

Table 3: MT-Bench scores for BoN and MBR decoding with various judge LLMs as utility metrics, along with the average score differences compared to greedy decoding across all models (denoted **Avg.**  $\Delta$ ). MBR decoding consistently outperforms BoN decoding across all comparable utility metrics.

as both reference-free and reference-based judges, allowing us to compare MBR and BoN decoding
 on a level playing field.<sup>3</sup>

343 We find that MBR decoding consistently outperforms BoN decoding across all selected judge mod-344 els. This difference is especially large for the JudgeLM models and for Llama3-70b-Instruct, where 345 BoN fails to significantly improve on greedy decoding. One explanation for this discrepancy is that our LLM judges are insufficiently good at reference-free evaluation for BoN decoding to be 346 effective. This idea is supported by prior studies comparing reference-free and reference-based 347 evaluation, which consistently show that reference-free methods tend to underperform, even when 348 using strong judge models like GPT-4 (Ye et al., 2024; Kim et al., 2024; Zheng et al., 2023). Another 349 explanation is that MBR decoding provides a smoothing effect that arises from our use of expected 350 utility in place of utility point estimates for output selection, tying back to our hypothesis that select-351 ing "high-consensus" outputs yields significant benefit. This averaging process reduces the impact 352 of individual mistakes made by the imperfect LLM judge, thereby providing for a more stable and 353 reliable measure of quality. We leave further exploration of these ideas to future work. 354

Notably, in Table 3, MBR performance improves by scaling the size of the LLM judge, with
Prometheus-2-8x7B outperforming Prometheus-2-7B, JudgeLM-33b outperforming JudgeLM-7b,
and Llama3-70b-Instruct outperforming Llama3-8b-Instruct. This suggests that improving the LLM
judge utility metric directly improves MBR decoding performance and that MBR decoding will
benefit as newer and better LLM judges are developed.

360 361

362

363

364

366

340

### 4 MBR DISTILLATION

Our results so far demonstrate the potential of MBR decoding to significantly improve the test-time performance of instruction-following models, but this comes at the cost of substantial inference-time compute costs due to the linear cost for generating  $N_{cand}$  candidate outputs, and the quadratic cost for computing utility across these candidates. To mitigate this, we explore distilling MBR-decoded outputs back into the model itself and aim to obtain MBR decoding-level (or better) performance without needing to perform MBR decoding at test time.

367368369370

371

372

- 4.1 EXPERIMENTAL SETUP
- 4.1.1 TRAINING AN SFT MODEL

We start by performing SFT on the base Llama2-7b and Llama2-13b models. This is necessary to instil instruction-following behaviour in the models so that they can be used to generate instructionfollowing self-training data. We choose not to use the official chat variants of these models as

this in the main text and report our findings in Appendix F instead.

 <sup>&</sup>lt;sup>3</sup>Because sequence-classifier reward models (Stiennon et al., 2022) do not support reference-based evalua tion, it is not possible to fairly compare BoN decoding with these methods to MBR. We therefore do not discuss

we wish to retain full control over the training procedure and avoid inheriting any biases introduced
through prior finetuning and alignment. We use 3000 randomly-drawn samples from UltraChat-200k
(Ding et al., 2023) for SFT. UltraChat is a highly diverse conversational instruction-following dataset
created using GPT-3.5-Turbo. Each sample consists of multi-turn prompts and responses, although
we only take the first turn of each sample in order to simplify experimentation. We designate our
SFT models as *sft*.

384 385

### 4.1.2 ITERATIVE DPO ON MBR-DECODED OUTPUTS

386 Having obtained SFT models, we now conduct DPO to improve the models on their own MBR-387 decoded outputs, an approach first proposed by Yang et al. (2024) for improving machine translation. 388 We start by randomly drawing a further 3000 prompts from UltraChat (excluding the samples that 389 have already been selected for SFT). Next, we generate  $N_{\rm cand} = 12$  candidate outputs from our sft 390 models using these prompts. Following Yang et al. (2024), we then score the candidate outputs using 391 a utility metric and form preference pairs from the highest-scoring and lowest-scoring outputs. This 392 preference pair dataset is then used for DPO on the sft models, yielding dpo-MBR-1. We extend 393 upon Yang et al. (2024) by iteratively repeating this process twice more, each time using the latest dpo models as the base model paired with a fresh set of 3000 prompts, yielding the models dpo-394 MBR-2 and dpo-MBR-3. See Appendix K for a summary of our SFT and DPO hyperparameters, 395 Appendix G.4 for experimental results from using another preference pair selection strategy, and 396 Appendix I.2 for mathematical and algorithmic overviews of MBR distillation. 397

398 399

414 415

422

423

424

425

426

427

428 429

430

431

### 4.1.3 UTILITY METRICS AND EVALUATION

400 We use Prometheus-2-7B as our utility metric, as this yielded the most promising results in our 401 earlier experiments (Section 3.2), although we also try MBR self-training with ROUGE as the utility 402 metric (Appendix G.3). We compare our *dpo* models with greedy decoding against the *sft* models 403 with greedy decoding, beam search and MBR decoding. For MBR decoding, we use  $N_{\text{cand}} = 30$ 404 and t = 0.3 with Prometheus-2-7B as the utility metric. We also baseline against models trained 405 with SFT on 12000 UltraChat samples that we call sft-full. Finally, we experiment with BoN selftraining, again using Prometheus-2-7B as the utility metric and following the same procedure as for 406 MBR self-training, which yields the models dpo-BoN-1, dpo-BoN-2 and dpo-BoN-3. 407

We evaluate our trained models using greedy decoding on AlpacaEval 2.0, once again reporting length-controlled win rates vs. GPT-4, and MT-Bench. As we only train our models to engage in single-turn conversations we evaluate only on the first turn of MT-Bench. We report additional evaluation results in the Appendix, including head-to-head results between various self-trained models and *sft* with greedy decoding (Appendix G.1), evaluation results on a selection of popular NLP benchmarks (Appendix G.2), and human study results (Appendix H).

### 4.2 Results

	AlpacaEval 2.0		MT-Bench		
	7B	13B	7B	13B	
sft w. Greedy	5.18	8.24	5.43	5.85	
sft w. MBR	9.99	13.6	5.78	6.31	
<i>sft</i> -full	6.35	9.40	5.55	6.26	
dpo-1-BoN	5.78	10.3	5.78	6.08	
dpo-2-BoN	6.22	11.2	5.91	6.41	
dpo-3-BoN	6.40	12.8	5.88	6.56	
dpo-1-MBR	5.68	10.8	5.78	6.48	
dpo-2-MBR	7.22	13.9	6.11	6.73	
dpo-3-MBR	8.86	15.3	6.14	6.75	

	AlpacaEval 2.0	MT-Bench
sft-1-MBR	5.52	5.48
sft-2-MBR	6.75	5.43
sft-3-MBR	6.48	5.51

Table 4: (Left) AlpacaEval 2.0 win rates (%) and MT-Bench scores for models self-trained using DPO. After three rounds of training, the self-trained models consistently outperform their BoN counterparts and SFT baselines. (Top) AlpacaEval 2.0 win rates (%) and MT-Bench scores for models self-trained using SFT. Self-training with SFT yields substantially worse results than self-training with DPO.

**DPO self-training significantly improves model performance** We report the results of our self-training experiment in the left subtable of Table 4. We find that three rounds of MBR self-training with DPO significantly improves model performance, with the 7B *dpo*-3-MBR model outperforming the 13B *sft* model with greedy decoding on both AlpacaEval 2.0 and MT-Bench. The improvements

445

446

457

458

460



Figure 4: Differences in AlpacaEval 2.0 win rates (%) between dpo-3-MBR models and their respective sft with greedy decoding baselines on different question categories. The largest improvements are seen in open-ended writing tasks, with less improvement on reasoning-focussed tasks (e.g. mathematical reasoning and coding).



456 Figure 5: (Left) AlpacaEval 2.0 win rate (%) vs. sft with greedy decoding against generation throughput. dpo-3-MBR with greedy decoding matches the performance of sft with MBR decoding, but with significantly higher throughput. (Right) Average generation and decoding times on the AlpacaEval dataset. The decoding step in MBR decoding takes disproportionately long. This problem can mitigated through MBR self-training. 459

461 saturate by the third round of self-training as measured on MT-Bench (6.11 vs. 6.14 (7B) and 6.73 462 vs. 6.75 (13B) in Table 4), although there appears to be room for further improvement on AlpacaEval 463 2.0 (7.22 vs. 8.86 (7B) and 13.9 vs. 15.3 (13B) in Table 4). Both dpo-3-MBR models outperform 464 their sft-full counterparts, which suggests that training on MBR-decoded outputs is more beneficial 465 than SFT on a larger split of UltraChat. The dpo-3-MBR models also generally outperform sft with MBR decoding, and this is especially prominent for MT-Bench, which suggests that DPO on MBR-466 decoded outputs enables models to recover and then exceed their MBR-decoding performances. We 467 find that DPO on BoN-decoded outputs also improves model performance, although less so than 468 DPO with MBR-decoded outputs. We attribute this to the relative strength of MBR decoding. 469

SFT self-training yields smaller gains than DPO self-training We experiment with iterative 470 SFT self-training, using the 7B sft model. We document our results in the right subtable of Table 4. 471 We use the same sets of prompts as for DPO and select as our SFT labels the highest-scoring sample 472 as determined by MBR, following Finkelstein et al. (2024). As before, we conduct three rounds of 473 iterative training, yielding sft-1-MBR, sft-2-MBR and sft-3-MBR. We find that SFT training yields 474 significantly less improvement than DPO. This indicates that MBR self-training benefits most from 475 preference learning, where the model learns to contrast its highest- and lowest-quality outputs. 476

**DPO self-trained model performance by question category** We repeat the analysis on perfor-477 mance by question category for dpo-3-MBR in Figure 4. Self-training improves performance on 478 almost all question categories, with generally larger improvement on writing-based categories and 479 smaller improvement on reasoning-based categories. We attribute this difference to the writing-480 skewed distribution of question categories in our UltraChat training data (see Appendix G.5). 481

482 **Analysis of compute costs** We illustrate the savings on compute introduced by self-training in Figure 5. We perform inference with various 7B models and decoding strategies on AlpacaEval 2.0, 483 using 2xA100 GPUs and vLLM (Kwon et al., 2023) as our inference engine. We use a generation 484 batch size of 1, a LLM judge utility calculation batch size of 32, and  $N_{\text{cand}} = 12$ . We find that MBR decoding imposes significant overhead, largely due to the quadratic  $\mathcal{O}(N_{\text{cand}}^2)$  cost incurred during 485

486 the utility calculation step. This overhead is removed through MBR self-training, which nonetheless 487 retains performance gains. Note that *dpo-3-MBR* generates longer outputs than *sft*, which explains 488 why its average generation time as seen in the right-hand plot of Figure 5 is higher. 489

#### 490 5 **RELATED WORK**

491

492 **MBR decoding** MBR decoding has been explored in the context of machine translation using a 493 variety of translation metrics such as COMET (Rei et al., 2020) and BLEURT (Sellam et al., 2020), 494 with promising results (Freitag et al., 2022; 2023; Farinhas et al., 2023; Stanojević & Sima'an, 495 2014). Prior works (Bertsch et al., 2023; Jinnai et al., 2023) also study MBR decoding for sum-496 marisation, using ROUGE and BERTScore as metrics. Suzgun et al. (2022) apply MBR decoding 497 to several tasks, including summarisation, machine translation and three different BIG-Bench tasks 498 (Srivastava et al., 2023). None of these works explore the use of MBR decoding in the more open-499 ended instruction-following domain, nor do they consider using LLM judges as utility metrics.

500 **LLM judges** Based on the strong instruction-following capabilities of LLMs, recent works ex-501 plore prompting LLMs to judge responses from other LLMs (Li et al., 2023; Zheng et al., 2023). 502 Follow-up works suggest that training on the evaluation traces of strong models may equip smaller 503 models with strong evaluation capabilities (Kim et al., 2023; 2024; Zhu et al., 2023b; Vu et al., 2024; 504 Wang et al., 2024b). These works focus on training LLMs to produce scoring decisions matching 505 those of humans. In our work, instead of viewing evaluation as an end goal, we explore utilising the evaluation capabilities of LLM judges as supervision to improve instruction-following LLMs. 506

507 Inference-time algorithms Many inference-time algorithms generate candidate outputs and se-508 lect a final output based on external criteria. In addition to MBR and BoN decoding, examples in-509 clude Self-Consistency (Wang et al., 2023), which selects the most self-consistent answers through 510 marginalisation of chain-of-thought reasoning paths and Universal Self-Consistency (USC) (Chen 511 et al., 2023), where the LLM is used to self-select consistent chain-of-thought reasoning paths from amongst many reasoning paths. Kuhn et al. (2023) propose an MBR-esque algorithm that uses 512 dense embedders and clustering to measure semantic uncertainty. Other inference-time algorithms 513 prompt the LLM to perform additional inference steps in a structured manner. Examples include 514 Tree-of-Thoughts (Yao et al., 2023) and Graph-of-Thoughts (Besta et al., 2024), as well as recursive 515 improvement strategies such as Self-Refine (Madaan et al., 2023) and Reflexion (Shinn et al., 2023). 516

517 Self-training Self-training is a promising avenue for model improvement as it enables training without labelled data. Gulcehre et al. (2023) introduce an algorithm that generates samples from a policy 518 and then updates the policy using offline RL. Yuan et al. (2024b) train models to score their own 519 outputs, and then use these scores to create preference datasets which for distillation. Huang et al. 520 (2022) train models on their own highest-confidence outputs as determined by majority voting. Self-521 training on MBR-decoded outputs has also been explored for machine translation. Finkelstein et al. 522 (2024) train models with SFT on their own MBR and quality estimation outputs for machine trans-523 lation and demonstrate that this yields improvements over baseline models. Wang et al. (2024a) use 524 MBR to generate targets for sequence-level distillation, again for machine translation. Yang et al. 525 (2024) are the first to use DPO to upweight the model's own MBR generations, allowing them to 526 recover much of their original MBR performances on translation using only greedy decoding.

527 528

529

### CONCLUSION 6

530 In this work, we investigate using LLM judges to supervise other LLMs on instruction-following 531 tasks through MBR decoding, and find that this yields significant and consistent improvements to 532 model performance relative to greedy decoding, beam search and BoN decoding. These benefits 533 persist across a wide range of question categories and are also consistent across models of various 534 sizes, demonstrating that small LLM judges can be used to improve much larger LLMs at inference time. To mitigate the significant inference-time costs associated with MBR decoding, we also ex-536 plore iterative self-training on MBR-decoded outputs. We find that MBR self-training using DPO, 537 but not SFT, enables models to recover and even exceed their base MBR decoding performance using only greedy decoding. We hope our work further highlights the potential of using LLM judges 538 for supervision and inspires future research into MBR decoding beyond its traditional domains and applications, particularly through the development of new utility metrics.

### 540 **Reproducibility Statement** 541

542 In an effort to make our work reproducible, we document all prompts (Appendices E and B), as well 543 as training and inference hyperparameters (Appendix K) used throughout our experiments. We also 544 include version information for all API-based LLMs (Appendix B), and choose to use open-source models (the Llama2, Llama3, Prometheus-2 and JudgeLM families) where possible.

547 References 548

546

567

568

569

570

576

577

- 549 Eneko Agirre, Daniel Cer, Mona Diab, and Aitor Gonzalez-Agirre. SemEval-2012 task 6: A pilot on semantic textual similarity. In Eneko Agirre, Johan Bos, Mona Diab, Suresh Manandhar, Yuval 550 Marton, and Deniz Yuret (eds.), \*SEM 2012: The First Joint Conference on Lexical and Com-551 putational Semantics - Volume 1: Proceedings of the main conference and the shared task, and 552 Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 553 2012), pp. 385–393, Montréal, Canada, 7-8 June 2012. Association for Computational Linguis-554 tics. URL https://aclanthology.org/S12-1051. 555
- 556 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020. URL https://arxiv.org/abs/2004.05150.
- Amanda Bertsch, Alex Xie, Graham Neubig, and Matthew R. Gormley. It's mbr all the way down: 559 Modern generation techniques through the lens of minimum bayes risk, 2023. 560
- 561 Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gi-562 aninazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten 563 Hoefler. Graph of thoughts: Solving elaborate problems with large language models. Proceedings of the AAAI Conference on Artificial Intelligence, 38(16):17682–17690, March 2024. 564 ISSN 2159-5399. doi: 10.1609/aaai.v38i16.29720. URL http://dx.doi.org/10.1609/ 565 aaai.v38i16.29720. 566
  - P.J. Bickel and K.A. Doksum. Mathematical Statistics: Basic Ideas and Selected Topics. Prentice Hall, 1977. ISBN 9780135641477. URL https://books.google.com.sg/books?id= ucMfAQAAIAAJ.
- Xinyun Chen, Renat Aksitov, Uri Alon, Jie Ren, Kefan Xiao, Pengcheng Yin, Sushant Prakash, 571 Charles Sutton, Xuezhi Wang, and Denny Zhou. Universal self-consistency for large language 572 model generation, 2023. 573
- 574 Zixiang Chen, Yihe Deng, Huizhuo Yuan, Kaixuan Ji, and Quanquan Gu. Self-play fine-tuning 575 converts weak language models to strong language models, 2024. URL https://arxiv. org/abs/2401.01335.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan 578 Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, Albert Webson, Shixiang Shane Gu, 579 Zhuyun Dai, Mirac Suzgun, Xinyun Chen, Aakanksha Chowdhery, Alex Castro-Ros, Marie Pel-580 lat, Kevin Robinson, Dasha Valter, Sharan Narang, Gaurav Mishra, Adams Yu, Vincent Zhao, 581 Yanping Huang, Andrew Dai, Hongkun Yu, Slav Petrov, Ed H. Chi, Jeff Dean, Jacob Devlin, 582 Adam Roberts, Denny Zhou, Quoc V. Le, and Jason Wei. Scaling instruction-finetuned language 583 models, 2022. URL https://arxiv.org/abs/2210.11416. 584
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and 585 Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge, 586 2018. URL https://arxiv.org/abs/1803.05457.
- 588 Ning Ding, Yulin Chen, Bokai Xu, Yujia Qin, Zhi Zheng, Shengding Hu, Zhiyuan Liu, Maosong 589 Sun, and Bowen Zhou. Enhancing chat language models by scaling high-quality instructional 590 conversations, 2023. URL https://arxiv.org/abs/2305.14233.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 592 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark,

594 Arun Rao, Aston Zhang, Aurelien Rodriguez, Austen Gregerson, Ava Spataru, Baptiste Roziere, 595 Bethany Biron, Binh Tang, Bobbie Chern, Charlotte Caucheteux, Chaya Nayak, Chloe Bi, Chris 596 Marra, Chris McConnell, Christian Keller, Christophe Touret, Chunyang Wu, Corinne Wong, 597 Cristian Canton Ferrer, Cyrus Nikolaidis, Damien Allonsius, Daniel Song, Danielle Pintz, Danny 598 Livshits, David Esiobu, Dhruv Choudhary, Dhruv Mahajan, Diego Garcia-Olano, Diego Perino, Dieuwke Hupkes, Egor Lakomkin, Ehab AlBadawy, Elina Lobanova, Emily Dinan, Eric Michael Smith, Filip Radenovic, Frank Zhang, Gabriel Synnaeve, Gabrielle Lee, Georgia Lewis Ander-600 son, Graeme Nail, Gregoire Mialon, Guan Pang, Guillem Cucurell, Hailey Nguyen, Hannah 601 Korevaar, Hu Xu, Hugo Touvron, Iliyan Zarov, Imanol Arrieta Ibarra, Isabel Kloumann, Ishan 602 Misra, Ivan Evtimov, Jade Copet, Jaewon Lee, Jan Geffert, Jana Vranes, Jason Park, Jay Ma-603 hadeokar, Jeet Shah, Jelmer van der Linde, Jennifer Billock, Jenny Hong, Jenya Lee, Jeremy 604 Fu, Jianfeng Chi, Jianyu Huang, Jiawen Liu, Jie Wang, Jiecao Yu, Joanna Bitton, Joe Spisak, 605 Jongsoo Park, Joseph Rocca, Joshua Johnstun, Joshua Saxe, Junteng Jia, Kalyan Vasuden Al-606 wala, Kartikeya Upasani, Kate Plawiak, Ke Li, Kenneth Heafield, Kevin Stone, Khalid El-Arini, 607 Krithika Iyer, Kshitiz Malik, Kuenley Chiu, Kunal Bhalla, Lauren Rantala-Yeary, Laurens van der 608 Maaten, Lawrence Chen, Liang Tan, Liz Jenkins, Louis Martin, Lovish Madaan, Lubo Malo, Lukas Blecher, Lukas Landzaat, Luke de Oliveira, Madeline Muzzi, Mahesh Pasupuleti, Mannat Singh, Manohar Paluri, Marcin Kardas, Mathew Oldham, Mathieu Rita, Maya Pavlova, 610 Melanie Kambadur, Mike Lewis, Min Si, Mitesh Kumar Singh, Mona Hassan, Naman Goyal, 611 Narjes Torabi, Nikolay Bashlykov, Nikolay Bogoychev, Niladri Chatterji, Olivier Duchenne, Onur 612 Çelebi, Patrick Alrassy, Pengchuan Zhang, Pengwei Li, Petar Vasic, Peter Weng, Prajjwal Bhar-613 gava, Pratik Dubal, Praveen Krishnan, Punit Singh Koura, Puxin Xu, Qing He, Qingxiao Dong, 614 Ragavan Srinivasan, Raj Ganapathy, Ramon Calderer, Ricardo Silveira Cabral, Robert Stojnic, 615 Roberta Raileanu, Rohit Girdhar, Rohit Patel, Romain Sauvestre, Ronnie Polidoro, Roshan Sum-616 baly, Ross Taylor, Ruan Silva, Rui Hou, Rui Wang, Saghar Hosseini, Sahana Chennabasappa, 617 Sanjay Singh, Sean Bell, Seohyun Sonia Kim, Sergey Edunov, Shaoliang Nie, Sharan Narang, 618 Sharath Raparthy, Sheng Shen, Shengye Wan, Shruti Bhosale, Shun Zhang, Simon Vandenhende, 619 Soumya Batra, Spencer Whitman, Sten Sootla, Stephane Collot, Suchin Gururangan, Sydney 620 Borodinsky, Tamar Herman, Tara Fowler, Tarek Sheasha, Thomas Georgiou, Thomas Scialom, Tobias Speckbacher, Todor Mihaylov, Tong Xiao, Ujjwal Karn, Vedanuj Goswami, Vibhor Gupta, 621 Vignesh Ramanathan, Viktor Kerkez, Vincent Gonguet, Virginie Do, Vish Vogeti, Vladan Petro-622 vic, Weiwei Chu, Wenhan Xiong, Wenyin Fu, Whitney Meers, Xavier Martinet, Xiaodong Wang, 623 Xiaoqing Ellen Tan, Xinfeng Xie, Xuchao Jia, Xuewei Wang, Yaelle Goldschlag, Yashesh Gaur, 624 Yasmine Babaei, Yi Wen, Yiwen Song, Yuchen Zhang, Yue Li, Yuning Mao, Zacharie Delpierre 625 Coudert, Zheng Yan, Zhengxing Chen, Zoe Papakipos, Aaditya Singh, Aaron Grattafiori, Abha 626 Jain, Adam Kelsey, Adam Shajnfeld, Adithya Gangidi, Adolfo Victoria, Ahuva Goldstand, Ajay 627 Menon, Ajay Sharma, Alex Boesenberg, Alex Vaughan, Alexei Baevski, Allie Feinstein, Amanda 628 Kallet, Amit Sangani, Anam Yunus, Andrei Lupu, Andres Alvarado, Andrew Caples, Andrew 629 Gu, Andrew Ho, Andrew Poulton, Andrew Ryan, Ankit Ramchandani, Annie Franco, Aparajita 630 Saraf, Arkabandhu Chowdhury, Ashley Gabriel, Ashwin Bharambe, Assaf Eisenman, Azadeh 631 Yazdan, Beau James, Ben Maurer, Benjamin Leonhardi, Bernie Huang, Beth Loyd, Beto De Paola, Bhargavi Paranjape, Bing Liu, Bo Wu, Boyu Ni, Braden Hancock, Bram Wasti, Bran-632 don Spence, Brani Stojkovic, Brian Gamido, Britt Montalvo, Carl Parker, Carly Burton, Catalina 633 Mejia, Changhan Wang, Changkyu Kim, Chao Zhou, Chester Hu, Ching-Hsiang Chu, Chris Cai, 634 Chris Tindal, Christoph Feichtenhofer, Damon Civin, Dana Beaty, Daniel Kreymer, Daniel Li, 635 Danny Wyatt, David Adkins, David Xu, Davide Testuggine, Delia David, Devi Parikh, Diana 636 Liskovich, Didem Foss, Dingkang Wang, Duc Le, Dustin Holland, Edward Dowling, Eissa Jamil, 637 Elaine Montgomery, Eleonora Presani, Emily Hahn, Emily Wood, Erik Brinkman, Esteban Ar-638 caute, Evan Dunbar, Evan Smothers, Fei Sun, Felix Kreuk, Feng Tian, Firat Ozgenel, Francesco 639 Caggioni, Francisco Guzmán, Frank Kanayet, Frank Seide, Gabriela Medina Florez, Gabriella 640 Schwarz, Gada Badeer, Georgia Swee, Gil Halpern, Govind Thattai, Grant Herman, Grigory 641 Sizov, Guangyi, Zhang, Guna Lakshminarayanan, Hamid Shojanazeri, Han Zou, Hannah Wang, Hanwen Zha, Haroun Habeeb, Harrison Rudolph, Helen Suk, Henry Aspegren, Hunter Gold-642 man, Ibrahim Damlaj, Igor Molybog, Igor Tufanov, Irina-Elena Veliche, Itai Gat, Jake Weissman, James Geboski, James Kohli, Japhet Asher, Jean-Baptiste Gaya, Jeff Marcus, Jeff Tang, Jennifer 644 Chan, Jenny Zhen, Jeremy Reizenstein, Jeremy Teboul, Jessica Zhong, Jian Jin, Jingyi Yang, Joe 645 Cummings, Jon Carvill, Jon Shepard, Jonathan McPhie, Jonathan Torres, Josh Ginsburg, Junjie 646 Wang, Kai Wu, Kam Hou U, Karan Saxena, Karthik Prasad, Kartikay Khandelwal, Katayoun 647 Zand, Kathy Matosich, Kaushik Veeraraghavan, Kelly Michelena, Keqian Li, Kun Huang, Kunal 674

683

686

687

688

689

690

691

699

648 Chawla, Kushal Lakhotia, Kyle Huang, Lailin Chen, Lakshya Garg, Lavender A, Leandro Silva, 649 Lee Bell, Lei Zhang, Liangpeng Guo, Licheng Yu, Liron Moshkovich, Luca Wehrstedt, Madian 650 Khabsa, Manav Avalani, Manish Bhatt, Maria Tsimpoukelli, Martynas Mankus, Matan Hasson, 651 Matthew Lennie, Matthias Reso, Maxim Groshev, Maxim Naumov, Maya Lathi, Meghan Ke-652 neally, Michael L. Seltzer, Michal Valko, Michelle Restrepo, Mihir Patel, Mik Vyatskov, Mikayel Samvelyan, Mike Clark, Mike Macey, Mike Wang, Miquel Jubert Hermoso, Mo Metanat, Mo-653 hammad Rastegari, Munish Bansal, Nandhini Santhanam, Natascha Parks, Natasha White, Navy-654 ata Bawa, Nayan Singhal, Nick Egebo, Nicolas Usunier, Nikolay Pavlovich Laptev, Ning Dong, 655 Ning Zhang, Norman Cheng, Oleg Chernoguz, Olivia Hart, Omkar Salpekar, Ozlem Kalinli, 656 Parkin Kent, Parth Parekh, Paul Saab, Pavan Balaji, Pedro Rittner, Philip Bontrager, Pierre Roux, 657 Piotr Dollar, Polina Zvyagina, Prashant Ratanchandani, Pritish Yuvraj, Qian Liang, Rachad Alao, 658 Rachel Rodriguez, Rafi Ayub, Raghotham Murthy, Raghu Nayani, Rahul Mitra, Raymond Li, 659 Rebekkah Hogan, Robin Battey, Rocky Wang, Rohan Maheswari, Russ Howes, Ruty Rinott, 660 Sai Jayesh Bondu, Samyak Datta, Sara Chugh, Sara Hunt, Sargun Dhillon, Sasha Sidorov, Sa-661 tadru Pan, Saurabh Verma, Seiji Yamamoto, Sharadh Ramaswamy, Shaun Lindsay, Shaun Lind-662 say, Sheng Feng, Shenghao Lin, Shengxin Cindy Zha, Shiva Shankar, Shuqiang Zhang, Shuqiang Zhang, Sinong Wang, Sneha Agarwal, Soji Sajuyigbe, Soumith Chintala, Stephanie Max, Stephen 663 Chen, Steve Kehoe, Steve Satterfield, Sudarshan Govindaprasad, Sumit Gupta, Sungmin Cho, Sunny Virk, Suraj Subramanian, Sy Choudhury, Sydney Goldman, Tal Remez, Tamar Glaser, 665 Tamara Best, Thilo Kohler, Thomas Robinson, Tianhe Li, Tianjun Zhang, Tim Matthews, Tim-666 othy Chou, Tzook Shaked, Varun Vontimitta, Victoria Ajayi, Victoria Montanez, Vijai Mohan, 667 Vinay Satish Kumar, Vishal Mangla, Vítor Albiero, Vlad Ionescu, Vlad Poenaru, Vlad Tiberiu 668 Mihailescu, Vladimir Ivanov, Wei Li, Wenchen Wang, Wenwen Jiang, Wes Bouaziz, Will Con-669 stable, Xiaocheng Tang, Xiaofang Wang, Xiaojian Wu, Xiaolan Wang, Xide Xia, Xilun Wu, 670 Xinbo Gao, Yanjun Chen, Ye Hu, Ye Jia, Ye Qi, Yenda Li, Yilin Zhang, Ying Zhang, Yossi Adi, 671 Youngjin Nam, Yu, Wang, Yuchen Hao, Yundi Qian, Yuzi He, Zach Rait, Zachary DeVito, Zef 672 Rosnbrick, Zhaoduo Wen, Zhenyu Yang, and Zhiwei Zhao. The llama 3 herd of models, 2024. 673 URL https://arxiv.org/abs/2407.21783.

- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B Hashimoto. Length-controlled al pacaeval: A simple way to debias automatic evaluators. *arXiv preprint arXiv:2404.04475*, 2024a.
- Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators, 2024b. URL https://arxiv.org/abs/2404.04475.
- Bryan Eikema and Wilker Aziz. Is map decoding all you need? the inadequacy of the mode in neural machine translation, 2020.
- Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation, 2018. URL https://arxiv.org/abs/1805.04833.
  - António Farinhas, José de Souza, and Andre Martins. An empirical study of translation hypothesis ensembling with large language models. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 11956–11970, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.733. URL https://aclanthology.org/2023. emnlp-main.733.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. Quality-aware decoding for neural machine translation. In Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 1396–1412, Seattle, United States, July 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.100. URL https://aclanthology.org/2022.naacl-main.100.
- Mara Finkelstein, Subhajit Naskar, Mehdi Mirzazadeh, Apurva Shah, and Markus Freitag. Mbr and qe finetuning: Training-time distillation of the best and most expensive decoding methods, 2024. URL https://arxiv.org/abs/2309.10966.

702 703 704 705	Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. <i>Transactions of the Association for Computational Linguistics</i> , 10:811–825, 2022. doi: 10.1162/tacl_a_00491. URL https://aclanthology.org/2022.tacl-1.47.
706 707 708	Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. Epsilon sampling rocks: Investigating sampling strategies for minimum bayes risk decoding for machine translation, 2023.
709 710 711	Alex Graves. Sequence transduction with recurrent neural networks, 2012. URL https: //arxiv.org/abs/1211.3711.
712 713 714 715	Caglar Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Ksenia Konyushkova, Lotte Weerts, Abhishek Sharma, Aditya Siddhant, Alex Ahern, Miaosen Wang, Chenjie Gu, Wolfgang Macherey, Arnaud Doucet, Orhan Firat, and Nando de Freitas. Reinforced self-training (rest) for language modeling, 2023. URL https://arxiv.org/abs/2308.08998.
716 717 718 719	Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja- cob Steinhardt. Measuring massive multitask language understanding, 2021. URL https: //arxiv.org/abs/2009.03300.
720 721 722	Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text degeneration, 2020. URL https://arxiv.org/abs/1904.09751.
723 724 725	Jiaxin Huang, Shixiang Shane Gu, Le Hou, Yuexin Wu, Xuezhi Wang, Hongkun Yu, and Jiawei Han. Large language models can self-improve, 2022. URL https://arxiv.org/abs/ 2210.11610.
726 727 728 729 730 731	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chap- lot, 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. Mistral 7b, 2023. URL https: //arxiv.org/abs/2310.06825.
732 733	Yuu Jinnai, Tetsuro Morimura, Ukyo Honda, Kaito Ariu, and Kenshi Abe. Model-based minimum bayes risk decoding, 2023.
734 735 736 737	Seungone Kim, Jamin Shin, Yejin Cho, Joel Jang, Shayne Longpre, Hwaran Lee, Sangdoo Yun, Seongjin Shin, Sungdong Kim, James Thorne, et al. Prometheus: Inducing fine-grained evaluation capability in language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2023.
739 740 741	Seungone Kim, Juyoung Suk, Shayne Longpre, Bill Yuchen Lin, Jamin Shin, Sean Welleck, Graham Neubig, Moontae Lee, Kyungjae Lee, and Minjoon Seo. Prometheus 2: An open source language model specialized in evaluating other language models, 2024.
742 743 744	Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. <i>arXiv preprint arXiv:2302.09664</i> , 2023.
745 746 747 748	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph E. Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the ACM SIGOPS 29th Symposium on Operating Systems Principles</i> , 2023.
749 750 751 752	Nathan Lambert, Valentina Pyatkin, Jacob Morrison, LJ Miranda, Bill Yuchen Lin, Khyathi Chandu, Nouha Dziri, Sachin Kumar, Tom Zick, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. Rewardbench: Evaluating reward models for language modeling, 2024. URL https://arxiv.org/abs/2403.13787.
754 755	Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catan- zaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding models, 2024. URL https://arxiv.org/abs/2405.17428.

773

786

100	Xuechen Li, Tianvi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulraiani, Carlos Guestrin, Percy
757	Liang, and Tatsunori B. Hashimoto. Alpacaeval: An automatic evaluator of instruction-following
758	models. https://github.com/tatsu-lab/alpaca_eval, 5 2023.
759	

- Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization* Branches Out, pp. 74–81, Barcelona, Spain, July 2004a. Association for Computational Linguistics. URL https://aclanthology.org/W04-1013.
- Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pp. 74–81, 2004b.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human falsehoods, 2022. URL https://arxiv.org/abs/2109.07958.
- Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler Hallinan, Luyu Gao, Sarah Wiegreffe, Uri Alon, Nouha Dziri, Shrimai Prabhumoye, Yiming Yang, Shashank Gupta, Bodhisattwa Prasad Majumder, Katherine Hermann, Sean Welleck, Amir Yazdanbakhsh, and Peter Clark. Selfrefine: Iterative refinement with self-feedback, 2023. URL https://arxiv.org/abs/ 2303.17651.
- Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih
   Yavuz. Sfr-embedding-mistral:enhance text retrieval with transfer learning. Salesforce
   AI Research Blog, 2024. URL https://blog.salesforceairesearch.com/
   sfr-embedded-mistral/.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le
  Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, Xiangru Tang, Dragomir
  Radev, Alham Fikri Aji, Khalid Almubarak, Samuel Albanie, Zaid Alyafeai, Albert Webson,
  Edward Raff, and Colin Raffel. Crosslingual generalization through multitask finetuning, 2023.
  URL https://arxiv.org/abs/2211.01786.
- Zach Nussbaum, John X. Morris, Brandon Duderstadt, and Andriy Mulyar. Nomic embed: Training a reproducible long context text embedder, 2024. URL https://arxiv.org/abs/2402. 01613.
- OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-787 cia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red 788 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-789 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-791 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, 792 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, 793 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey 794 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, 796 Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-797 son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan 798 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-799 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan 800 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, 801 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun 802 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook 804 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel 805 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,

810 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel 811 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-812 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, 813 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel 814 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, 815 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, 816 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra 817 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, 818 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-819 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, 820 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, 821 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, 822 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-823 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-824 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, 825 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming 827 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao 828 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL 829 https://arxiv.org/abs/2303.08774. 830

- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
  evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pp. 311–318, 2002.
- Yiwei Qin, Weizhe Yuan, Graham Neubig, and Pengfei Liu. T5score: Discriminative fine-tuning of generative evaluation metrics. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 15185–15202, 2023.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D. Manning, and
   Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model,
   2024. URL https://arxiv.org/abs/2305.18290.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. Comet: A neural framework for mt evaluation, 2020. URL https://arxiv.org/abs/2009.09025.
- Thibault Sellam, Dipanjan Das, and Ankur P. Parikh. Bleurt: Learning robust metrics for text generation, 2020. URL https://arxiv.org/abs/2004.04696.
- Noah Shinn, Federico Cassano, Edward Berman, Ashwin Gopinath, Karthik Narasimhan, and
   Shunyu Yao. Reflexion: Language agents with verbal reinforcement learning, 2023. URL
   https://arxiv.org/abs/2303.11366.
- Charlie Snell, Jaehoon Lee, Kelvin Xu, and Aviral Kumar. Scaling llm test-time compute optimally can be more effective than scaling model parameters. *arXiv preprint arXiv:2408.03314*, 2024.
- Yifan Song, Guoyin Wang, Sujian Li, and Bill Yuchen Lin. The good, the bad, and the greedy:
   Evaluation of llms should not ignore non-determinism, 2024. URL https://arxiv.org/
   abs/2407.10457.
- 855 Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam 856 Fisch, Adam R. Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, Agnieszka Kluska, Aitor Lewkowycz, Akshat Agarwal, Alethea Power, Alex Ray, Alex Warstadt, Alexander W. 858 Kocurek, Ali Safaya, Ali Tazarv, Alice Xiang, Alicia Parrish, Allen Nie, Aman Hussain, 859 Amanda Askell, Amanda Dsouza, Ambrose Slone, Ameet Rahane, Anantharaman S. Iyer, Anders Andreassen, Andrea Madotto, Andrea Santilli, Andreas Stuhlmüller, Andrew Dai, Andrew La, Andrew Lampinen, Andy Zou, Angela Jiang, Angelica Chen, Anh Vuong, Animesh 861 Gupta, Anna Gottardi, Antonio Norelli, Anu Venkatesh, Arash Gholamidavoodi, Arfa Tabas-862 sum, Arul Menezes, Arun Kirubarajan, Asher Mullokandov, Ashish Sabharwal, Austin Her-863 rick, Avia Efrat, Aykut Erdem, Ayla Karakaş, B. Ryan Roberts, Bao Sheng Loe, Barret Zoph,

864 Bartłomiej Bojanowski, Batuhan Özyurt, Behnam Hedayatnia, Behnam Neyshabur, Benjamin 865 Inden, Benno Stein, Berk Ekmekci, Bill Yuchen Lin, Blake Howald, Bryan Orinion, Cameron 866 Diao, Cameron Dour, Catherine Stinson, Cedrick Argueta, César Ferri Ramírez, Chandan Singh, 867 Charles Rathkopf, Chenlin Meng, Chitta Baral, Chiyu Wu, Chris Callison-Burch, Chris Waites, 868 Christian Voigt, Christopher D. Manning, Christopher Potts, Cindy Ramirez, Clara E. Rivera, Clemencia Siro, Colin Raffel, Courtney Ashcraft, Cristina Garbacea, Damien Sileo, Dan Garrette, Dan Hendrycks, Dan Kilman, Dan Roth, Daniel Freeman, Daniel Khashabi, Daniel Levy, 870 Daniel Moseguí González, Danielle Perszyk, Danny Hernandez, Danqi Chen, Daphne Ippolito, 871 Dar Gilboa, David Dohan, David Drakard, David Jurgens, Debajyoti Datta, Deep Ganguli, De-872 nis Emelin, Denis Kleyko, Deniz Yuret, Derek Chen, Derek Tam, Dieuwke Hupkes, Diganta 873 Misra, Dilvar Buzan, Dimitri Coelho Mollo, Diyi Yang, Dong-Ho Lee, Dylan Schrader, Eka-874 terina Shutova, Ekin Dogus Cubuk, Elad Segal, Eleanor Hagerman, Elizabeth Barnes, Eliza-875 beth Donoway, Ellie Pavlick, Emanuele Rodola, Emma Lam, Eric Chu, Eric Tang, Erkut Erdem, 876 Ernie Chang, Ethan A. Chi, Ethan Dyer, Ethan Jerzak, Ethan Kim, Eunice Engefu Manyasi, Ev-877 genii Zheltonozhskii, Fanyue Xia, Fatemeh Siar, Fernando Martínez-Plumed, Francesca Happé, 878 Francois Chollet, Frieda Rong, Gaurav Mishra, Genta Indra Winata, Gerard de Melo, Germán Kruszewski, Giambattista Parascandolo, Giorgio Mariani, Gloria Wang, Gonzalo Jaimovitch-879 López, Gregor Betz, Guy Gur-Ari, Hana Galijasevic, Hannah Kim, Hannah Rashkin, Hannaneh 880 Hajishirzi, Harsh Mehta, Hayden Bogar, Henry Shevlin, Hinrich Schütze, Hiromu Yakura, Hongming Zhang, Hugh Mee Wong, Ian Ng, Isaac Noble, Jaap Jumelet, Jack Geissinger, Jackson Kernion, Jacob Hilton, Jaehoon Lee, Jaime Fernández Fisac, James B. Simon, James Koppel, 883 James Zheng, James Zou, Jan Kocoń, Jana Thompson, Janelle Wingfield, Jared Kaplan, Jarema Radom, Jascha Sohl-Dickstein, Jason Phang, Jason Wei, Jason Yosinski, Jekaterina Novikova, 885 Jelle Bosscher, Jennifer Marsh, Jeremy Kim, Jeroen Taal, Jesse Engel, Jesujoba Alabi, Jiacheng Xu, Jiaming Song, Jillian Tang, Joan Waweru, John Burden, John Miller, John U. Balis, Jonathan Batchelder, Jonathan Berant, Jörg Frohberg, Jos Rozen, Jose Hernandez-Orallo, Joseph 888 Boudeman, Joseph Guerr, Joseph Jones, Joshua B. Tenenbaum, Joshua S. Rule, Joyce Chua, 889 Kamil Kanclerz, Karen Livescu, Karl Krauth, Karthik Gopalakrishnan, Katerina Ignatyeva, Katja 890 Markert, Kaustubh D. Dhole, Kevin Gimpel, Kevin Omondi, Kory Mathewson, Kristen Chiafullo, Ksenia Shkaruta, Kumar Shridhar, Kyle McDonell, Kyle Richardson, Laria Reynolds, Leo 891 Gao, Li Zhang, Liam Dugan, Lianhui Qin, Lidia Contreras-Ochando, Louis-Philippe Morency, 892 Luca Moschella, Lucas Lam, Lucy Noble, Ludwig Schmidt, Luheng He, Luis Oliveros Colón, 893 Luke Metz, Lütfi Kerem Şenel, Maarten Bosma, Maarten Sap, Maartje ter Hoeve, Maheen Fa-894 rooqi, Manaal Faruqui, Mantas Mazeika, Marco Baturan, Marco Marelli, Marco Maru, Maria 895 Jose Ramírez Quintana, Marie Tolkiehn, Mario Giulianelli, Martha Lewis, Martin Potthast, 896 Matthew L. Leavitt, Matthias Hagen, Mátyás Schubert, Medina Orduna Baitemirova, Melody 897 Arnaud, Melvin McElrath, Michael A. Yee, Michael Cohen, Michael Gu, Michael Ivanitskiy, Michael Starritt, Michael Strube, Michał Swędrowski, Michele Bevilacqua, Michihiro Yasunaga, 899 Mihir Kale, Mike Cain, Mimee Xu, Mirac Suzgun, Mitch Walker, Mo Tiwari, Mohit Bansal, 900 Moin Aminnaseri, Mor Geva, Mozhdeh Gheini, Mukund Varma T, Nanyun Peng, Nathan A. 901 Chi, Nayeon Lee, Neta Gur-Ari Krakover, Nicholas Cameron, Nicholas Roberts, Nick Doiron, Nicole Martinez, Nikita Nangia, Niklas Deckers, Niklas Muennighoff, Nitish Shirish Keskar, 902 Niveditha S. Iyer, Noah Constant, Noah Fiedel, Nuan Wen, Oliver Zhang, Omar Agha, Omar El-903 baghdadi, Omer Levy, Owain Evans, Pablo Antonio Moreno Casares, Parth Doshi, Pascale Fung, 904 Paul Pu Liang, Paul Vicol, Pegah Alipoormolabashi, Peiyuan Liao, Percy Liang, Peter Chang, Pe-905 ter Eckersley, Phu Mon Htut, Pinyu Hwang, Piotr Miłkowski, Piyush Patil, Pouya Pezeshkpour, 906 Priti Oli, Qiaozhu Mei, Qing Lyu, Qinlang Chen, Rabin Banjade, Rachel Etta Rudolph, Raefer 907 Gabriel, Rahel Habacker, Ramon Risco, Raphaël Millière, Rhythm Garg, Richard Barnes, Rif A. 908 Saurous, Riku Arakawa, Robbe Raymaekers, Robert Frank, Rohan Sikand, Roman Novak, Ro-909 man Sitelew, Ronan LeBras, Rosanne Liu, Rowan Jacobs, Rui Zhang, Ruslan Salakhutdinov, 910 Ryan Chi, Ryan Lee, Ryan Stovall, Ryan Teehan, Rylan Yang, Sahib Singh, Saif M. Moham-911 mad, Sajant Anand, Sam Dillavou, Sam Shleifer, Sam Wiseman, Samuel Gruetter, Samuel R. 912 Bowman, Samuel S. Schoenholz, Sanghyun Han, Sanjeev Kwatra, Sarah A. Rous, Sarik Ghazarian, Sayan Ghosh, Sean Casey, Sebastian Bischoff, Sebastian Gehrmann, Sebastian Schus-913 ter, Sepideh Sadeghi, Shadi Hamdan, Sharon Zhou, Shashank Srivastava, Sherry Shi, Shikhar 914 Singh, Shima Asaadi, Shixiang Shane Gu, Shubh Pachchigar, Shubham Toshniwal, Shyam Upad-915 hyay, Shyamolima, Debnath, Siamak Shakeri, Simon Thormeyer, Simone Melzi, Siva Reddy, 916 Sneha Priscilla Makini, Soo-Hwan Lee, Spencer Torene, Sriharsha Hatwar, Stanislas Dehaene, 917 Stefan Divic, Stefano Ermon, Stella Biderman, Stephanie Lin, Stephen Prasad, Steven T. Pianta918 dosi, Stuart M. Shieber, Summer Misherghi, Svetlana Kiritchenko, Swaroop Mishra, Tal Linzen, 919 Tal Schuster, Tao Li, Tao Yu, Tariq Ali, Tatsu Hashimoto, Te-Lin Wu, Théo Desbordes, Theodore 920 Rothschild, Thomas Phan, Tianle Wang, Tiberius Nkinyili, Timo Schick, Timofei Kornev, Ti-921 tus Tunduny, Tobias Gerstenberg, Trenton Chang, Trishala Neeraj, Tushar Khot, Tyler Shultz, 922 Uri Shaham, Vedant Misra, Vera Demberg, Victoria Nyamai, Vikas Raunak, Vinay Ramasesh, Vinay Uday Prabhu, Vishakh Padmakumar, Vivek Srikumar, William Fedus, William Saun-923 ders, William Zhang, Wout Vossen, Xiang Ren, Xiaoyu Tong, Xinran Zhao, Xinyi Wu, Xudong 924 Shen, Yadollah Yaghoobzadeh, Yair Lakretz, Yangqiu Song, Yasaman Bahri, Yejin Choi, Yichi 925 Yang, Yiding Hao, Yifu Chen, Yonatan Belinkov, Yu Hou, Yufang Hou, Yuntao Bai, Zachary 926 Seid, Zhuoye Zhao, Zijian Wang, Zijie J. Wang, Zirui Wang, and Ziyi Wu. Beyond the im-927 itation game: Quantifying and extrapolating the capabilities of language models, 2023. URL 928 https://arxiv.org/abs/2206.04615. 929

- Miloš Stanojević and Khalil Sima'an. Fitting sentence level translation evaluation with many dense features. In Alessandro Moschitti, Bo Pang, and Walter Daelemans (eds.), *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 202–206, Doha, Qatar, October 2014. Association for Computational Linguistics. doi: 10.3115/v1/D14-1025. URL https://aclanthology.org/D14-1025.
- Nisan Stiennon, Long Ouyang, Jeff Wu, Daniel M. Ziegler, Ryan Lowe, Chelsea Voss, Alec Radford,
   Dario Amodei, and Paul Christiano. Learning to summarize from human feedback, 2022. URL
   https://arxiv.org/abs/2009.01325.
- Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. Follow the wisdom of the crowd: Effective text generation via minimum bayes risk decoding, 2022.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023a. URL https://arxiv.org/abs/2302.13971.
- 945 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-946 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy 947 Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, 948 Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel 949 Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, 950 Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, 951 Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, 952 Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh 953 Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen 954 Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, 955 Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 956 2023b. URL https://arxiv.org/abs/2307.09288.
- Tu Vu, Kalpesh Krishna, Salaheddin Alzubi, Chris Tar, Manaal Faruqui, and Yun-Hsuan Sung.
   Foundational autoraters: Taming large language models for better automatic evaluation. *arXiv* preprint arXiv:2407.10817, 2024.
- Jun Wang, Eleftheria Briakou, Hamid Dadkhahi, Rishabh Agarwal, Colin Cherry, and Trevor
   Cohn. Don't throw away data: Better sequence knowledge distillation, 2024a. URL https:
   //arxiv.org/abs/2407.10456.
- Peifeng Wang, Austin Xu, Yilun Zhou, Caiming Xiong, and Shafiq Joty. Direct judgement preference optimization. *arXiv preprint arXiv:2409.14664*, 2024b.
- Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models, 2023.
- Jason Wei, Maarten Bosma, Vincent Y. Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du,
   Andrew M. Dai, and Quoc V. Le. Finetuned language models are zero-shot learners, 2022. URL <a href="https://arxiv.org/abs/2109.01652">https://arxiv.org/abs/2109.01652</a>.

- 972 Sean Welleck, Amanda Bertsch, Matthew Finlayson, Hailey Schoelkopf, Alex Xie, Graham Neubig, 973 Ilia Kulikov, and Zaid Harchaoui. From decoding to meta-generation: Inference-time algorithms 974 for large language models, 2024. URL https://arxiv.org/abs/2406.16838. 975 Tianhao Wu, Weizhe Yuan, Olga Golovneva, Jing Xu, Yuandong Tian, Jiantao Jiao, Jason Weston, 976 and Sainbayar Sukhbaatar. Meta-rewarding language models: Self-improving alignment with 977 llm-as-a-meta-judge. arXiv preprint arXiv:2407.19594, 2024. 978 979 Jing Xu, Andrew Lee, Sainbayar Sukhbaatar, and Jason Weston. Some things are more cringe than 980 others: Preference optimization with the pairwise cringe loss. arXiv preprint arXiv:2312.16682, 981 2023. 982 Guangyu Yang, Jinghong Chen, Weizhe Lin, and Bill Byrne. Direct preference optimization for 983 neural machine translation with minimum bayes risk decoding, 2024. URL https://arxiv. 984 org/abs/2311.08380. 985 Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik 986 Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023. 987 URL https://arxiv.org/abs/2305.10601. 988 989 Seonghyeon Ye, Doyoung Kim, Sungdong Kim, Hyeonbin Hwang, Seungone Kim, Yongrae Jo, 990 James Thorne, Juho Kim, and Minjoon Seo. Flask: Fine-grained language model evaluation 991 based on alignment skill sets, 2024. URL https://arxiv.org/abs/2307.10928. 992 Weizhe Yuan, Graham Neubig, and Pengfei Liu. Bartscore: Evaluating generated text as text gener-993 ation. Advances in Neural Information Processing Systems, 34:27263–27277, 2021. 994 995 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, 996 and Jason Weston. Self-rewarding language models, 2024a. 997 Weizhe Yuan, Richard Yuanzhe Pang, Kyunghyun Cho, Xian Li, Sainbayar Sukhbaatar, Jing Xu, 998 and Jason Weston. Self-rewarding language models, 2024b. URL https://arxiv.org/ 999 abs/2401.10020. 1000 Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a ma-1002 chine really finish your sentence?, 2019. URL https://arxiv.org/abs/1905.07830. 1003 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. Bertscore: Evaluat-1004 ing text generation with bert. In International Conference on Learning Representations. 1005 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 1008 Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. Banghua Zhu, Evan Frick, Tianhao Wu, Hanlin Zhu, and Jiantao Jiao. Starling-7b: Improving llm 1010 helpfulness & harmlessness with rlaif, November 2023a. 1011 1012 Lianghui Zhu, Xinggang Wang, and Xinlong Wang. Judgelm: Fine-tuned large language models 1013 are scalable judges. 2023b. 1014 Lianghui Zhu, Xinggang Wang, and Xinlong Wang. Judgelm: Fine-tuned large language models 1015 are scalable judges, 2023c. URL https://arxiv.org/abs/2310.17631. 1016 1017 1018 1020 1021 1023 1024
- 1025

# <sup>1026</sup> A MBR DECODING: ADDITIONAL RESULTS

1027 1028 1029

1030

1031 1032 1033

1035 1036 This section presents additional results from our experiments exploring the use of MBR decoding to improve test-time performance.

### 1034 A.1 ADDITIONAL MBR UTILITY METRICS

1037 1038 1039

> 1040 1041

> 1043 1044 1045

	Llama2-7B	Llama2-70B	Avg. $\Delta$
Greedy	14.4	22.8	0
ROUGE-1 MBR	16.2	24.7	1.85
SFR-Embedder MBR	12.1	22.2	-1.45
Prometheus MBR	17.7	26.2	3.35
ROUGE-2 MBR	16.6	24.6	2.00
ROUGE-L MBR	15.5	24.7	1.50
NVEmbed-Embedder MBR	14.1	22.1	-0.50
Nomic-Embedder MBR	16.3	24.1	1.60

Table 5: AlpacaEval 2.0 win rates (%) for additional MBR decoding experiments, along with the average win rate differences compared to greedy decoding across all models (denoted **Avg.**  $\Delta$ ).

1049

We experiment with two additional ROUGE variants and two additional dense embedders as utility metrics. The two ROUGE variants are ROUGE-2 and ROUGE-L, which detect bigram overlap and longest co-occurring n-gram overlap respectively. The two dense embedders are NVEmbed-v2 (Lee et al., 2024) and Nomic-Text-v1.5 (Nussbaum et al., 2024), which, like SFR Embedder, are strong, long-context dense embedders that rank highly on the MTEB leaderboard Muennighoff et al. (2023).

We find that MBR decoding with ROUGE-2 performs slightly better than MBR decoding with ROUGE-1, while MBR decoding with ROUGE-L performs slightly worse. MBR decoding with NVEmbed and Nomic both perform better than MBR decoding with SFR Embedder, with Nomic showing some improvement relative to greedy decoding. The latter result suggests that dense embedders could potentially be used for MBR decoding, although further work is required to understand what properties make for good MBR embedders. Overall, none of our additional utility metrics provide comparable improvements to MBR with Prometheus.

- 1062
- 1063 1064

### A.2 COMPARISON WITH UNIVERSAL SELF-CONSISTENCY

1066

1067We compare MBR decoding to Universal Self-Consistency (USC) (Chen et al., 2023). In USC,1068 $N_{cand}$  outputs are sampled from the LLM and passed directly to the LLM for consistency detection.1069This entails prompting the LLM to choose the most consistent output. In Chen et al. (2023), the1070authors demonstrate that USC improves over greedy decoding for mathematical reasoning, code1071generation, summarisation and question-answering.

1072 The limited context lengths of LLMs poses a significant challenge when using USC, as it requires 1073 fitting all  $N_{\text{cand}} = 30$  samples into a single prompt. In contrast, MBR decoding only requires fitting 1074 two outputs into a single prompt as utility is computed pairwise. As our existing choice of models 1075 have limited context lengths (4096 tokens for Llama2, 8192 tokens for Llama3) and our outputs can be long (up to 1024 tokens), we are unable to assess USC on equal footing with MBR decoding using these models without significantly reducing  $N_{\text{cand}}$ . In order to facilitate a fair comparison, we 1077 therefore use the Llama-3.1 models (Dubey et al., 2024) in place of the Llama2 and Llama3 models 1078 for this experiment. The Llama-3.1 models possess context lengths of 128k, thereby allowing us to 1079 fit all  $N_{cand}$  samples into a single to prompt as required. Our USC prompt is as follows:

1080 You are given a collection of 30 responses to a prompt. Select 1081 the most consistent response based on majority consensus. The 1082 most concistent response should be the most representative of all the responses provided. You should consider a variety of 1084 factors when evaluating consistency, including content, arguments and examples employed, style, structure and final answer, if relevant. Do not pass judgement on the quality or correctness 1086 of the response. Consider only consistency. 1087

Provide a short explanation of your choice, followed by your choice. Your choice should follow this format: "Most Consistent Response: [[Response ID]]", for example: "Most Consistent Response: [[15]]" if response 15 is the most consistent amongst all responses.

Responses: {responses}

 Llama3.1-8B
 Llama3.1-70B
 Avg. Δ

 Greedy
 34.2
 42.1
 0

 USC
 37.3
 43.7
 2.35

 MBR Prometheus
 40.2
 45.8
 4.85

1101Table 6: AlpacaEval 2.0 win rates (%) for Llama3.1 models with greedy decoding, USC and MBR decoding1102with Prometheus, along with the average win rate differences compared to greedy decoding across all models1103(denoted Avg.  $\Delta$ ).

We document our findings in Table 6. We find while that USC provides some improvements over greedy decoding, this improvement is smaller than the improvement provided by MBR decoding with Prometheus.

1108

1125

1104

1088

1089

1090

1091

1093 1094

1095

1098

1099

1100

1109 A.3 SAMPLING BASELINES

1110 We conduct additional baseline experiments with top-p (Holtzman et al., 2020) and top-k (Fan et al., 2018) sampling. We use these sampling methods to directly obtain a final output, and do not explore using replacing temperature sampling with these sampling strategies for generating the hypothesis set - we leave this to future work.

1114				
1115		Llama2-7B	Llama2-70B	Avg. $\Delta$
1116	Greedy	14.4	22.8	0
1110	ROUGE-1 MBR	16.2	24.7	1.85
1117	Prometheus MBR	17.7	26.2	3.35
1118	top- $p (p = 0.9, t = 0.3)$	14.3	23.7	0.40
1119	top- $p$ ( $p = 0.9, t = 0.7$ )	14.7	23.9	0.70
1120	top- $p (p = 0.5, t = 0.3)$	15.3	24.9	1.50
1121	top- $p (p = 0.5, t = 0.7)$	15.0	24.9	1.35
1100	top- $k (k = 50, t = 0.3)$	14.7	23.6	0.55
1122	top- $k (k = 50, t = 0.7)$	15.0	23.5	0.65
1123	top- $k$ ( $k = 20, t = 0.3$ )	15.0	24.7	1.25
1124	top- $k$ ( $k = 20, t = 0.7$ )	14.7	24.0	0.75

Table 7: AlpacaEval 2.0 win rates for top-p and top-k decoding, along with the average score differences compared to greedy decoding across all models (denoted **Avg.**  $\Delta$ ). Top-p and top-k sampling improve over greedy decoding, but do not match the performance improvements of Prometheus MBR decoding.

We find that top-*p* and top-*k* sampling improves over greedy decoding, and can achieve performance
close to that of MBR decoding with ROUGE. The improvements are nonetheless much smaller
than MBR decoding with Prometheus. Nonetheless, these results demonstrate that top-*p* and top-*k*sampling could be used to produce hypothesis sets containing higher quality candidates, which could
in turn improve downstream MBR decoding performance. We believe this is an exciting avenue for future work.

#### A.4 **GENERATION LENGTHS**



Figure 6: Average generation lengths (in words) for various decoding strategies, averaged over all five generator LLMs. MBR with Prometheus produces slightly longer outputs than most baseline methods, although these outputs are still far shorter than those produced by MBR with SFR-embedder and by the longest decoding baseline.



Figure 7: Average generation lengths (in words) for various models and decoding strategies on AlpacaEval. MBR with Prometheus produces slightly longer outputs than most baseline methods, although these outputs are still far shorter than those produced by MBR with SFR-embedder and by the longest decoding baseline.

We compute the average generation lengths (in words) of our five generator LLMs on AlpacaEval 2.0 with different decoding strategies and plot the results in Figures 6 and 7. We find that MBR decoding with Prometheus yields outputs that are on average longer than those yielded by greedy decoding, beam search or Prometheus BoN decoding. They are nonetheless much shorter than the outputs yielded by MBR decoding with SFR-Embedder and decoding by selecting the longest output. We hypothesise that MBR decoding with Prometheus encourages selection of more detailed responses which improves model performance, although we note that verbosity alone is not enough to improve performance on our benchmarks as the *longest* baseline fails to yield any gains, and as we use length-controlled metrics for evaluation.

# 1188 B GPT-40-JUDGE DETAILS

1190

We use GPT-40 as a judge for both AlpacaEval 2.0 and MT-Bench. More specifically, we use gpt-40-2024-05-13, accessed through an Azure endpoint.

For AlpacaEval 2.0, we use the official AlpacaEval implementation<sup>4</sup> to conduct evaluation. We use the weighted\_alpaca\_eval\_gpt4\_turbo evaluator and baseline results against the default gpt-4-1106-preview generations unless other specified.

For MT-Bench, we use the single- and multi-turn judge prompts provided by LLMSYS FastChat<sup>5</sup> as our judge prompts. Due to stochasticity in the outputs of the judge models, even with temperature set to zero, we generate three judgements per sample and take as the final score the median of the three judgement scores.

1201

1202 1203

1204 1205 1206

1207

1208

## C PROMETHEUS SCORING RUBRICS

Prometheus takes as input a scoring rubric that defines scoring criteria to be used during evaluation. We use a single, generic scoring rubric for all our experiments:

[Consider a wide range of factors such as the helpfulness, 1209 relevance, accuracy, depth, creativity, and level of detail of 1210 the response.] 1211 Score 1: The answer is completely unhelpful and incorrect. 1212 Nothing useful can be learned from it. 1213 Score 2: The answer contains some helpful and useful information, 1214 but major flaws, in terms of facuality, accuracy and relevance, 1215 are also present. 1216 The answer is mostly helpful and relevant, although Score 3: 1217 minor flaws exist. 1218 Score 4: The answer is accurate, relevant and helpful, although 1219 there are some clear improvements that can be made with respect to depth, creativity and detail. 1220 Score 5: The answer is excellent. It is completely accurate and relevant, and demonstrates a high degree of depth and creativity. 1222

1223 1224

1225

1226

We hypothesise that the performance of MBR and BoN decoding with Prometheus could be improved through further optimisation of the scoring rubric, particularly with question-specific adjustments, where unique rubrics are tailored to each question. We leave this to future work.

1227 1228 1229

# <sup>1230</sup> D ALPACAEVAL CATEGORIES

1231 1232

1233 We classify questions from AlpacaEval and MT-Bench and then evaluate performance by category.

For AlpacaEval, we perform manual inspection on the dataset and identify ten common question categories. We then use GPT-40 to classify questions based on these categories, using the following prompt:

1238

1239

1240 1241

<sup>&</sup>lt;sup>4</sup>github.com/tatsu-lab/alpaca\_eval

<sup>&</sup>lt;sup>5</sup>github.com/lm-sys/FastChat

```
1242
        Categorise an instruction based on the following list of
1243
        categories. Only choose one category, and only return the
1244
        category, nothing else.
1945
1246
        - Creative writing
        - Business, technical and scientific writing
1247
        - Argumentation, debate and persuasion
1248
        - Mathematical reasoning
1249
        - Puzzles and logical reasoning
1250
        - Coding
1251
        - How-to and other guides
1252
        - Recommendations and advice
1253
        - Factual question-answering
1254
        - Other
1255
1256
                                 Bob has 5 sisters and 1 brother.
        Example Instruction:
                                                                         How many
1257
        siblings does one of Bob's sisters have?
        Category: Puzzles and logical reasoning
        Example Instruction:
                                 Write me a resignation email for my job as
1259
        an accountant, explaining that I am leaving to pursue my dream of
1260
        becoming a lion tamer.
1261
        Category: Business, technical and scientific writing
1262
        Example Instruction: What factors gave rise to the English Civil
1263
        War. Category: Factual question-answering
1264
        Example Instruction: I am visiting Kyoto next April.
                                                                         Recommend
1265
       me 10 things to do!
1266
        Category: Recommendations and advice
1267
        Example Instruction: Pretend to be Donald Trump - write a speech
1268
        announcing that you are becoming a Democrat.
        Category: Creative Writing
1269
        Instruction:
                       {instruction}
1270
        Category:
1271
1272
1273
      The percentage of questions assigned to them are listed in Table 8.
1274
1275
                                                    Percentage of Questions
                                Category
1276
                          Factual question-answering
                                                            24.3
1277
                           How-to and other guides
                                                            20.6
1278
                                                            12.1
                         Recommendations and advice
1279
                           Mathematical reasoning
                                                             3.2
                                                            22
                                 Other
1280
                              Creative writing
                                                            11.8
1281
                     Business, technical and scientific writing
                                                            12.7
1282
                         Puzzles and logical reasoning
                                                             3.5
1283
                                 Coding
                                                             6.7
1284
                                                             2.7
                      Argumentation, debate and persuasion
1285
                        Table 8: AlpacaEval question categories identified by GPT-40.
1286
1287
      The results of our performance by category analysis for MBR decoding with Prometheus on Al-
1288
      pacaEval are illustrated in the main text, in Figure 3.
1289
1290
1291
```

1292

### E LLAMA3 AS AN MBR AND BON UTILITY METRIC

We experiment with using Llama3-70b-Instruct as an MBR and BoN utility metric in Section 3.3.1. This entails prompting the model to act as either a reference-based or reference-free evaluator.

Our prompt for single-turn reference-based evaluation is as follows:

1 1 1	
Please	a act as an impartial judge and evaluate the quality of
the r	esponse provided by an AI assistant to the user question
displa	aved below. In addition to the user question, you are
also (	given a reference answer. This is the best possible answ
rovi	ded by a human expert. You should evaluate the assistant
espoi	nse based on this. A good assistant's answer should shar
the c	ontent and style of the reference answer. Begin your
evalu	ation by providing a short explanation. Be as objective
possil	ole. After providing your explanation, you must rate the
respo	nse on a scale of 1 to 10 by strictly following this form
"[[ra	ting]]", for example: "Rating: [[5]]".
[011es:	tionl
ques	tion}
[Refe	rence Answer]
{refe	rence}
[m]	Stand of Tasistant ( Tanana)
[The	Start of Assistant's Answer]
[The ] [answe	Start of Assistant's Answer] er} End of Assistant's Answer]
The answe	Start of Assistant's Answer] er} End of Assistant's Answer] opt for single-turn reference-free evaluation is as follows:
[The : {answe [The ] ur prom [Inst: Please the re displa explai explai by st: "Rati;	Start of Assistant's Answer] er} End of Assistant's Answer] mpt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]".
[The : {answe [The ] ur prom [Inst: Please the re displa explan explan oy st: [Ques]	Start of Assistant's Answer] er} End of Assistant's Answer] upt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]".
[The : {answe [The ] ur prom [Inst: Please the re displa explai explai by st: [Ques: {ques:	Start of Assistant's Answer] er} End of Assistant's Answer] mpt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]". tion] tion]
[The : {answe [The ] ur prom [Inst: Please the re displa explase explase by st: [Quese {quese	Start of Assistant's Answer] er} End of Assistant's Answer] mpt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]". tion] tion}
[The : {answe [The ] ur prom [Inst: Please the re displation explation explation guess {quess {quess [The :	Start of Assistant's Answer] er} End of Assistant's Answer] mpt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]". tion] tion} Start of Assistant's Answer]
[The : {answe [The ] ur prom [Inst: Please the re displate explate explate y st: [Quess {quess [The : {answe	<pre>Start of Assistant's Answer] er} End of Assistant's Answer] mpt for single-turn reference-free evaluation is as follows: ruction] e act as an impartial judge and evaluate the quality of esponse provided by an AI assistant to the user question ayed below. Begin your evaluation by providing a short nation. Be as objective as possible. After providing you nation, you must rate the response on a scale of 1 to 10 rictly following this format: "[[rating]]", for example: ng: [[5]]". tion] tion} Start of Assistant's Answer] er}</pre>

1350 Please act as an impartial judge and evaluate the quality of 1351 the response provided by an AI assistant to the user question 1352 displayed below. Your evaluation should focus on the assistant's 1353 answer to the second user question. In addition to the user 1354 question and conversation history, you are also given a reference answer. This is the best possible answer to the second user 1355 question provided by a human expert. You should evaluate the 1356 assistant's response based on this. A good assistant's answer 1357 should share the content and style of the reference answer. Begin 1358 your evaluation by providing a short explanation. Be as objective 1359 as possible. After providing your explanation, you must rate the 1360 response on a scale of 1 to 10 by strictly following this format: 1361 "[[rating]]", for example: "Rating: [[5]]". 1362

1363 1364

1365

1366

1367

1368

1369

1370

1371

For multi-turn, reference-free evaluation:

```
Please act as an impartial judge and evaluate the quality of
the response provided by an AI assistant to the user question
displayed below. Your evaluation should focus on the assistant's
answer to the second user question. Begin your evaluation by
providing a short explanation. Be as objective as possible.
After providing your explanation, you must rate the response on a
scale of 1 to 10 by strictly following this format: "[[rating]]",
for example: "Rating: [[5]]".
```

The prompt template for both reference-based and reference-free multi-turn evaluation is:

```
1376
       <|The Start of Assistant A's Conversation with User|>
1377
1378
       ### User:
1379
       {question_1}
1380
1381
       ### Assistant A:
1382
       {answer_1}
       ### User:
1384
        {question_2}
1385
1386
       ### Assistant A:
1387
       {answer 2}
1388
1389
       <|The End of Assistant A's Conversation with User|>
```

1390 1391 1392

1393

1394 1395

### F REWARD MODEL AS A BON UTILITY METRIC

	2-7B	2-13B	2-70B	3-8B	3-70B	Avg. $\Delta$
Greedy	5.72	5.90	6.50	7.54	8.29	0
BS	5.58	5.95	6.49	7.30	8.20	-0.09
Prometheus BoN	5.77	6.08	6.65	7.66	8.42	0.13
Starling-RM-7b BoN	<u>5.99</u>	6.49	6.85	7.88	<u>8.46</u>	0.34
Prometheus MBR	6.10	6.26	<u>6.79</u>	7.69	8.50	0.28

Table 9: MT-Bench scores for various models and decoding strategies, along with the average score differences compared to greedy decoding across all models (denoted **Avg.**  $\Delta$ ). BoN decoding with reward model StarlingRM-7b marginally outperforms MBR decoding with Prometheus.

26

We evaluate BoN decoding using Starling-RM-7B-alpha (Zhu et al., 2023a) as the utility metric.
Starling-RM is a strong 7B sequence classifier reward model trained to facilitate RLHF (Stiennon et al., 2022) that achieves a similar overall score to Prometheus-2-7B on RewardBench (Lambert et al., 2024). It takes as inputs a prompt and a single candidate and outputs a scalar reward. As with reward models in general, Starling-RM does not support reference-based evaluation.

1409 We find that BoN decoding with Starling-RM as a utility metric outperforms MBR decoding with 1410 Prometheus by a small margin. We have two possible explanations for this discrepancy. Firstly, 1411 Starling-RM achieves a much higher score than Prometheus-2-7B on RewardBench's Chat task 1412 (likely due to the distribution of its training data), suggesting that it may simply be more suited to 1413 our particular benchmarks. Secondly, by providing continuous scalar rewards instead of discrete integer scores, Starling-RM enables more fine-grained evaluation of candidate outputs, allowing it 1414 to distinguish between outputs that Prometheus might have rated equally. We believe that using a 1415 reference-based reward model as the metric for MBR decoding could combine the advantages of 1416 fine-grained scoring with the increased reliability of consensus-based output selection. We leave the 1417 development of such models to future work. 1418

## G SELF-TRAINING: ADDITIONAL RESULTS



### G.1 HEAD-TO-HEAD RESULTS

Figure 8: AlpacaEval 2.0 win rates (%) for self-trained models and various SFT baselines against *sft* with greedy decoding. Generation for all *dpo* models is done with greedy decoding. We find that MBR self-training with DPO allows models to match their MBR-decoding performance.

We conduct head-to-head evaluation of our DPO self-trained models and various SFT baselines against *sft* with greedy decoding using the AlpacaEval 2.0 and illustrate our findings in Figure 8. Our head-to-head results show that our MBR self-trained models outperform our BoN self-trained models, the *sft* with beam search baseline, and the full *sft* with greedy decoding baseline. Our MBR self-trained models match the performance of the *sft* with MBR decoding baseline.

### G.2 RESULTS ON NLP BENCHMARKS

	Model	MMLU $(\uparrow)$	ARC challenge $(\uparrow)$	HellaSwag (↑)	TruthfulQA $(\uparrow)$
	sft	47.2	57.3	80.6	51.6
7B	dpo-3-BoN	47.2	57.5	80.6	53.5
	dpo-3-MBR	47.3	57.1	80.7	52.6
		•			
	sft	56.1	62.3	83.5	48.4
13B	dpo-3-BoN	56.3	62.5	83.5	48.2
	dpo-3-MBR	56.1	62.6	83.5	47.4

<sup>1452</sup>Table 10: Evaluation results of Prometheus self-trained models on four different NLP benchmarks. We find that1453MBR and BoN self-training maintains performance on across all four datasets compared with the *sft* models.

1442

1419 1420

1421 1422

We assess our self-trained models on four different NLP benchmarks: MMLU (Hendrycks et al., 2021), ARC challenge (Clark et al., 2018), HellaSwag (Zellers et al., 2019) and TruthfulQA (Lin et al., 2022), and report our results in Table 10. We find that self-training maintains performance across all four benchmarks despite us using a training dataset that is irrelevant for these tasks. This

27

<sup>1454</sup> 

shows that MBR self-training can be used to improve the instruction-following abilities of models without jeopardising other skills. 

### G.3 MBR SELF-TRAINING WITH ROUGE

	AlpacaEval 2.0		MT-Bench	
	7B	13B	7B	13B
sft	5.18	8.24	5.43	5.85
dpo-1-MBR-Prometheus	5.68	10.8	5.78	6.48
dpo-2-MBR-Prometheus	7.22	13.9	6.11	6.73
dpo-3-MBR-Prometheus	8.86	15.3	6.14	6.75
dpo-1-MBR-ROUGE	4.66	5.61	7.65	5.98
dpo-2-MBR-ROUGE	5.83	5.78	9.01	6.06
dpo-3-MBR-ROUGE	5.42	5.67	8.31	5.91

Table 11: AlpacaEval 2.0 win rates (%) and MT-Bench scores for models self-trained using DPO with Prometheus and with ROUGE-1 as utility metrics. MBR self-training with ROUGE fails to yield substan-tial gains. 

We explore using ROUGE-1 as the utility metric for MBR self-training. We find that MBR self-training with ROUGE fails to yield substantial gains. Improvements also saturate quickly, with model performance decreasing after the third training iteration. These results highlight the impor-tance of choosing the correct MBR utility metric for self-training.

### G.4 ALTERNATIVE PAIR SELECTION STRATEGIES

1482		AlpacaEval 2.0	MT-Bench
1483	sft	5.18	5.43
1484	$dpo-1-MBR_{BW}$	5.68	5.78
1485	$dpo$ -2-MBR $_{BW}$	7.22	6.11
1400	$dpo$ -3-MBR $_{BW}$	8.86	6.14
1400	$dpo-1-MBR_{BMW}$	4.95	5.78
1487	$dpo-2-MBR_{BMW}$	8.06	5.96
1488	dpo-3-MBR <sub>BMW</sub>	8.88	5.97
1489			

Table 12: AlpacaEval 2.0 win rates (%) and MT-Bench scores for models self-trained using DPO with BW and BMW preference pair selection strategies with sft-7b as the base model. Generation is done using greedy decoding for all models. We find that our BMW models perform similarly to the original BW models on AlpacaEval 2.0 and slightly worse on MT-Bench. 

We investigate using an alternative MBR preference pair selection strategy. Following Yang et al. (2024), we create two preference pairs from each sample, one formed from the highest-scoring (best) and median-scoring (mid) outputs, and another formed from the median-scoring and lowest-scoring (worst) outputs. We follow the exact same DPO training procedure as before, but replace our original BW training set with this new BMW training set. 

We document our results in Table 12. We find that our BMW models perform similarly to the original BW models on AlpacaEval 2.0 and slightly worse on MT-Bench. We do not pursue this line of work further and leave additional investigations to future work. 

G.5 PERFORMANCE BY CATEGORY

We replicate the question category analysis described in Section 3.3 and Appendix D for our Prometheus MBR self-trained models and report results in Figure 4. We find that performance rel-ative to the *sft* models improves across almost all question categories, with performance on writing based categories improving more than on reasoning based categories. Performance on mathematical reasoning remains unchanged for the 7B model and decreases for the 13B model. 

To better understand this discrepancy, we sample 1000 prompts from our 12000 UltraChat training prompts and categorise them using GPT-40, following the same procedure described in Section 3.3 and Appendix D. We document our results in Table 13. We find that reasoning-based questions

1512	Category	Percentage of Questions
1513	Factual question-answering	26.4
1514	How-to and other guides	21.1
1515	Recommendations and advice	4.4
1510	Mathematical reasoning	0.1
1516	Other	1.3
1517	Creative writing	18.6
1518	Business, technical and scientific writing	19.0
1519	Puzzles and logical reasoning	0.1
1520	Coding	6.3
1520	Argumentation, debate and persuasion	2.7

Table 13: UltraChat-200k question categories identified by GPT-40. We perform this analysis on a subsample of 1000 prompts sampled randomly from our 12000 training prompts.

(coding, puzzles and logical reasoning, mathematical reasoning) are underrepresented in this dataset,
 with mathematical reasoning and puzzles and logical reasoning especially underrepresented. We
 attribute our models' inconsistent improvements in these areas to this lack of data.

# 1529 G.6 GENERATION LENGTHS

st the strain of the strain of

Figure 9: Average generation lengths (in words) after iterative Prometheus MBR and BoN self-training on AlpacaEval. MBR self-training with DPO teaches the model to generate more detailed responses.

We measure the generation lengths of our Prometheus MBR and BoN self-trained models on AlpacaEval. We find that MBR self-training with DPO teaches the model to generate longer and more
detailed responses at every iteration. In contrast, BoN self-training and MBR self-training with SFT
only results in small increases in generation lengths.

### 1553 H HUMAN STUDY

	Win	Draw	Loss
Prometheus MBR vs. Greedy	30.0	53.3	16.7
Prometheus BoN vs. Greedy	21.7	60.0	18.3

1559Table 14: Head-to-head evaluation of Prometheus MBR and Prometheus BoN vs. greedy decoding for Llama2-<br/>70b conducted on the AlpacaEval dataset by human evaluators.

We conduct a human study for key MBR inference (Section 3) and MBR distillation (Section 4)
experiments. The objective of this study is to verify that our LLM evaluation results (AlpacaEval 2.0 and MT-Bench) align with real human judgements.

1565 We recruited 4 volunteers, each of whom were given 60 samples in total to evaluate. Each sample consisted of a randomly-sampled AlpacaEval prompt and two corresponding generations displayed

1566		Win	Draw	Loss
1567	dpo-3-MBR vs. sft	46.7	43.3	10.0
1568	dpo-3-BoN vs. sft	33.3	55.0	11.7

1569 Table 15: Head-to-head evaluation of 13B dpo-3-MBR and dpo-3-BoN vs. sft conducted on the AlpacaEval 1570 dataset by human evaluators. Greedy decoding used for all models.

1571 1572

in a random order. Volunteers were asked to select their preferred generation and, if neither genera-1573 tion was preferred, to then rate the generations as equal. 1574

1575 We conducted four head-to-head experiments. The first two experiments, corresponding to experiments in Section 3, were between Prometheus MBR vs. greedy decoding and Prometheus BoN 1576 vs. greedy decoding with Llama2-70b. The next two experiments, corresponding to experiments in Section 4, were between 13B dpo-3-MBR vs. sft and dpo-3-BoN vs. sft. We used 60 samples for 1578 each experiment. Volunteers were given an even distribution of samples from each experiment. 1579

1580 Our results show good alignment with our LLM evaluation results. We find that Prometheus MBR 1581 decoding performs well against greedy decoding, with its win rate higher than the corresponding win rate for Prometheus BoN decoding vs. greedy decoding. We also find that dpo-3-MBR significantly outperforms sft, and its margin of victory is greater than the margin of victory of dpo-3-BoN vs. sft. Furthermore, the margin of improvement associated with MBR distillation is larger than the 1584 corresponding margin for MBR inference. These findings align with our findings from our automatic 1585 LLM evaluation experiments. 1586

1587

1589

1591

#### Ι ALGORITHMS

1590 L1 MBR INFERENCE

We provide the algorithm for MBR inference, complementing the mathematical overview provided 1592 in Section 2. 1593

### 1594 Algorithm 1 MBR Inference 1595

**Inputs:** Prompt x, generator model p, reference-based utility metric u, number of candidates  $N_{\text{cand}}$ , 1596 sampling temperature t. 1597 **Output:** MBR output  $\hat{y}$ 

1598

Initialise  $\mathcal{H}_{hyp} \leftarrow \emptyset$ for  $i \in \{1, 2, \dots, \widetilde{N_{\text{cand}}}\}$  do Sample  $y^{(i)} \sim p(\cdot|x)$  with temperature t Add  $y^{(i)}$  to  $\mathcal{H}_{hyp}$ // Form hypothesis set end for for  $y^{(i)} \in \mathcal{H}_{hyp}$  do *Compute*  $\tilde{u}(y^{(i)}) = \frac{1}{N_{cand}} \sum_{j=1}^{N_{cand}} u(y^{(i)}, y^{(j)})$ 1604 // Compute expected utility end for Select  $\hat{y} = \arg \max_{y \in \mathcal{H}_{hvp}} \tilde{u}(y)$ 1608 return  $\hat{y}$ 1609

1610 1611

1612

### I.2 MBR DISTILLATION

1613 In MBR distillation, we first gather a preference dataset using MBR decoding over a set of input prompts. Given input prompts  $X_k$  and a base model  $\pi_{\theta_{k-1}}$ , we first sample  $N_{\text{cand}}$  outputs per prompt 1614 1615  $y^{(i)} \sim \pi_{\theta_{k-1}}(\cdot|x)$ , where  $x \in X_k$ , forming hypothesis set  $\mathcal{H}_{hyp} = \{y^{(1)}, y^{(2)}, \dots, y^{(N_{cand})}\}$ . 1616 We then compute expected utility for elements in  $\mathcal{H}_{hyp}$ 1617 1618  $\tilde{u}(y^{(i)}) = \frac{1}{N_{\text{cand}}} \sum_{i=1}^{N_{\text{cand}}} u(y^{(i)}, y^{(j)})$ 1619 (4) and form preference pairs using the output that maximises expected utility  $\hat{y}^+$  and the output that minimises it  $\hat{y}^-$ .

 $\hat{y}^+ = rg\max_{y \in \mathcal{H}_{\mathrm{hyp}}} ilde{u}(y) \ \hat{y}^- = rg\min_{y \in \mathcal{H}_{\mathrm{hyp}}} ilde{u}(y)$ 

1627 We collect these preference pairs and their prompt  $(x, \hat{y}^+, \hat{y}^-)$  in a dataset we denote  $\mathcal{Y}_k$ .

We then use these preference pairs for DPO (Rafailov et al., 2024) training. In DPO, we minimise the following policy objective:

$$\mathcal{L}_{\text{DPO}} = -\mathbb{E}_{(x,\hat{y}^+,\hat{y}^- \sim \mathcal{Y}_k)} \left[ \log \sigma \left( \beta \log \frac{\pi_{\theta}(\hat{y}^+|x)}{\pi_{\text{ref}}(\hat{y}^+|x)} - \beta \log \frac{\pi_{\theta}(\hat{y}^-|x)}{\pi_{\text{ref}}(\hat{y}^-|x)} \right) \right]$$
(5)

where  $\pi_{\theta}$  is the policy and  $\pi_{ref}$  the reference model. We repeat this process iteratively with k = 1, 2, ..., K. The initial model  $\pi_{\theta_0}$  is the base *sft* model. We choose K = 3 in our experiments.

### Algorithm 2 MBR Distillation with DPO

1623

1624

1625 1626

1635

1659

1661 1662

Inputs: Prompt sets X<sub>1</sub>, X<sub>2</sub>,..., X<sub>K</sub>, sft model π<sub>θ0</sub>, reference-based utility metric u, number of candidates N<sub>cand</sub>, sampling temperature t, number of self-training iterations K.
 Output: Self-trained model π<sub>θK</sub>

1640		
1641	for $k \in \{1, 2, \dots, K\}$ do	
1642	Initialise $\mathcal{Y}_k \leftarrow \emptyset$	
1643	for $x \in X_k$ do	
1644	Initialise $\mathcal{H}_{hyp} \leftarrow \emptyset$	
1645	for $i \in \{1, 2, \dots, N_{ ext{cand}}\}$ do	
1646	Sample $y^{(i)} \sim \pi_{\theta_{k-1}}(\cdot x)$ with temperature $t$	
1647	Add $y^{(i)}$ to $\mathcal{H}_{\mathrm{hyp}}$	// Form hypothesis set
1648	end for	
16/0	for $y^{(i)} \in \mathcal{H}_{ ext{hyp}}$ do	
1650	Compute $\tilde{u}(y^{(i)}) = \frac{1}{N-i} \sum_{i=1}^{N_{\text{cand}}} u(y^{(i)}, y^{(j)})$	// Compute expected utility
1651	end for	
1652	Select $\hat{y}^+ = \arg \max_{y \in \mathcal{H}_{\text{byp}}} \tilde{u}(y)$	// Select highest scoring output
1653	Select $\hat{y}^- = \arg\min_{y \in \mathcal{H}_{hom}} \tilde{u}(y)$	// Select lowest scoring output
1654	Add $(\hat{y}^+, \hat{y}^-)$ to $\mathcal{Y}_k$	// Form preference pairs
1655	end for	L L
1656	Update $\pi_{\theta_k} \leftarrow DPO(\pi_{\theta_{k-1}}, \mathcal{Y}_k)$	// DPO training on preference pairs
1657	end for	
1658	return $\pi_{\theta_K}$	

### J LIMITATIONS

While our work demonstrates the significant potential of MBR decoding, there are limitations that 1663 should be addressed in future research. Firstly, although we demonstrate using existing judge LLMs 1664 utility metrics that MBR decoding consistently outperforms BoN decoding, this does not preclude 1665 the existence of reference-free metrics that *are* powerful enough to match or surpass the performance of their direct reference-based counterparts. This relates to a possible broader limitation on the benefits of using consensus quality for output selection, as the consensus solution may not always be the optimal one. We encourage future work to train better utility metrics in order to better understand 1669 these limitations. A second limitation of our work is that we do not study the biases introduced by 1670 the utility metric. One particularly pernicious form of bias is "reward-hacking" behavior, where the utility metric (likely as a result of its own training) selects outputs that evaluate well on our 1671 benchmarks but that are actually worse in quality. While we preclude this from being the case in our 1672 experiments via our human study (Appendix H), this does not mean that such pernicious behavior 1673 cannot arise in other settings. Finally, we do not study limitations on scalability. Although we show that small judge LLMs (7B) can serve as utility metrics for much larger models (70B), it is likely that
 weaker utility metrics cease being useful for very strong LLMs and on very complex tasks. Further
 research is needed to determine when this breakdown occurs.

