

THEORETICAL ASPECTS OF BIAS AND DIVERSITY IN MINIMUM BAYES RISK DECODING

Anonymous authors

Paper under double-blind review

ABSTRACT

Text generation commonly relies on greedy and beam decoding that limit the search space and degrade output quality. Minimum Bayes Risk (MBR) decoding can mitigate this problem by utilizing automatic evaluation metrics and model-generated pseudo-references. Previous studies have conducted empirical analyses to reveal the improvement by MBR decoding, and reported various observations. However, despite these observations, the theoretical relationship between them remains uncertain. To address this, we present a novel theoretical interpretation of MBR decoding from the perspective of bias-diversity decomposition. We decompose errors in the estimated quality of generated hypotheses in MBR decoding into two key factors: *bias*, which reflects the closeness between utility functions and human evaluations, and *diversity*, which represents the variation in the estimated quality of utility functions. Our theoretical analysis reveals the difficulty in simultaneously improving both bias and diversity, and highlights the effectiveness of increasing diversity to enhance MBR decoding performance. This analysis verifies the alignment between our theoretical insights and the empirical results reported in previous work. Furthermore, to support our theoretical findings, we propose a new metric, pseudo-bias, which approximates the bias term using gold references. We also introduce a new MBR approach, Metric-augmented MBR (MAMBR), which increases diversity by adjusting the behavior of utility functions without altering the pseudo-references. Experimental results across multiple NLP tasks show that the decomposed terms in the bias-diversity decomposition correlate well with performance, and that MAMBR improves text generation quality by modifying utility function behavior. Our code will be available at [https://github.com/\[Anonymized\]](https://github.com/[Anonymized]).

1 INTRODUCTION

As demonstrated by the success of large language models (LLMs) (Brown et al., 2020; OpenAI et al., 2024), text generation is one of the most fundamental tasks in Natural Language Processing (NLP). Text generation commonly relies on greedy and beam searches, which heavily restrict the search space when decoding texts from a model. This procedure, which only considers the model’s predictions within a limited search space, can sometimes degrade the quality of the generated text.

Minimum Bayes Risk (MBR) decoding (Goel & Byrne, 2000) can mitigate this problem by using a utility function, essentially an automatic evaluation metric, along with pseudo-references generated by the model. MBR decoding was initially applied to speech recognition (Goel & Byrne, 2000) and later to statistical machine translation (SMT) (Kumar & Byrne, 2002; 2004; Duan et al., 2011). Following these successes, MBR decoding has been expanded to various text generation tasks, including neural machine translation (NMT) (Stahlberg et al., 2017), text summarization (Bertsch et al., 2023), and image captioning (Borgeaud & Emerson, 2020).

Since MBR decoding has become an important inference technique in text generation, various empirical studies have explored its characteristics. Müller & Sennrich (2021); Freitag et al. (2022a); Fernandes et al. (2022); Amrhein & Sennrich (2022) highlight the importance of using high quality evaluation metrics that is robust and correlate well with human evaluations as utility functions. Jinnai et al. (2024a); Heineman et al. (2024) emphasize the importance of high-quality pseudo-references that closely resemble human-created ones, while also stressing the significance of pseudo-reference

054 diversity. Although these empirical findings cover various aspects in detail, a unified interpretation
 055 remains challenging due to the lack of theoretical frameworks explaining the relationships behind
 056 them.

057 To address this gap, we provide theoretical interpretations of MBR decoding through bias-diversity
 058 decomposition (Krogh & Vedelsby, 1994; Wood et al., 2024). Our theoretical interpretation focuses
 059 on errors in the estimated quality of hypotheses in MBR decoding. These errors are decomposed
 060 into two critical factors: *bias* and *diversity*. The bias term represents the closeness between the
 061 estimated quality produced by utility functions and human evaluations. The diversity term reflects
 062 the variance in the estimated quality across different utility functions. Based on this interpretation,
 063 we theoretically demonstrate the difficulty in improving both the bias and diversity terms simulta-
 064 neously, and we highlight the effectiveness of increasing diversity in MBR decoding, verifying the
 065 correspondence with empirically induced results from previous work.

066 To empirically verify our theoretical findings, we propose *pseudo-bias*, which approximates the bias
 067 term using gold references. Furthermore, to explore the potential for increasing the diversity term
 068 by adjusting the behavior of utility functions without varying the pseudo-references, we introduce a
 069 new MBR approach: Metric-augmented MBR (MAMBR).

070 Our empirical analysis on machine translation, text summarization, and image captioning—using
 071 pseudo-references generated by five different sampling methods—shows that the decomposed bias
 072 and diversity terms correlate with performance, consistent with our theoretical analysis. Moreover,
 073 using MAMBR demonstrates performance improvements by simply modifying the behavior of util-
 074 ity functions.

076 2 MINIMUM BAYES RISK (MBR) DECODING

077 Minimum Bayes Risk (MBR) decoding (Goel & Byrne, 2000) estimates the quality of each hypoth-
 078 esis h in a set \mathcal{H} by comparing it with each pseudo-reference (evidence sample) y in a set of all
 079 sequences Ω and its model predicted probability $P(y|x)$ for a given input sequence x . This pro-
 080 cess uses an evaluation metric, treated as a utility function $f_\theta(h, y)$, which calculates the similarity
 081 between h and y to choose the best hypothesis \hat{h}_{best} in \mathcal{H} as follows:

$$082 \hat{h}_{best} = \arg \max_{h \in \mathcal{H}} \sum_{y \in \Omega} f_\theta(h, y) P(y|x), \quad (1)$$

083 where θ represents the parameters of the evaluation metric used as the utility function $f_\theta(h, y)$.
 084 Since calculating Ω is intractable, Eikema & Aziz (2020; 2022) replace Ω with $|\mathcal{Y}|$, a set of sampled
 085 y as follows:

$$086 \hat{h}_{mbr} = \arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_\theta(h, y), \quad y \sim P(y|x). \quad (2)$$

087 Here, instead of using the utility function $f_\theta(h, y)$, we can assume human-estimated quality (Naskar
 088 et al., 2023; Suzgun et al., 2023; Jinnai et al., 2024a; Ohashi et al., 2024) denoted as $\hat{f}_\theta(h)$. Under
 089 this assumption, the ideal decoding, which fully relies on human-estimated results, is represented as
 090 follows:

$$091 \hat{h}_{human} = \arg \max_{h \in \mathcal{H}} \hat{f}_\theta(h). \quad (3)$$

092 In this paper, we focus on analyzing the differences between the internally estimated qualities for
 093 each hypothesis by MBR decoding and those estimated by humans to better understand the charac-
 094 teristics of MBR decoding (§3).

095 3 THEORETICAL ANALYSIS BASED ON BIAS-DIVERSITY DECOMPOSITION

096 3.1 EVALUATION DISCREPANCY

097 To measure the discrepancy between the human estimated quality, $\hat{f}_\theta(h)$ and the MBR decoding es-
 098 timated quality, $\frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_\theta(h, y)$, we define a $|\mathcal{H}|$ -dimensional vector \mathbf{u}^j that represents estimated
 099

quality for each hypothesis based on the j -th pseudo-reference and also define $\bar{\mathbf{u}}$, the average vector of all \mathbf{u}^j as follows:

$$\mathbf{u}^j = \begin{bmatrix} u_1^j \\ \cdots \\ u_{|\mathcal{H}|}^j \end{bmatrix}, \quad u_i^j = f_\theta(h_i, y_j), \quad \bar{\mathbf{u}} = \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} \mathbf{u}^j. \quad (4)$$

Similarly, we can define a $|\mathcal{H}|$ -dimensional vector, $\hat{\mathbf{u}}$ that represents the human estimated quality for each hypothesis as follows:

$$\hat{\mathbf{u}} = \begin{bmatrix} \hat{u}_1 \\ \cdots \\ \hat{u}_{|\mathcal{H}|} \end{bmatrix}, \quad \hat{u}_i = \hat{f}_\theta(h_i). \quad (5)$$

Here, by using Equations 4 and 5, we can reformulate MBR decoding in Equation 2 and the ideal decoding in Equation 3 as follows:

$$(2) \equiv \hat{h}_{mbr} = \arg \max_{h_i} \bar{u}_i, \quad (3) \equiv \hat{h}_{human} = \arg \max_{h_i} \hat{u}_i. \quad (6)$$

Therefore, based on Equation 6, we can investigate the discrepancy between the estimated quality by MBR decoding and human through the comparison of $\bar{\mathbf{u}}$ and $\hat{\mathbf{u}}$. In our work, to estimate the discrepancy, we consider the prediction error of $\bar{\mathbf{u}}$ to $\hat{\mathbf{u}}$ by using Mean Squared Error (MSE) as follows:

$$MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}}) = \frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} (\hat{u}_i - \bar{u}_i)^2. \quad (7)$$

3.2 BIAS-DIVERSITY DECOMPOSITION

Our goal is to reveal the characteristics of MBR decoding through theoretical analysis. To achieve this, we focus on the bias and diversity underlying Equation 7. Based on this approach, we can induce the following decomposition:

Theorem 1. *The quality estimation error for each hypothesis in MBR decoding, $(\hat{u}_i - \bar{u}_i)^2$, can be decomposed to bias and diversity (ambiguity) terms (Krogh & Vedelsby, 1994) as follows:*

$$(\hat{u}_i - \bar{u}_i)^2 = \underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_\theta(h_i, y_j))^2}_{\text{Bias}} - \underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i - f_\theta(h_i, y_j))^2}_{\text{Diversity}}. \quad (8)$$

Proof. See Appendix A. □

In Equation 8, two terms represent bias and diversity. **Unlike the well-known bias-variance decomposition (Geman et al., 1992) that targets a single estimator which is \mathbf{u} in our case¹, the second term is negative, which is why it is referred to as diversity rather than variance (Wood et al., 2024).** The bias term indicates how closely the utility function’s estimated quality for a hypothesis matches human estimation. The diversity term reflects how different the utility function’s estimated qualities are for each other. This decomposition emphasizes the importance of increasing the diversity term while reducing the bias term to improve the quality estimation error, $(\hat{u}_i - \bar{u}_i)^2$, for each hypothesis.

While MBR decoding considers all hypotheses to rank and select the best one, Theorem 1 addresses only the quality estimation for each hypothesis. To bridge this gap, we decompose $MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}})$ that accounts for all hypotheses. The following theorem addresses this broader perspective:

Theorem 2. *The quality estimation error for all hypotheses in MBR decoding, $MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}})$, can be decomposed into bias and diversity terms for all hypotheses as follows:*

$$MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}}) = \underbrace{\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_\theta(h_i, y_j))^2}_{\text{Bias for all hypotheses}} - \underbrace{\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i - f_\theta(h_i, y_j))^2}_{\text{Diversity for all hypotheses}}. \quad (9)$$

¹This becomes $MSE(\hat{\mathbf{u}}, \mathbf{u}) = \frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} ((\hat{u}_i - (\frac{1}{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} u_j))^2 + ((\frac{1}{|\mathcal{H}|} \sum_{j=1}^{|\mathcal{H}|} u_j) - u_i)^2)$.

162 *Proof.* See Appendix B. □

163
164
165 As in Theorem 1, Theorem 2 highlights the importance of increasing the diversity term while de-
166 creasing the bias term to improve the quality estimation in MBR decoding.

168 3.3 INTERPRETATION

169
170 The decompositions presented in Theorems 1 and 2 allow us to provide theoretical interpretations
171 for the empirically analyzed characteristics of MBR decoding and its extensions discussed in prior
172 studies.

174 3.3.1 CORRELATION TO HUMAN EVALUATION RESULTS

175
176 The bias term of the decomposition, $(\hat{u}_i - f_\theta(h_i, y_j))^2$, highlights the importance of considering the
177 closeness between the human-estimated quality, \hat{u}_i , and the quality estimated by the utility function,
178 $f_\theta(h_i, y_j)$, for improving the performance of MBR decoding. Specifically, since the utility function,
179 $f_\theta(h_i, y_j)$, is influenced by the pseudo-reference y_j , the bias term underscores the significance of
180 considering the utility function’s correlation to human evaluation and the closeness between pseudo-
181 references and human-created references. Therefore, it emphasizes the importance of examining
182 both utility functions and sampling strategies for generating pseudo-references.

183 **Quality of Evaluation Metrics.** Our theoretical insight is supported by empirical findings (Müller
184 & Sennrich, 2021; Freitag et al., 2022a; Fernandes et al., 2022; Amrhein & Sennrich, 2022), in which
185 the quality of evaluation metrics used as utility functions is crucial for performance improvement.

186 **Quality of Pseudo-References.** Ohashi et al. (2024); Jinnai et al. (2024a) empirically show the
187 importance of selecting appropriate pseudo-references. Thus, our findings theoretically support
188 these empirically driven insights of these previous studies.

189 **Challenges in the Real World.** Our theoretical findings emphasize the necessity of directly reduc-
190 ing the bias term. However, this requires human evaluation of the combination of pseudo-references
191 and evaluation metrics, used as utility functions, for each hypothesis. This task is clearly challenging
192 due to the high cost of human evaluation. As a solution, we propose a method to approximate this
193 in §4.1 and evaluated its correlation with task-specific performance in §5.3.

195 3.3.2 DIVERSITY OF AUTOMATIC EVALUATION RESULTS

196
197 Based on the diversity term of the decomposition, increasing diversity can contribute to performance
198 improvements by reducing each prediction error $(\hat{u}_i - \bar{u}_i)^2$ of $MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}})$ in MBR decoding. A key
199 insight here is that the diversity expressed by $(\bar{u}_i - f_\theta(h_i, y_j))^2$ stems from the different estimated
200 qualities produced by each utility function $f_\theta(h_i, y_j)$. Thus, this diversity can be influenced by the
201 pseudo-reference y_j and/or the model parameters θ of the evaluation metric.

202 **Diversity of Pseudo-references.** This finding supports the previous studies (Freitag et al., 2023a;
203 Jinnai et al., 2024a; Heineman et al., 2024) that conclude the diversity of sampling methods is es-
204 sential for performance improvement of MBR decoding considering that the diversity of the pseudo-
205 references can indirectly contribute to increasing the diversity of $f_\theta(h_i, y_j)$ by each y_j . **Basically,**
206 **as Jinnai et al. (2024a) introduce, reranking algorithms are more effective when the candidates are**
207 **diverse (Gimpel et al., 2013; Li & Jurafsky, 2016; Li et al., 2016) owing to their diverse information**
208 **to make a consensus.**

209 **Diversity of Evaluation Metrics.** We anticipate performance improvements by combining multiple
210 different evaluation metrics as utility functions to increase diversity. While this approach has been
211 shown to improve the quality estimation of generated texts (Glushkova et al., 2023), to the best of
212 our knowledge, it has not yet been applied to MBR decoding.

213 **Unexplored Aspect.** Furthermore, the effect of increasing the diversity of estimated qualities from
214 utility functions by varying the evaluation metric’s model parameters θ remains uncertain. To inves-
215 tigate this, we propose a method to adjust the diversity of estimated qualities by modifying θ in §4.2
and compare its behavior with that of varying pseudo-references in §5.4.

3.3.3 MBR DECODING AS ENSEMBLE LEARNING

Our decomposition of MBR decoding aligns with ensemble learning, which is induced by Krogh & Vedelsby (1994). Thus, we can understand that the quality estimation by MBR decoding is a kind of ensemble learning.

Quality Estimation. Our decomposition starts from the definition $MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}})$ in Eq. 7, the error between the estimated qualities from human evaluation and MBR decoding. We can actually observe the reduction of errors as the improvement in quality score estimation of (Naskar et al., 2023; Cheng & Vlachos, 2024) by ensembling utility functions that are similar to MBR decoding.

Weighted-voting. Furthermore, this viewpoint supports the validity of the previous work (Suzgun et al., 2023; Bertsch et al., 2023) that shows the interpretation of MBR decoding as soft-weighted voting, a variant of ensemble learning. Different from our setting, soft-weighted voting restricts the value range of voters (utility functions) from 0 to 1. Wood et al. (2024) shows that soft-weighted voting can be converted to the decomposition of Krogh & Vedelsby (1994), equivalent to our decomposition in Eq. 8. Therefore, weighted voting-based MBR decoding is similarly explained in our decomposition.

Number of Pseudo-references. Generally, increasing the number of pseudo-references improves performance but demands additional computational cost. DeNero et al. (2009); Eikema & Aziz (2022); Cheng & Vlachos (2023); Deguchi et al. (2024b); Vamvas & Sennrich (2024); Trabelsi et al. (2024) prune samples to speedup inference and maintain the original quality similar to the case of pruning estimators in ensemble learning (Liu et al., 2004; Bonab & Can, 2016; 2019).

Considering an ensemble learning method, such as the Bayes optimal classifier (Mitchell, 1997), and assuming that Eq. 2 approximates the expectation by sampling y_j , we can explain the performance improvement of increased pseudo-references by the law of large numbers and the success of the pruning and weighted utility functions (Jinnai et al., 2024b) through importance sampling (Kloek & Van Dijk, 1978). (See Appendix C for more details.)

3.3.4 BIAS AND DIVERSITY TRADE-OFF

At first glance, based on the interpretation in §3.3.1 and §3.3.2, decreasing bias while increasing diversity seems to be the best strategy to improve quality estimation performance in MBR decoding, which was investigated by Jinnai et al. (2024a). To understand the validity of this strategy, we need to focus on the bias-diversity trade-off (Krogh & Vedelsby, 1994).

Limitation of MBR Decoding. The bias-diversity trade-off highlights the difficulty of decreasing bias while increasing diversity. In Eqs. 8 and 9, when the bias term approaches zero, i.e., each $f_\theta(h_i, y_j)$ approaching to \hat{u}_i , the diversity term also approaches zero. This theoretical fact indicates that even if we can prepare high-quality evaluation metrics and high-quality pseudo-references that correlate well with human behavior, there may be no performance improvement due to diminished diversity.

Diversity Assists Inferior Methods. Conversely, when the evaluation metrics and pseudo-references are inferior, we can expect performance improvements through increased diversity at the cost of increased bias. This phenomenon can explain the sometimes competitive performance of BLEU (Papineni et al., 2002) against COMET (Rei et al., 2020) in Freitag et al. (2022b), and that of ancestral sampling (Robert, 1999) against other sampling methods in Freitag et al. (2023a); Ohashi et al. (2024) using MBR decoding. However, unlike the case where decreased bias leads to decreased diversity, increased bias does not guarantee increased diversity. Therefore, we must carefully assess their diversity when using low-quality evaluation metrics and pseudo-references in MBR decoding.

4 REMAINING PROBLEMS & SOLUTIONS

Our theoretical analysis covers various aspects of MBR decoding. However, for a comprehensive analysis, we should investigate empirical results not addressed in previous work and bridge the gap between theory and real-world applications. To this end, we provide the following solutions.

4.1 PSEUDO-BIAS

As discussed in §3.3.1, the bias term suggests the importance of considering the correlation between the results of human evaluation and the evaluation metric’s decisions based on pseudo-references to improve the performance of MBR decoding. However, calculating the bias term requires human evaluation, and conducting human evaluations for each setting is unrealistic and difficult. To address this issue, we introduce *pseudo-bias*, an approximation of the bias term in our decomposition. By using $|\hat{\mathcal{Y}}|$, the number of gold references \hat{y} , pseudo-bias is defined as follows:

$$\underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_{\theta}(h_i, y_j))^2}_{\text{Bias}} \approx \underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\tilde{u}_i - f_{\theta}(h_i, y_j))^2}_{\text{Pseudo-bias}}, \quad \tilde{u}_i = \frac{1}{|\hat{\mathcal{Y}}|} \sum_{j=1}^{|\hat{\mathcal{Y}}|} f_{\theta}(h_i, \hat{y}_j). \quad (10)$$

This formulation is based on the premise that automatic evaluation metrics correlate to human evaluation when receiving human-created references.² Since we can calculate the diversity term without any approximation, we compare pseudo-bias with diversity in terms of how they correlate with performance.

4.2 METRIC-AUGMENTED MBR

The discussion in §3.3.2 shows the possibility of increasing the diversity of the utility function, $f_{\theta}(h_i, y_j)$, by changing the evaluation metric’s model parameters, θ , as well as by introducing diversity through pseudo-references. To this end, we propose a new method called Metric-augmented Minimum Bayes Risk (MAMBR) decoding. In MAMBR, we employ different parameters for the evaluation metric to enhance the diversity of utility functions. Letting Θ be a set of model parameters, MAMBR is defined as follows:

$$\hat{h}_{\text{mambr}} = \arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}| |\Theta|} \sum_{\theta \in \Theta} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y). \quad (11)$$

When using MAMBR, we train evaluation metrics with different initial random seeds to generate Θ as a set of diverse model parameters.

5 EMPIRICAL ANALYSIS

We conduct empirical analysis corresponding to our theoretical analysis through experiments for a comprehensive understanding of MBR decoding.

5.1 OVERALL SETTINGS

We target three different text generation tasks, machine translation, text summarization, and image captioning to investigate the general performance of MBR decoding. In all tasks, we followed the settings of Jinnai et al. (2024b) for generating samples. We used epsilon sampling (Hewitt et al., 2022) to generate hypotheses.³ For the generation of pseudo-references, we used various sampling approaches: beam decoding, nucleus sampling (Holtzman et al., 2020) with $p = 0.9$, ancestral sampling, top- k sampling (Fan et al., 2018) with $k = 10$, and epsilon sampling with $\epsilon = 0.02$. We set the sampling size for hypotheses to 64. We chose the sampling size for pseudo-references from $\{4, 8, 16, 32, 64\}$. We used the following datasets, models⁴, and evaluation metrics for each task:

Machine Translation We used the WMT19 English to German (En-De) and WMT19 English to Russian (En-Ru) datasets (Barrault et al., 2019). We used `facebook/wmt19-en-de` for En-De

²For the pseudo-bias, we used COMET (Unbabel/wmt22-comet-da) and BERTScore with `microsoft/deberta-xlarge-mnli` whose pearson correlations are 0.990 on the system-level task for English to German (Freitag et al., 2023b) and 0.7781 (https://github.com/Tiiiger/bert_score) on WMT16 to English (Bojar et al., 2016), respectively.

³Appendix G.1 includes the results with hypotheses generated by different sampling methods.

⁴We used all models from <https://huggingface.co/models> (Wolf et al., 2020).



Figure 1: Correlation between measures in our decomposition and performance for each dataset. The underlined scores indicate statistically significant results ($p < 0.05$).⁷Note that the italic scores at Avg. are not the target of the significance test.

and facebook/wmt19-en-ru for En-Ru, respectively. As the utility function and evaluation metric, we used COMET with the model Unbabel/wmt22-comet-da.

Text Summarization We used the SAMSum (Gliwa et al., 2019) and XSum (Narayan et al., 2018) datasets, and used philschmid/bart-large-cnn-samsum and facebook/bart-large-xsum for generation in SAMSum and XSum, respectively. As the utility function and evaluation metric, we used BERTScore (Zhang* et al., 2020) with the model microsoft/deberta-xlarge-mnli.

Image Captioning We used the MSCOCO dataset (Lin et al., 2014) with the split of Karpathy & Fei-Fei (2015) and the NoCaps dataset (Agrawal et al., 2019). We used Salesforce/blip2-flan-t5-xl-coco and Salesforce/blip2-flan-t5-xl for generation in MSCOCO and NoCaps, respectively. As the utility function and evaluation metric, we used BERTScore with the model microsoft/deberta-xlarge-mnli. Since both datasets have multiple references, we report their average scores.

Our implementation of the generation part is based on the released code of Jinnai et al. (2024b)⁵ and the MBR decoding part is based on the toolkit, mbrs by Deguchi et al. (2024a)⁶. We generate samples on NVIDIA GeForce RTX 3090 and perform MBR decoding on an NVIDIA RTX A6000.

5.2 CORRELATION OF BIAS AND DIVERSITY TO PERFORMANCE

To verify our theoretical decomposition, we investigate the correlation of bias and diversity to performance on each dataset. For this purpose, we approximately compute the bias term by using our pseudo-bias in §4.1. Furthermore, we investigate the importance of whether to consider the entire candidate (Eq. 9) or one best candidate (Eq. 8).

Settings We compared the following measures in our decomposition: OVERALL BIAS, the first term of Eq. 9 for all hypotheses; ONE BEST BIAS, the first term of Eq. 8 for the one best result by MBR decoding; OVERALL DIVERSITY, the second term of Eq. 9 for all hypotheses; ONE BEST DIVERSITY, the second term of Eq. 8 for the one best result by MBR decoding; OVERALL MSE, indicating errors for all hypotheses based on Eq. 9; and ONE BEST MSE, indicating errors for the one best result by MBR decoding based on Eq. 8. For the comparison, we calculated Spearman’s rank correlation and Pearson correlation between these measures and the performance based on the results of five different sampling methods with five different sampling sizes on each dataset (see §5.1). Since lower bias and MSE are better for performance, we took their negative values in the correlation calculation. Moreover, we report averaged correlation across all datasets by Fisher z-transformation (Corey et al., 1998).

Results Figure 1 shows the correlation between the measures and performance for each dataset. These results show that MSE for both overall and one best results correlates well with the performance for each dataset in Spearman’s rank correlation, indicating the importance of considering quality estimation in MBR decoding, as in Eqs. 8 and 9. On the other hand, the decomposed bias

⁵<https://github.com/CyberAgentAILab/model-based-mbr>

⁶<https://github.com/naist-nlp/mbrs>

⁷We used Student’s t-test (Student, 1908) for both spearman and pearson correlations.

378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431

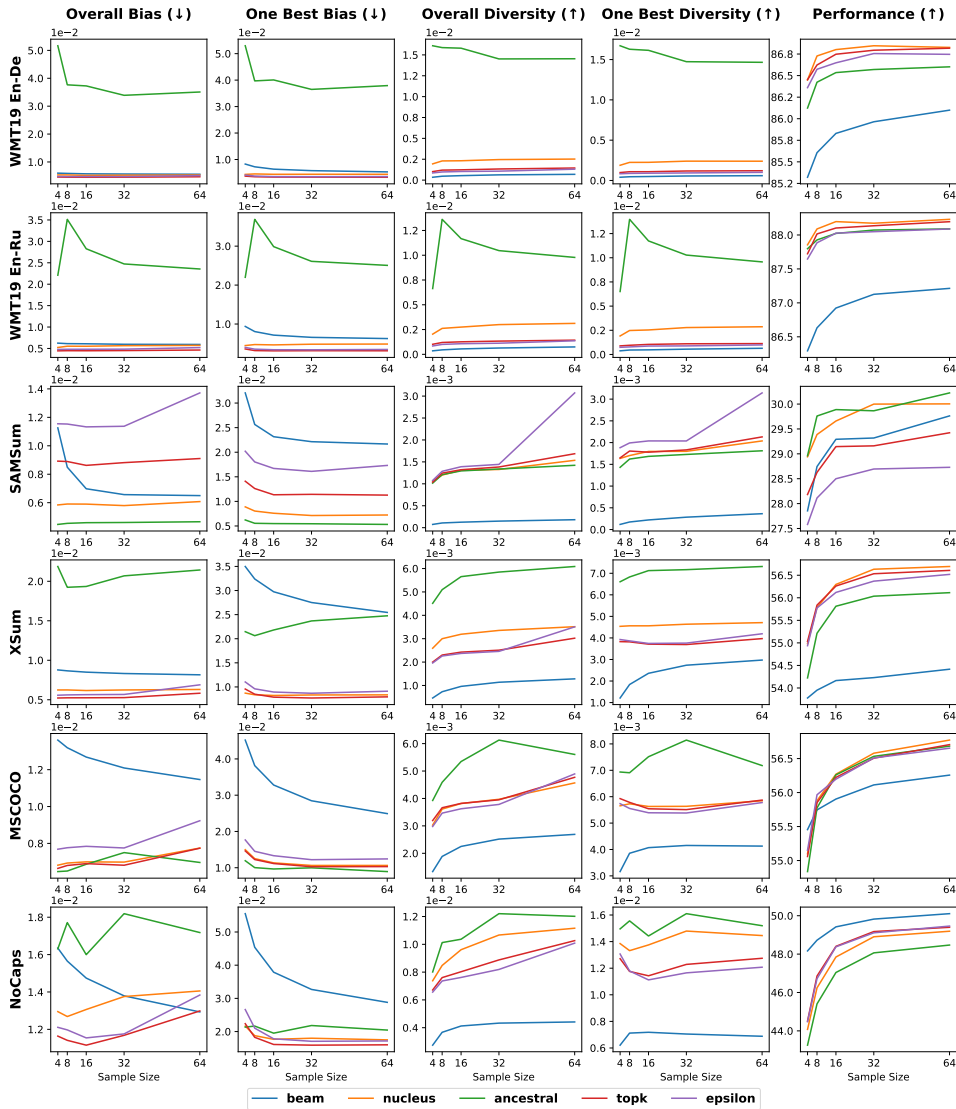


Figure 2: The relationship between bias, diversity, and performance in MBR decoding. The x-axis shows the number of used pseudo-references. (↑) indicates higher scores are better whereas (↓) indicates lower scores are better.

and diversity show different tendencies. ONE BEST BIAS, which considers the one best result, is important for bias, whereas OVERALL DIVERSITY, which considers overall results, is important for diversity. This result is reasonable given an assumption that MBR decoding aims to select texts that are close to human-created ones. Based on this assumption, we can say that diversity supports the selection by considering the importance of all hypotheses not covered by One Best Bias. In contrast to the results in Spearman’s rank correlation, the coefficients of Pearson’s correlation decrease. Based on these results, we can conclude that the measures, i.e., ONE BEST BIAS, OVERALL DIVERSITY, ONE BEST MSE and OVERALL MSE, correlate well with the rank in performance, but they are difficult to precisely capture subtle differences of values. (See Appendix D for further details.)

5.3 BIAS AND DIVERSITY TRADE-OFF

To investigate the bias-diversity trade-off in more detail, we followed the setup described in §5.2. We plotted the results for each dataset using different sampling methods in Figure 2. The results show that while ancestral sampling exhibits the highest bias, except in the case of the SAMSum

Table 1: Results of MAMBR with ancestral sampling. Bold font indicates the best result.

		WMT19 En-De					WMT19 En-Ru					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		85.7	85.9	85.9	85.9	85.9	87.4	87.4	87.5	87.5	87.5	
Num. of Models		2	85.7	86.0	86.0	85.9	85.9	87.4	87.4	87.5	87.5	87.6
4		85.8	86.0	86.0	86.0	86.0	87.4	87.4	87.5	87.5	87.6	
8		85.8	86.0	86.0	86.0	86.1	87.4	87.5	87.6	87.6	87.6	
		SAMSum					XSum					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		28.6	29.1	29.5	29.5	29.7	54.2	55.2	55.7	56.0	56.1	
Num. of Models		2	28.8	29.6	29.9	29.9	30.1	54.2	55.2	55.7	56.1	56.2
4		28.7	29.5	29.9	29.8	30.2	54.2	55.2	55.8	56.1	56.2	
8		28.7	29.5	29.8	29.9	30.1	54.3	55.3	55.8	56.1	56.2	
		MSCOCO					NoCaps					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		54.9	55.8	56.3	56.5	56.8	42.9	45.3	46.8	47.8	48.6	
Num. of Models		2	54.9	55.8	56.4	56.6	56.8	43.2	45.6	47.2	48.3	48.9
4		54.9	56.0	56.4	56.7	56.9	43.3	45.6	47.3	48.4	49.0	
8		54.9	56.0	56.5	56.8	56.9	43.5	45.7	47.4	48.5	49.0	

Table 2: Results of MAMBR with epsilon sampling. Notations are the same as Table 1.

		WMT19 En-De					WMT19 En-Ru					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		85.9	86.1	86.2	86.2	86.2	87.3	87.6	87.7	87.7	87.7	
Num. of Models		2	86.0	86.1	86.1	86.2	86.3	87.3	87.6	87.7	87.7	87.7
4		86.0	86.1	86.2	86.2	86.3	87.4	87.6	87.7	87.7	87.7	
8		86.0	86.2	86.3	86.3	86.4	87.4	87.7	87.7	87.7	87.8	
		SAMSum					XSum					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		27.5	27.9	28.3	28.4	28.5	54.9	55.7	56.1	56.3	56.4	
Num. of Models		2	27.7	28.1	28.5	28.6	28.7	54.9	55.7	56.1	56.3	56.5
4		27.7	28.2	28.5	28.6	28.6	54.9	55.7	56.1	56.4	56.5	
8		27.6	28.2	28.6	28.6	28.7	54.9	55.8	56.2	56.4	56.5	
		MSCOCO					NoCaps					
Num. of Samples		4	8	16	32	64	4	8	16	32	64	
1		55.2	55.9	56.3	56.5	56.7	44.4	46.7	48.5	49.1	49.5	
Num. of Models		2	55.2	55.9	56.3	56.5	56.7	44.4	46.8	48.6	49.2	49.6
4		55.2	56.0	56.3	56.6	56.8	44.5	46.9	48.7	49.3	49.7	
8		55.3	56.1	56.3	56.6	56.8	44.6	47.0	48.7	49.4	49.7	

dataset, it sometimes outperforms other sampling methods owing to its greater diversity. Focusing on top-k sampling, which has the lowest bias, again excluding the SAMSum dataset, we can observe that the reduction in bias tends to limit the increase in diversity. This finding supports our previously noted bias-diversity trade-off in MBR decoding. However, as evidenced by the performance of beam decoding, which has the lowest diversity, the importance of bias and diversity varies depending on the target dataset. Therefore, while our theoretical analysis effectively explains the performance tendencies in MBR decoding, it remains essential to consider task-specific features carefully to achieve further performance improvements. (See Appendix E for further details.)

5.4 EFFECTIVENESS OF METRIC-AUGMENTED MBR

We investigate the possibility of improving performance of MAMBR in Eq. 11 by changing automatic evaluation metric’s model parameters.

Table 3: Results of MAMBR with beam decoding. Notations are the same as Table 1.

		WMT19 En-De					WMT19 En-Ru				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	85.2	85.4	85.6	85.7	85.8	86.5	86.8	87.0	87.1	87.1
Num. of Models	2	85.3	85.5	85.7	85.8	85.8	86.5	86.8	86.9	87.1	87.1
	4	85.3	85.5	85.7	85.8	85.8	86.6	86.9	87.0	87.1	87.2
	8	85.3	85.5	85.7	85.8	85.9	86.5	86.8	87.0	87.1	87.2
		SAMSum					XSum				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	27.6	28.7	29.2	29.3	29.7	53.8	54.0	54.2	54.2	54.4
Num. of Models	2	27.8	28.9	29.2	29.4	29.7	53.8	53.9	54.1	54.2	54.4
	4	27.8	28.9	29.2	29.4	29.7	53.8	54.0	54.2	54.2	54.4
	8	27.8	28.9	29.2	29.4	29.7	53.8	54.0	54.2	54.2	54.4
		MSCOCO					NoCaps				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	55.4	55.7	55.9	56.1	56.3	48.2	48.8	49.4	49.9	50.2
Num. of Models	2	55.4	55.8	55.9	56.1	56.3	48.2	48.8	49.4	49.9	50.3
	4	55.5	55.7	55.9	56.1	56.3	48.2	48.8	49.5	49.9	50.2
	8	55.5	55.7	55.9	56.1	56.3	48.2	48.8	49.5	49.9	50.2

Settings To prepare the set of model parameters, we trained eight models by varying their initial seeds. We trained `Unbabel/wmt22-comet-da` on the Direct Assessments (DA) task (Graham et al., 2013), using the WMT 2017 to 2020 datasets (Bojar et al., 2017; 2018; Barrault et al., 2019; 2020) for training and the WMT 2021 dataset (Akhbardeh et al., 2021) for validation in COMET. Additionally, we trained `microsoft/deberta-large` on the MNLI dataset from GLUE (Wang et al., 2018) for BERTScore. During inference, to control model diversity, we selected the top- n models based on their proximity to the median validation scores, with n chosen from 1, 2, 4, 8. We aimed to explore the relationship between the diversity caused by pseudo-references and model parameters, focusing on sampling strategies with varying levels of diversity: ancestral sampling, epsilon sampling, and beam decoding.

Results Tables 1, 2, and 3 respectively show the MAMBR results with pseudo-references from ancestral sampling, epsilon sampling, and beam decoding. In ancestral and epsilon sampling, the best and moderately diversified sampling strategies (as shown in Figure 2), we observe performance improvement as the number of models increases. On the other hand, in the lowest diversity method, beam decoding, performance improvement is limited. **These results suggest that MAMBR can improve performance by enhancing the diversity of evaluation metrics, although the diversity of the sampling strategy itself remains important.** (See Appendix F for further details.)

6 CONCLUSION

This work provides a unified theoretical interpretation of Minimum Bayes Risk (MBR) decoding through the lens of bias-diversity decomposition. By decomposing the errors in quality estimation in MBR decoding into bias and diversity, we highlight the trade-off between improving these two factors, with an emphasis on the benefits of increasing diversity. Our theoretical insights align with previous empirical results, and we further investigate aspects not covered by these empirical findings through the introduction of the pseudo-bias metric and the development of the Metric-augmented MBR (MAMBR) approach. Experimental results across multiple tasks demonstrate the validity of our theoretical findings and the effectiveness of our approach in improving text generation quality. These findings bridge the gap between empirical observations and theoretical understanding of MBR decoding, offering new insights for optimizing text generation. (See Appendix H for the limitations of our work.)

7 REPRODUCIBILITY STATEMENT

We performed our experiments by running publicly available models, `facebook/wmt19-en-de`, `facebook/wmt19-en-ru`, `philschmid/bart-large-cnn-samsum`, `facebook/bart-large-xsum`, `Salesforce/blip2-flan-t5-xl-coco`, and `Salesforce/blip2-flan-t5-xl` in HuggingFace Transformers (Wolf et al., 2020) on the publicly available datasets, WMT19 English to German (Barrault et al., 2019), WMT19 English to Russian (Barrault et al., 2019), SAMSum (Gliwa et al., 2019), XSum (Narayan et al., 2018), MSCOCO (Lin et al., 2014; Karpathy & Fei-Fei, 2015), and NoCaps (Agrawal et al., 2019), respectively with utilizing the publicly available MBR decoding toolkit, `mbrs` (Deguchi et al., 2024a) as described in §5.1. In addition, we will release our code at <https://github.com/naist-nlp/mbr-bias-diversity>.

REFERENCES

- Harsh Agrawal, Karan Desai, Yufei Wang, Xinlei Chen, Rishabh Jain, Mark Johnson, Dhruv Batra, Devi Parikh, Stefan Lee, and Peter Anderson. `nocaps`: novel object captioning at scale. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 8948–8957, 2019.
- Farhad Akhbardeh, Arkady Arkhangorodsky, Magdalena Biesialska, Ondřej Bojar, Rajen Chatterjee, Vishrav Chaudhary, Marta R. Costa-jussa, Cristina España-Bonet, Angela Fan, Christian Federmann, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Leonie Harter, Kenneth Heafield, Christopher Homan, Matthias Huck, Kwabena Amponsah-Kaakyire, Jungo Kasai, Daniel Khashabi, Kevin Knight, Tom Kocmi, Philipp Koehn, Nicholas Lourie, Christof Monz, Makoto Morishita, Masaaki Nagata, Ajay Nagesh, Toshiaki Nakazawa, Matteo Negri, Santanu Pal, Allahsera Auguste Tapo, Marco Turchi, Valentin Vydrin, and Marcos Zampieri. Findings of the 2021 conference on machine translation (WMT21). In Loic Barrault, Ondrej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussa, Christian Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grundkiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Tom Kocmi, Andre Martins, Makoto Morishita, and Christof Monz (eds.), *Proceedings of the Sixth Conference on Machine Translation*, pp. 1–88, Online, November 2021. Association for Computational Linguistics. URL <https://aclanthology.org/2021.wmt-1.1>.
- Chantal Amrhein and Rico Sennrich. Identifying weaknesses in machine translation metrics through minimum Bayes risk decoding: A case study for COMET. In Yulan He, Heng Ji, Sujian Li, Yang Liu, and Chua-Hui Chang (eds.), *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pp. 1125–1141, Online only, November 2022. Association for Computational Linguistics. URL <https://aclanthology.org/2022.aacl-main.83>.
- Loïc Barrault, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Philipp Koehn, Shervin Malmasi, Christof Monz, Mathias Müller, Santanu Pal, Matt Post, and Marcos Zampieri. Findings of the 2019 conference on machine translation (WMT19). In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, André Martins, Christof Monz, Matteo Negri, Aurélie Névélol, Mariana Neves, Matt Post, Marco Turchi, and Karin Verspoor (eds.), *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pp. 1–61, Florence, Italy, August 2019. Association for Computational Linguistics. doi: 10.18653/v1/W19-5301. URL <https://aclanthology.org/W19-5301>.
- Loïc Barrault, Magdalena Biesialska, Ondřej Bojar, Marta R. Costa-jussà, Christian Federmann, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Matthias Huck, Eric Joanis, Tom Kocmi, Philipp Koehn, Chi-kiu Lo, Nikola Ljubešić, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Santanu Pal, Matt Post, and Marcos Zampieri. Findings of the 2020 conference on machine translation (WMT20). In Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian Federmann, Mark Fishel, Alexander Fraser, Yvette Graham, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes,

- 594 Philipp Koehn, André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki
595 Nakazawa, and Matteo Negri (eds.), *Proceedings of the Fifth Conference on Machine Trans-*
596 *lation*, pp. 1–55, Online, November 2020. Association for Computational Linguistics. URL
597 <https://aclanthology.org/2020.wmt-1.1>.
- 598
599 Amanda Bertsch, Alex Xie, Graham Neubig, and Matthew Gormley. It’s MBR all the way down:
600 Modern generation techniques through the lens of minimum Bayes risk. In Yanai Elazar, Allyson
601 Ettinger, Nora Kassner, Sebastian Ruder, and Noah A. Smith (eds.), *Proceedings of the Big Pic-*
602 *ture Workshop*, pp. 108–122, Singapore, December 2023. Association for Computational Linguis-
603 tics. doi: 10.18653/v1/2023.bigpicture-1.9. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.bigpicture-1.9)
604 [bigpicture-1.9](https://aclanthology.org/2023.bigpicture-1.9).
- 605 Ondřej Bojar, Yvette Graham, Amir Kamran, and Miloš Stanojević. Results of the WMT16 met-
606 rics shared task. In Ondřej Bojar, Christian Buck, Rajen Chatterjee, Christian Federmann, Liane
607 Guillou, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Aurélie Névéol, Mariana Neves,
608 Pavel Pecina, Martin Popel, Philipp Koehn, Christof Monz, Matteo Negri, Matt Post, Lucia Spe-
609 cia, Karin Verspoor, Jörg Tiedemann, and Marco Turchi (eds.), *Proceedings of the First Confer-*
610 *ence on Machine Translation: Volume 2, Shared Task Papers*, pp. 199–231, Berlin, Germany,
611 August 2016. Association for Computational Linguistics. doi: 10.18653/v1/W16-2302. URL
612 <https://aclanthology.org/W16-2302>.
- 613 Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Shujian
614 Huang, Matthias Huck, Philipp Koehn, Qun Liu, Varvara Logacheva, Christof Monz, Matteo
615 Negri, Matt Post, Raphael Rubino, Lucia Specia, and Marco Turchi. Findings of the 2017 confer-
616 ence on machine translation (WMT17). In Ondřej Bojar, Christian Buck, Rajen Chatterjee, Chris-
617 tian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp
618 Koehn, and Julia Kreutzer (eds.), *Proceedings of the Second Conference on Machine Translation*,
619 pp. 169–214, Copenhagen, Denmark, September 2017. Association for Computational Linguis-
620 tics. doi: 10.18653/v1/W17-4717. URL <https://aclanthology.org/W17-4717>.
- 621 Ondřej Bojar, Christian Federmann, Mark Fishel, Yvette Graham, Barry Haddow, Matthias Huck,
622 Philipp Koehn, and Christof Monz. Findings of the 2018 conference on machine translation
623 (WMT18). In Ondřej Bojar, Rajen Chatterjee, Christian Federmann, Mark Fishel, Yvette Graham,
624 Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Philipp Koehn, Christof Monz, Matteo
625 Negri, Aurélie Névéol, Mariana Neves, Matt Post, Lucia Specia, Marco Turchi, and Karin Ver-
626 spoor (eds.), *Proceedings of the Third Conference on Machine Translation: Shared Task Papers*,
627 pp. 272–303, Belgium, Brussels, October 2018. Association for Computational Linguistics. doi:
628 10.18653/v1/W18-6401. URL <https://aclanthology.org/W18-6401>.
- 629 Hamed Bonab and Fazli Can. Less is more: A comprehensive framework for the number of compo-
630 nents of ensemble classifiers. *IEEE Transactions on Neural Networks and Learning Systems*, 30
631 (9):2735–2745, 2019. doi: 10.1109/TNNLS.2018.2886341.
- 632
633 Hamed R. Bonab and Fazli Can. A theoretical framework on the ideal number of classifiers for
634 online ensembles in data streams. In *Proceedings of the 25th ACM International on Conference*
635 *on Information and Knowledge Management, CIKM ’16*, pp. 2053–2056, New York, NY, USA,
636 2016. Association for Computing Machinery. ISBN 9781450340731. doi: 10.1145/2983323.
637 2983907. URL <https://doi.org/10.1145/2983323.2983907>.
- 638 Sebastian Borgeaud and Guy Emerson. Leveraging sentence similarity in natural language gen-
639 eration: Improving beam search using range voting. In Alexandra Birch, Andrew Finch, Hi-
640 roaki Hayashi, Kenneth Heafield, Marcin Junczys-Dowmunt, Ioannis Konstas, Xian Li, Graham
641 Neubig, and Yusuke Oda (eds.), *Proceedings of the Fourth Workshop on Neural Generation and*
642 *Translation*, pp. 97–109, Online, July 2020. Association for Computational Linguistics. doi:
643 10.18653/v1/2020.ngt-1.11. URL <https://aclanthology.org/2020.ngt-1.11>.
- 644 Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhari-
645 wal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agar-
646 wal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh,
647 Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz
Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec

- 648 Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In
649 H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (eds.), *Advances in Neu-*
650 *ral Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc.,
651 2020. URL [https://proceedings.neurips.cc/paper_files/paper/2020/](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf)
652 [file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfc4967418bfb8ac142f64a-Paper.pdf).
- 653 Julius Cheng and Andreas Vlachos. Faster minimum Bayes risk decoding with confidence-based
654 pruning. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Confer-*
655 *ence on Empirical Methods in Natural Language Processing*, pp. 12473–12480, Singapore, De-
656 cember 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.
657 767. URL <https://aclanthology.org/2023.emnlp-main.767>.
- 658 Julius Cheng and Andreas Vlachos. Measuring uncertainty in neural machine translation with
659 similarity-sensitive entropy. In Yvette Graham and Matthew Purver (eds.), *Proceedings of the*
660 *18th Conference of the European Chapter of the Association for Computational Linguistics (Vol-*
661 *ume 1: Long Papers)*, pp. 2115–2128, St. Julian’s, Malta, March 2024. Association for Compu-
662 tational Linguistics. URL <https://aclanthology.org/2024.eacl-long.129>.
- 663 David M Corey, William P Dunlap, and Michael J Burke. Averaging correlations: Expected values
664 and bias in combined pearson rs and fisher’s z transformations. *The Journal of general psychology*,
665 125(3):245–261, 1998.
- 666 Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe. mbrs: A library for
667 minimum bayes risk decoding, 2024a. URL <https://arxiv.org/abs/2408.04167>.
- 668 Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, Taro Watanabe, Hideki Tanaka, and Masao
669 Utiyama. Centroid-based efficient minimum Bayes risk decoding. In Lun-Wei Ku, Andre Martins,
670 and Vivek Srikumar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*,
671 pp. 11009–11018, Bangkok, Thailand and virtual meeting, August 2024b. Association for Com-
672 putational Linguistics. URL [https://aclanthology.org/2024.findings-acl.](https://aclanthology.org/2024.findings-acl.654)
673 654.
- 674 John DeNero, David Chiang, and Kevin Knight. Fast consensus decoding over translation forests. In
675 Keh-Yih Su, Jian Su, Janyce Wiebe, and Haizhou Li (eds.), *Proceedings of the Joint Conference*
676 *of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural*
677 *Language Processing of the AFNLP*, pp. 567–575, Suntec, Singapore, August 2009. Association
678 for Computational Linguistics. URL <https://aclanthology.org/P09-1064>.
- 679 Nan Duan, Mu Li, Ming Zhou, and Lei Cui. Improving phrase extraction via MBR phrase
680 scoring and pruning. In *Proceedings of Machine Translation Summit XIII: Papers*, Xi-
681 amen, China, September 19-23 2011. URL [https://aclanthology.org/2011.](https://aclanthology.org/2011.mtsummit-papers.20)
682 [mtsummit-papers.20](https://aclanthology.org/2011.mtsummit-papers.20).
- 683 Bryan Eikema and Wilker Aziz. Is MAP decoding all you need? the inadequacy of the mode
684 in neural machine translation. In Donia Scott, Nuria Bel, and Chengqing Zong (eds.), *Pro-*
685 *ceedings of the 28th International Conference on Computational Linguistics*, pp. 4506–4520,
686 Barcelona, Spain (Online), December 2020. International Committee on Computational Linguis-
687 tics. doi: 10.18653/v1/2020.coling-main.398. URL [https://aclanthology.org/2020.](https://aclanthology.org/2020.coling-main.398)
688 [coling-main.398](https://aclanthology.org/2020.coling-main.398).
- 689 Bryan Eikema and Wilker Aziz. Sampling-based approximations to minimum Bayes risk decod-
690 ing for neural machine translation. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.),
691 *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,
692 pp. 10978–10993, Abu Dhabi, United Arab Emirates, December 2022. Association for Computa-
693 tional Linguistics. doi: 10.18653/v1/2022.emnlp-main.754. URL [https://aclanthology.](https://aclanthology.org/2022.emnlp-main.754)
694 [org/2022.emnlp-main.754](https://aclanthology.org/2022.emnlp-main.754).
- 695 Angela Fan, Mike Lewis, and Yann Dauphin. Hierarchical neural story generation. In Iryna
696 Gurevych and Yusuke Miyao (eds.), *Proceedings of the 56th Annual Meeting of the Associa-*
697 *tion for Computational Linguistics (Volume 1: Long Papers)*, pp. 889–898, Melbourne, Aus-
698 tralia, July 2018. Association for Computational Linguistics. doi: 10.18653/v1/P18-1082. URL
699 <https://aclanthology.org/P18-1082>.
- 700
701

- 702 Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham
703 Neubig, and Andre Martins. Quality-aware decoding for neural machine translation. In Ma-
704 rine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (eds.), *Proceed-*
705 *ings of the 2022 Conference of the North American Chapter of the Association for Computa-*
706 *tional Linguistics: Human Language Technologies*, pp. 1396–1412, Seattle, United States, July
707 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.100. URL
708 <https://aclanthology.org/2022.naacl-main.100>.
- 709 Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. High Quality Rather than High
710 Model Probability: Minimum Bayes Risk Decoding with Neural Metrics. *Transactions of the*
711 *Association for Computational Linguistics*, 10:811–825, 07 2022a. ISSN 2307-387X. doi: 10.
712 1162/tacl.a.00491. URL https://doi.org/10.1162/tacl_a_00491.
- 713 Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis,
714 Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. Results of WMT22 metrics
715 shared task: Stop using BLEU – neural metrics are better and more robust. In Philipp Koehn,
716 Loïc Barrault, Ondřej Bojar, Fethi Bougares, Rajen Chatterjee, Marta R. Costa-jussà, Christian
717 Federmann, Mark Fishel, Alexander Fraser, Markus Freitag, Yvette Graham, Roman Grund-
718 kiewicz, Paco Guzman, Barry Haddow, Matthias Huck, Antonio Jimeno Yepes, Tom Kocmi,
719 André Martins, Makoto Morishita, Christof Monz, Masaaki Nagata, Toshiaki Nakazawa, Matteo
720 Negri, Aurélie Névél, Mariana Neves, Martin Popel, Marco Turchi, and Marcos Zampieri (eds.),
721 *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pp. 46–68, Abu Dhabi,
722 United Arab Emirates (Hybrid), December 2022b. Association for Computational Linguistics.
723 URL <https://aclanthology.org/2022.wmt-1.2>.
- 724 Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. Epsilon sampling rocks: Investi-
725 gating sampling strategies for minimum Bayes risk decoding for machine translation. In
726 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Compu-*
727 *tational Linguistics: EMNLP 2023*, pp. 9198–9209, Singapore, December 2023a. Association
728 for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.617. URL <https://aclanthology.org/2023.findings-emnlp.617>.
- 729 Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson,
730 Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho,
731 Alon Lavie, and George Foster. Results of WMT23 metrics shared task: Metrics might be guilty
732 but references are not innocent. In Philipp Koehn, Barry Haddow, Tom Kocmi, and Christof Monz
733 (eds.), *Proceedings of the Eighth Conference on Machine Translation*, pp. 578–628, Singapore,
734 December 2023b. Association for Computational Linguistics. doi: 10.18653/v1/2023.wmt-1.51.
735 URL <https://aclanthology.org/2023.wmt-1.51>.
- 736 Stuart Geman, Elie Bienenstock, and René Doursat. Neural networks and the bias/variance dilemma.
737 *Neural computation*, 4(1):1–58, 1992.
- 738 Kevin Gimpel, Dhruv Batra, Chris Dyer, and Gregory Shakhnarovich. A systematic exploration of
739 diversity in machine translation. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen
740 Livescu, and Steven Bethard (eds.), *Proceedings of the 2013 Conference on Empirical Methods*
741 *in Natural Language Processing*, pp. 1100–1111, Seattle, Washington, USA, October 2013. As-
742 sociation for Computational Linguistics. URL <https://aclanthology.org/D13-1111>.
- 743 Bogdan Gliwa, Iwona Mochol, Maciej Biesek, and Aleksander Wawer. SAMSum corpus: A human-
744 annotated dialogue dataset for abstractive summarization. In Lu Wang, Jackie Chi Kit Cheung,
745 Giuseppe Carenini, and Fei Liu (eds.), *Proceedings of the 2nd Workshop on New Frontiers in Sum-*
746 *marization*, pp. 70–79, Hong Kong, China, November 2019. Association for Computational Lin-
747 guistics. doi: 10.18653/v1/D19-5409. URL <https://aclanthology.org/D19-5409>.
- 748 Taisiya Glushkova, Chrysoula Zerva, and André F. T. Martins. BLEU meets COMET: Combin-
749 ing lexical and neural metrics towards robust machine translation evaluation. In Mary Nurmi-
750 nen, Judith Brenner, Maarit Koponen, Sirkku Latomaa, Mikhail Mikhailov, Frederike Schierl,
751 Tharindu Ranasinghe, Eva Vanmassenhove, Sergi Alvarez Vidal, Nora Aranberri, Mara Nunzi-
752 atini, Carla Parra Escartín, Mikel Forcada, Maja Popovic, Carolina Scarton, and Helena Mo-
753 niz (eds.), *Proceedings of the 24th Annual Conference of the European Association for Machine*
754 *Learning*, pp. 100–110, London, UK, September 2022. Association for Computational Linguis-
755 tics. doi: 10.18653/v1/E22-100. URL <https://aclanthology.org/E22-100>.

- 756 *Translation*, pp. 47–58, Tampere, Finland, June 2023. European Association for Machine Trans-
757 lation. URL <https://aclanthology.org/2023.eamt-1.6>.
- 758
- 759 Vaibhava Goel and William J Byrne. Minimum bayes-risk automatic speech recognition. *Com-
760 puter Speech & Language*, 14(2):115–135, 2000. ISSN 0885-2308. doi: [https://doi.org/10.1006/
761 csla.2000.0138](https://doi.org/10.1006/csla.2000.0138). URL [https://www.sciencedirect.com/science/article/pii/
762 S0885230800901384](https://www.sciencedirect.com/science/article/pii/S0885230800901384).
- 763 Yvette Graham, Timothy Baldwin, Alistair Moffat, and Justin Zobel. Continuous measurement
764 scales in human evaluation of machine translation. In Antonio Pareja-Lora, Maria Liakata, and
765 Stefanie Dipper (eds.), *Proceedings of the 7th Linguistic Annotation Workshop and Interoper-
766 ability with Discourse*, pp. 33–41, Sofia, Bulgaria, August 2013. Association for Computational
767 Linguistics. URL <https://aclanthology.org/W13-2305>.
- 768 David Heineman, Yao Dou, and Wei Xu. Improving minimum bayes risk decoding with multi-
769 prompt, 2024. URL <https://arxiv.org/abs/2407.15343>.
- 770 John Hewitt, Christopher Manning, and Percy Liang. Truncation sampling as language model
771 desmoothing. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), *Findings of the Asso-
772 ciation for Computational Linguistics: EMNLP 2022*, pp. 3414–3427, Abu Dhabi, United Arab
773 Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.
774 findings-emnlp.249. URL [https://aclanthology.org/2022.findings-emnlp.
775 249](https://aclanthology.org/2022.findings-emnlp.249).
- 776 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text
777 degeneration. In *International Conference on Learning Representations*, 2020. URL [https://
778 //openreview.net/forum?id=rygGQyrFvH](https://openreview.net/forum?id=rygGQyrFvH).
- 779
- 780 Yuu Jinnai, Ukyo Honda, Tetsuro Morimura, and Peinan Zhang. Generating diverse and high-
781 quality texts by minimum Bayes risk decoding. In Lun-Wei Ku, Andre Martins, and Vivek Sriku-
782 mar (eds.), *Findings of the Association for Computational Linguistics ACL 2024*, pp. 8494–8525,
783 Bangkok, Thailand and virtual meeting, August 2024a. Association for Computational Linguis-
784 tics. URL <https://aclanthology.org/2024.findings-acl.503>.
- 785 Yuu Jinnai, Tetsuro Morimura, Ukyo Honda, Kaito Ariu, and Kenshi Abe. Model-based mini-
786 mum Bayes risk decoding for text generation. In Ruslan Salakhutdinov, Zico Kolter, Katherine
787 Heller, Adrian Weller, Nuria Oliver, Jonathan Scarlett, and Felix Berkenkamp (eds.), *Proceed-
788 ings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings
789 of Machine Learning Research*, pp. 22326–22347. PMLR, 21–27 Jul 2024b. URL [https://
790 proceedings.mlr.press/v235/jinnai24a.html](https://proceedings.mlr.press/v235/jinnai24a.html).
- 791 Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descrip-
792 tions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
793 3128–3137, 2015.
- 794 Teun Kloek and Herman K Van Dijk. Bayesian estimates of equation system parameters: an ap-
795 plication of integration by monte carlo. *Econometrica: Journal of the Econometric Society*, pp.
796 1–19, 1978.
- 797
- 798 Anders Krogh and Jesper Vedelsby. Neural network ensembles, cross validation, and active learning.
799 In G. Tesauro, D. Touretzky, and T. Leen (eds.), *Advances in Neural Information Processing Sys-
800 tems*, volume 7. MIT Press, 1994. URL [https://proceedings.neurips.cc/paper_
801 files/paper/1994/file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/1994/file/b8c37e33defde51cf91e1e03e51657da-Paper.pdf).
- 802 Shankar Kumar and William Byrne. Minimum Bayes-risk word alignments of bilingual texts.
803 In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Process-
804 ing (EMNLP 2002)*, pp. 140–147. Association for Computational Linguistics, July 2002. doi:
805 10.3115/1118693.1118712. URL <https://aclanthology.org/W02-1019>.
- 806 Shankar Kumar and William Byrne. Minimum Bayes-risk decoding for statistical machine transla-
807 tion. In *Proceedings of the Human Language Technology Conference of the North American
808 Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pp. 169–176,
809 Boston, Massachusetts, USA, May 2 - May 7 2004. Association for Computational Linguistics.
URL <https://aclanthology.org/N04-1022>.

- 810 Jiwei Li and Dan Jurafsky. Mutual information and diverse decoding improve neural machine trans-
811 lation, 2016. URL <https://arxiv.org/abs/1601.00372>.
812
- 813 Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. A diversity-promoting objec-
814 tive function for neural conversation models. In Kevin Knight, Ani Nenkova, and Owen Rambow
815 (eds.), *Proceedings of the 2016 Conference of the North American Chapter of the Association for*
816 *Computational Linguistics: Human Language Technologies*, pp. 110–119, San Diego, California,
817 June 2016. Association for Computational Linguistics. doi: 10.18653/v1/N16-1014. URL
818 <https://aclanthology.org/N16-1014>.
- 819 Tsung-Yi Lin, Michael Maire, Serge J. Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
820 Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *European*
821 *Conference on Computer Vision*, 2014. URL [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:14113767)
822 [CorpusID:14113767](https://api.semanticscholar.org/CorpusID:14113767).
- 823 Huan Liu, Amit Mandvikar, and Jigar Mody. An empirical study of building compact ensembles.
824 In *Advances in Web-Age Information Management: 5th International Conference, WAIM 2004,*
825 *Dalian, China, July 15-17, 2004 5*, pp. 622–627. Springer, 2004.
826
- 827 Tom Mitchell. Introduction to machine learning. *Machine learning*, 7:2–5, 1997.
828
- 829 Mathias Müller and Rico Sennrich. Understanding the properties of minimum Bayes risk decoding
830 in neural machine translation. In Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (eds.),
831 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the*
832 *11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*,
833 pp. 259–272, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/
834 v1/2021.acl-long.22. URL <https://aclanthology.org/2021.acl-long.22>.
- 835 Shashi Narayan, Shay B. Cohen, and Mirella Lapata. Don’t give me the details, just the sum-
836 mary! topic-aware convolutional neural networks for extreme summarization. In Ellen Riloff,
837 David Chiang, Julia Hockenmaier, and Jun’ichi Tsujii (eds.), *Proceedings of the 2018 Confer-*
838 *ence on Empirical Methods in Natural Language Processing*, pp. 1797–1807, Brussels, Bel-
839 gium, October–November 2018. Association for Computational Linguistics. doi: 10.18653/v1/
840 D18-1206. URL <https://aclanthology.org/D18-1206>.
- 841 Subhajit Naskar, Daniel Deutsch, and Markus Freitag. Quality estimation using minimum Bayes
842 risk. In Philipp Koehn, Barry Haddow, Tom Kocmi, and Christof Monz (eds.), *Proceedings*
843 *of the Eighth Conference on Machine Translation*, pp. 806–811, Singapore, December 2023.
844 Association for Computational Linguistics. doi: 10.18653/v1/2023.wmt-1.67. URL <https://aclanthology.org/2023.wmt-1.67>.
845
- 846 Atsumoto Ohashi, Ukyo Honda, Tetsuro Morimura, and Yuu Jinnai. On the true distribution approx-
847 imation of minimum Bayes-risk decoding. In Kevin Duh, Helena Gomez, and Steven Bethard
848 (eds.), *Proceedings of the 2024 Conference of the North American Chapter of the Association for*
849 *Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pp. 459–
850 468, Mexico City, Mexico, June 2024. Association for Computational Linguistics. doi: 10.18653/
851 v1/2024.naacl-short.38. URL <https://aclanthology.org/2024.naacl-short.38>.
852
- 853 OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Floren-
854 cia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red
855 Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Moham-
856 mad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher
857 Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brock-
858 man, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann,
859 Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis,
860 Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey
861 Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux,
862 Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila
863 Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix,
Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gib-
son, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan

- 864 Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hal-
865 lacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan
866 Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu,
867 Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun
868 Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Ka-
869 mali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook
870 Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel
871 Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kopic, Gretchen
872 Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel
873 Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez,
874 Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv
875 Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney,
876 Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick,
877 Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel
878 Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Ra-
879 jeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe,
880 Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel
881 Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe
882 de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny,
883 Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl,
884 Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra
885 Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders,
886 Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Sel-
887 sam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor,
888 Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky,
889 Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang,
890 Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-
891 ston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vi-
892 jayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan
893 Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng,
894 Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Work-
895 man, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming
896 Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao
897 Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024. URL
898 <https://arxiv.org/abs/2303.08774>.
- 899 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic
900 evaluation of machine translation. In Pierre Isabelle, Eugene Charniak, and Dekang Lin (eds.),
901 *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pp.
902 311–318, Philadelphia, Pennsylvania, USA, July 2002. Association for Computational Linguis-
903 tics. doi: 10.3115/1073083.1073135. URL <https://aclanthology.org/P02-1040>.
- 904 Maja Popović. chrF: character n-gram F-score for automatic MT evaluation. In Ondřej Bojar,
905 Rajan Chatterjee, Christian Federmann, Barry Haddow, Chris Hokamp, Matthias Huck, Varvara
906 Logacheva, and Pavel Pecina (eds.), *Proceedings of the Tenth Workshop on Statistical Machine
907 Translation*, pp. 392–395, Lisbon, Portugal, September 2015. Association for Computational Lin-
908 guistics. doi: 10.18653/v1/W15-3049. URL <https://aclanthology.org/W15-3049>.
- 909 Maja Popović. chrF++: words helping character n-grams. In Ondřej Bojar, Christian Buck, Rajen
910 Chatterjee, Christian Federmann, Yvette Graham, Barry Haddow, Matthias Huck, Antonio Ji-
911 meno Yepes, Philipp Koehn, and Julia Kreutzer (eds.), *Proceedings of the Second Conference
912 on Machine Translation*, pp. 612–618, Copenhagen, Denmark, September 2017. Association for
913 Computational Linguistics. doi: 10.18653/v1/W17-4770. URL <https://aclanthology.org/W17-4770>.
- 914 Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. COMET: A neural framework for MT
915 evaluation. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu (eds.), *Proceedings of the
916 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 2685–
917 2702, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/
2020.emnlp-main.213. URL <https://aclanthology.org/2020.emnlp-main.213>.

- 918 CP Robert. Monte carlo statistical methods, 1999.
919
- 920 Thibault Sellam, Dipanjan Das, and Ankur Parikh. BLEURT: Learning robust metrics for text
921 generation. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault (eds.), *Proceedings*
922 *of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 7881–7892,
923 Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.acl-main.
924 704. URL <https://aclanthology.org/2020.acl-main.704>.
- 925 Felix Stahlberg, Adrià de Gispert, Eva Hasler, and Bill Byrne. Neural machine translation by
926 minimising the Bayes-risk with respect to syntactic translation lattices. In Mirella Lapata,
927 Phil Blunsom, and Alexander Koller (eds.), *Proceedings of the 15th Conference of the Eu-*
928 *ropean Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*,
929 pp. 362–368, Valencia, Spain, April 2017. Association for Computational Linguistics. URL
930 <https://aclanthology.org/E17-2058>.
- 931 Student. Probable error of a correlation coefficient. *Biometrika*, pp. 302–310, 1908.
932
- 933 Mirac Suzgun, Luke Melas-Kyriazi, and Dan Jurafsky. Follow the wisdom of the crowd: Effective
934 text generation via minimum Bayes risk decoding. In Anna Rogers, Jordan Boyd-Graber,
935 and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL*
936 *2023*, pp. 4265–4293, Toronto, Canada, July 2023. Association for Computational Linguistics.
937 doi: 10.18653/v1/2023.findings-acl.262. URL [https://aclanthology.org/2023.](https://aclanthology.org/2023.findings-acl.262)
938 [findings-acl.262](https://aclanthology.org/2023.findings-acl.262).
- 939 Firas Trabelsi, David Vilar, Mara Finkelstein, and Markus Freitag. Efficient minimum bayes risk
940 decoding using low-rank matrix completion algorithms, 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2406.02832)
941 [abs/2406.02832](https://arxiv.org/abs/2406.02832).
- 942 Jannis Vamvas and Rico Sennrich. Linear-time minimum Bayes risk decoding with reference ag-
943 gregation. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), *Proceedings of the 62nd*
944 *Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pp.
945 790–801, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL
946 <https://aclanthology.org/2024.acl-short.71>.
- 947 Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. GLUE:
948 A multi-task benchmark and analysis platform for natural language understanding. In Tal Linzen,
949 Grzegorz Chrupała, and Afra Alishahi (eds.), *Proceedings of the 2018 EMNLP Workshop Black-*
950 *boxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–355, Brussels, Belgium,
951 November 2018. Association for Computational Linguistics. doi: 10.18653/v1/W18-5446. URL
952 <https://aclanthology.org/W18-5446>.
953
- 954 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
955 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick
956 von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gug-
957 ger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art
958 natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in*
959 *Natural Language Processing: System Demonstrations*, pp. 38–45, Online, October 2020. As-
960 sociation for Computational Linguistics. URL [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2020.emnlp-demos.6)
961 [2020.emnlp-demos.6](https://www.aclweb.org/anthology/2020.emnlp-demos.6).
- 962 Danny Wood, Tingting Mu, Andrew M. Webb, Henry W. J. Reeve, Mikel Luján, and Gavin Brown.
963 A unified theory of diversity in ensemble learning. *J. Mach. Learn. Res.*, 24(1), mar 2024. ISSN
964 1532-4435.
- 965 Tianyi Zhang*, Varsha Kishore*, Felix Wu*, Kilian Q. Weinberger, and Yoav Artzi. Bertscore:
966 Evaluating text generation with bert. In *International Conference on Learning Representations*,
967 2020. URL <https://openreview.net/forum?id=SkeHuCVFDr>.
968
969
970
971

A PROOF FOR THEOREM 1

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

$$(\hat{u}_i - \bar{u}_i)^2 \quad (12)$$

$$=(\hat{u}_i)^2 - 2\hat{u}_i\bar{u}_i + (\bar{u}_i)^2 \quad (13)$$

$$=(\hat{u}_i)^2 - 2\hat{u}_i\bar{u}_i + 2(\bar{u}_i)^2 - (\bar{u}_i)^2 \quad (14)$$

$$=(\hat{u}_i)^2 - 2\hat{u}_i\bar{u}_i + 2\bar{u}_i\bar{u}_i - (\bar{u}_i)^2 \quad (15)$$

$$=(\hat{u}_i)^2 - \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} 2\hat{u}_i u_i^j + \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} 2\bar{u}_i u_i^j - (\bar{u}_i)^2 \quad (16)$$

$$=\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i)^2 - \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} 2\hat{u}_i u_i^j + \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} 2\bar{u}_i u_i^j - \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i)^2 \quad (17)$$

$$=\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} ((\hat{u}_i)^2 - 2\hat{u}_i u_i^j + 2\bar{u}_i u_i^j - (\bar{u}_i)^2) \quad (18)$$

$$=\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} ((\hat{u}_i)^2 - 2\hat{u}_i u_i^j + (u_i^j)^2 - (u_i^j)^2 + 2\bar{u}_i u_i^j - (\bar{u}_i)^2) \quad (19)$$

$$=\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} ((\hat{u}_i)^2 - 2\hat{u}_i u_i^j + (u_i^j)^2 - ((u_i^j)^2 - 2\bar{u}_i u_i^j + (\bar{u}_i)^2)) \quad (20)$$

$$=\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} ((\hat{u}_i - u_i^j)^2 - (\bar{u}_i - u_i^j)^2) \quad (21)$$

$$=\underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_\theta(h_i, y_j))^2}_{\text{Bias}} - \underbrace{\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i - f_\theta(h_i, y_j))^2}_{\text{Diversity}} \quad (22)$$

B PROOF FOR THEOREM 2

$$MSE(\hat{\mathbf{u}}, \bar{\mathbf{u}}) \quad (23)$$

$$=\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} (\hat{u}_i - \bar{u}_i)^2 \quad (24)$$

$$=\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} \left(\frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_\theta(h_i, y_j))^2 - \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i - f_\theta(h_i, y_j))^2 \right) \quad (25)$$

$$=\underbrace{\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\hat{u}_i - f_\theta(h_i, y_j))^2}_{\text{Bias for all hypotheses}} - \underbrace{\frac{1}{|\mathcal{H}|} \sum_{i=1}^{|\mathcal{H}|} \frac{1}{|\mathcal{Y}|} \sum_{j=1}^{|\mathcal{Y}|} (\bar{u}_i - f_\theta(h_i, y_j))^2}_{\text{Diversity for all hypotheses}} \quad (26)$$

C INTERPRETATION AS ENSEMBLE LEARNING

When $|\mathcal{Y}|$ is large enough to satisfy the law of large numbers, we can induce the following expectation in MBR decoding by using a model’s prediction, $P(y|x)$:

$$\arg \max_{h \in \mathcal{H}} \sum_{y \in \Omega} f_{\theta}(h, y) P(y|x) \quad (27)$$

$$= \arg \max_{h \in \mathcal{H}} \mathbb{E}_{P(y|x)} [f_{\theta}(h, y)] \quad (28)$$

$$\approx \arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P(y|x) \quad (29)$$

Since this expectation is based on $P(y|x)$, we can understand the importance of increasing the number of pseudo-references to induce a reliable $P(y|x)$.

Theorem 3. *When $f_{\theta}(h, y)$ is normalized as a probability $P_{\theta}(h|y)$, Equation 27 is equivalent to Bayes Optimal Classifier (BOC) in Mitchell (1997).*

Proof. Self-evident by the following reformulation:

$$\arg \max_{h \in \mathcal{H}} \sum_{y \in \Omega} f_{\theta}(h, y) P(y|x) = \arg \max_{h \in \mathcal{H}} \sum_{y \in \Omega} P_{\theta}(h|y) P(y|x) \quad (30)$$

□

Theorem 4. *When $f_{\theta}(h, y)$ is normalized as a probability $P_{\theta}(h|y)$, Equation 29 is equivalent to the Gibbs algorithm in Mitchell (1997) that approximates BOC by sampling.*

Proof. Self-evident by the following reformulation:

$$\arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P(y|x) \quad (31)$$

$$= \arg \max_{h \in \mathcal{H}} \sum_{y \in \mathcal{Y}} P_{\theta}(h|y), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P(y|x) \quad (32)$$

□

Hence, we can understand that MBR decoding represented as Eqs. 27 and 29 approximates the ensemble learning method, BOC. In this interpretation, since $P(y|x)$ is a prior of BOC, we can also understand that MBR approximately uses the model-predicted probability as its prior.

When pruning unnecessary y in the BOC formulation of Equation 30, because the sum of $P_{\theta}(h|y)$ for all h is always 1, we can determine the importance of y based solely on $P(y|x)$. Since we can arbitrarily choose $P(y|x)$ during sampling, we understand that pruning methods select the importance of each y as a prior in BOC. Note that utility functions are not always normalized; therefore, there is a gap between this interpretation and the actual MBR decoding. Addressing this gap remains an open problem.

In practice, directly drawing samples from $P(y|x)$ is intractable. Therefore, we must use approximate search methods, which are commonly influenced by left-to-right decoding and threshold values. These factors can lead to unreachable states and biases, as seen in greedy or beam decoding and other sampling approaches. Letting $P'(y|x)$ denote the model’s prediction with the approximate search, we can similarly induce the following expectation:

$$\arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P'(y|x) \quad (33)$$

$$\approx \arg \max_{h \in \mathcal{H}} \mathbb{E}_{P'(y|x)} [f_{\theta}(h, y)] \quad (34)$$

Unfortunately, due to $P'(y|x)$, Equation 34 deviates from Equation 28. To precisely predict Eq. 27 using samples from $P'(y|x)$, we can consider the following theorem:

Theorem 5. When $|\mathcal{Y}|$ is enough large to satisfy the law of large numbers, by using importance sampling, we can induce Eq. 28 from $P'(y|x)$.

Proof.

$$\arg \max_{h \in \mathcal{H}} \mathbb{E}_{P(y|x)} [f_{\theta}(h, y)] \quad (35)$$

$$= \arg \max_{h \in \mathcal{H}} \sum_y P(y|x) f_{\theta}(h, y) \quad (36)$$

$$= \arg \max_{h \in \mathcal{H}} \sum_y P(y|x) f_{\theta}(h, y) \frac{P'(y|x)}{P'(y|x)} \quad (37)$$

$$= \arg \max_{h \in \mathcal{H}} \sum_y P'(y|x) f_{\theta}(h, y) \frac{P(y|x)}{P'(y|x)} \quad (38)$$

$$\approx \arg \max_{h \in \mathcal{H}} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y) \frac{P(y|x)}{P'(y|x)}, \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P'(y|x) \quad (39)$$

□

Apart from the fact that even precisely calculating $P'(y|x)$ is also difficult, we can induce the following theorem:

Theorem 6. When $|\mathcal{Y}|$ is enough large to satisfy the law of large numbers and $P'(y|x)$ equals a discrete uniform distribution $\mathcal{U}(0, |\mathcal{Y}|)$, Equation 39 is equivalent to Model-based MBR (MBMBR) of Jinnai et al. (2024b).

Proof.

$$\arg \max_{h \in \mathcal{H}} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y) \frac{P(y|x)}{P'(y|x)}, \quad y_1, \dots, y_{|\mathcal{Y}|} \sim P'(y|x) \quad (40)$$

$$= \arg \max_{h \in \mathcal{H}} \frac{1}{|\mathcal{Y}|} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y) P(y|x), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim \mathcal{U}(0, |\mathcal{Y}|) \quad (41)$$

$$= \arg \max_{h \in \mathcal{H}} \sum_{y \in \mathcal{Y}} f_{\theta}(h, y) P(y|x), \quad y_1, \dots, y_{|\mathcal{Y}|} \sim \mathcal{U}(0, |\mathcal{Y}|) \quad (42)$$

□

From Theorem 6, we can understand that MBMBR is an effective approach when sampling methods are unreliable. Based on the interpretation from the viewpoint of BOC, Equation 42 estimates the importance for each y through prior $P(y|x)$, which can be used for pruning y .

Even though our interpretation can explain the pruning of pseudo-references based on priors in BOC, pruning hypotheses are out-of-scope of this interpretation.

D CORRELATION OF BIAS AND DIVERSITY TO PERFORMANCE

We further investigate whether our analysis in §5.2 is consistent when metrics used in MBR decoding and performance evaluation are different.

Settings Based on the inherited settings from §5.2, we changed the performance evaluation metrics, COMET and BERTScore to BLEURT (Sellam et al., 2020) and chrF++ (Popović, 2015; 2017). We used BLEURT on single sentence generation tasks, WMT19 En-De and En-Ru, XSum, MSCOCO, and NoCaps. Since SAMSUM is a multiple-sentence generation task and BLEURT cannot handle it, we used chrF++ instead.

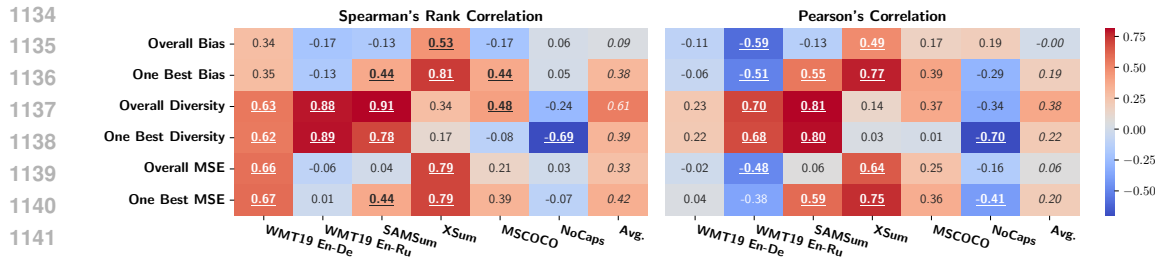


Figure 3: Correlation between measures in our decomposition and performance for each dataset when using different metrics in decoding and performance evaluation. The notations are the same as Figure 1.

Results Figure 3 shows the correlation. Similar to the results in §5.2, the measures, i.e., ONE BEST BIAS, OVERALL DIVERSITY, and ONE BEST MSE in Spearman’s rank correlation correlate well with the rank in performance, even though these correlation values are degraded by different evaluation metrics from decoding time. The lower correlation values in Pearson’s correlation than Spearman’s rank correlation also show similar tendencies in §5.2 and indicate the difficulty of precisely estimating the performance values from these measures. From these results, we can confirm that correlation tendencies are consistent when changing the performance evaluation metrics.

E BIAS AND DIVERSITY TRADE-OFF

Similar to Appendix D, we further investigate whether our analysis in §5.3 is consistent when metrics used in MBR decoding and performance evaluation are different.

Settings We inherited the setting of Appendix D. Thus, COMET and BERTScore used in MBR decoding are replaced with BLEURT and chrF++ in performance evaluation.

Results Figure 4 shows the results. We can see the changed performances in the subfigures of the rightmost column. The entire tendencies of beam decoding are almost the same as Figure 2, excluding the case of the performance drop in SAMSum, whose evaluation metric is changed from BERTScore to chrF++. However, this behavior is reasonable considering the highest One Best Bias and lowest Overall Diversity of beam decoding in SAMSum. This result shows the possibility of adopting bias and diversity in a metric to the estimation of performance in other evaluation metrics. On the other hand, these relationships are not always consistent as represented in the uncorrelated values on NoCaps that permit diversified generation by its 10 gold references.

F BIAS AND DIVERSITY OF MAMBR

Figures 5, 6, and 7 show the bias and diversity corresponding to the results in Tables 1, 2, and 3, respectively. The results show that MAMBR actually increases the diversities in WMT19 En-De and En-Ru and SAMSum but not in the other dataset. Thus, this improvement depends on the datasets. On the other hand, we can see the improvement of bias in some cases. This is reasonable because using multiple metric models itself is an ensembling approach and can contribute to performance improvement.

G EXPERIMENTAL RESULTS ON THE FIRST 1000 EXAMPLES

To consider more detailed configurations and reveal the possibility of more efficient investigation, we conducted additional evaluation using only the first 1000 examples for each dataset based on the setting of Jinnai et al. (2024b).

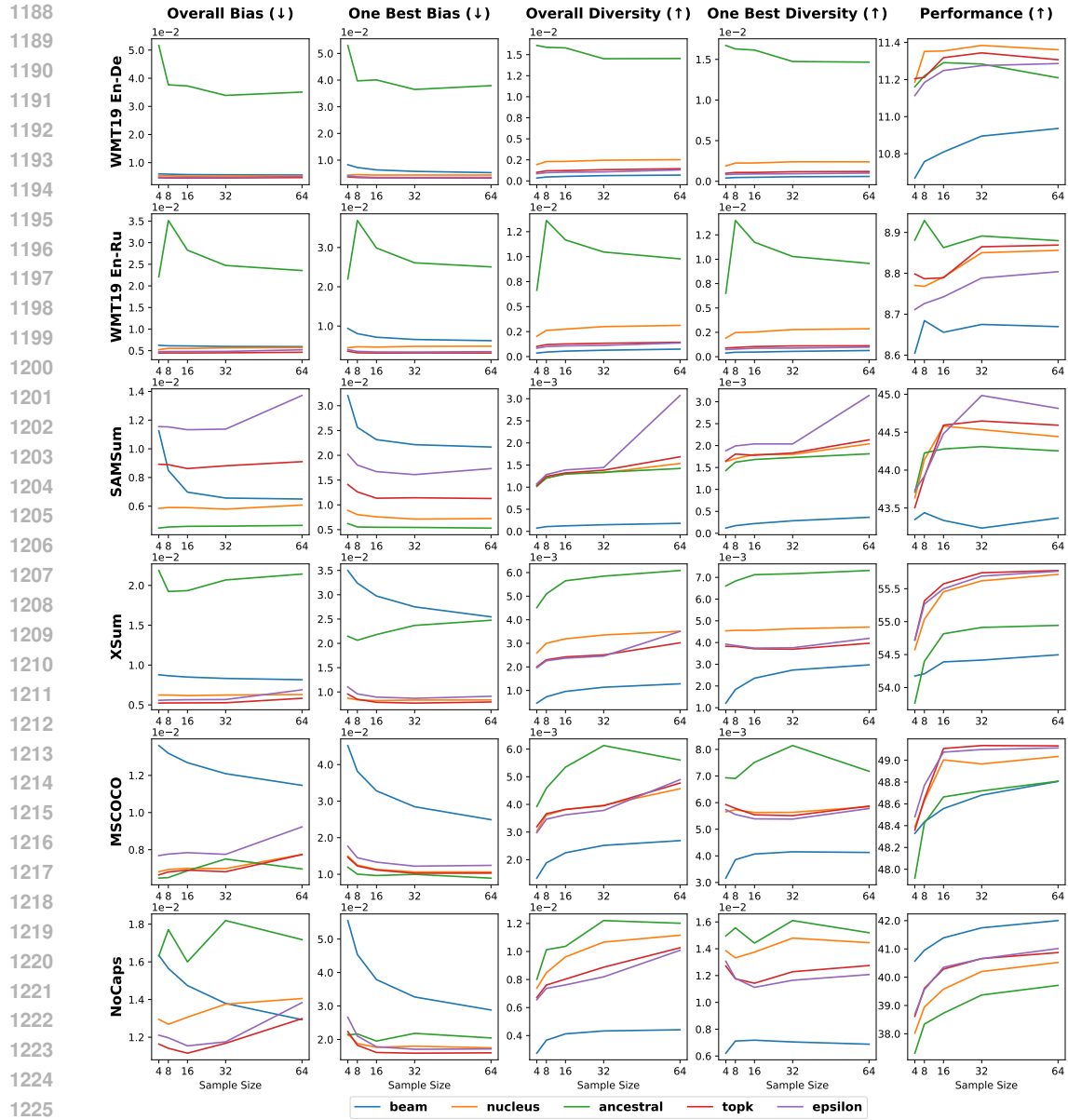


Figure 4: The relationship between bias, diversity, and performance in MBR decoding when using different metrics in decoding and performance evaluation. The notations are the same as Figure 2.

G.1 HYPOTHESES GENERATED BY DIFFERENT SAMPLING STRATEGIES

Figures 8 to 12 present the bias and diversity decomposition plots for different hypothesis generation strategies. The results indicate that differences in the generated hypotheses influence performance in some cases, whereas the overall tendencies of the sampling strategy used for generating pseudo-references remain similar despite these variations.

G.2 MAMBR

Tables 4 to 6 show the MAMBR results for the first 1000 lines. From these results, we observe a similar trend to those obtained when the dataset is fully used, as described in §5.4. Similarly, Figures 13 to 15 demonstrate that the results are nearly identical to those obtained when the dataset is fully utilized.

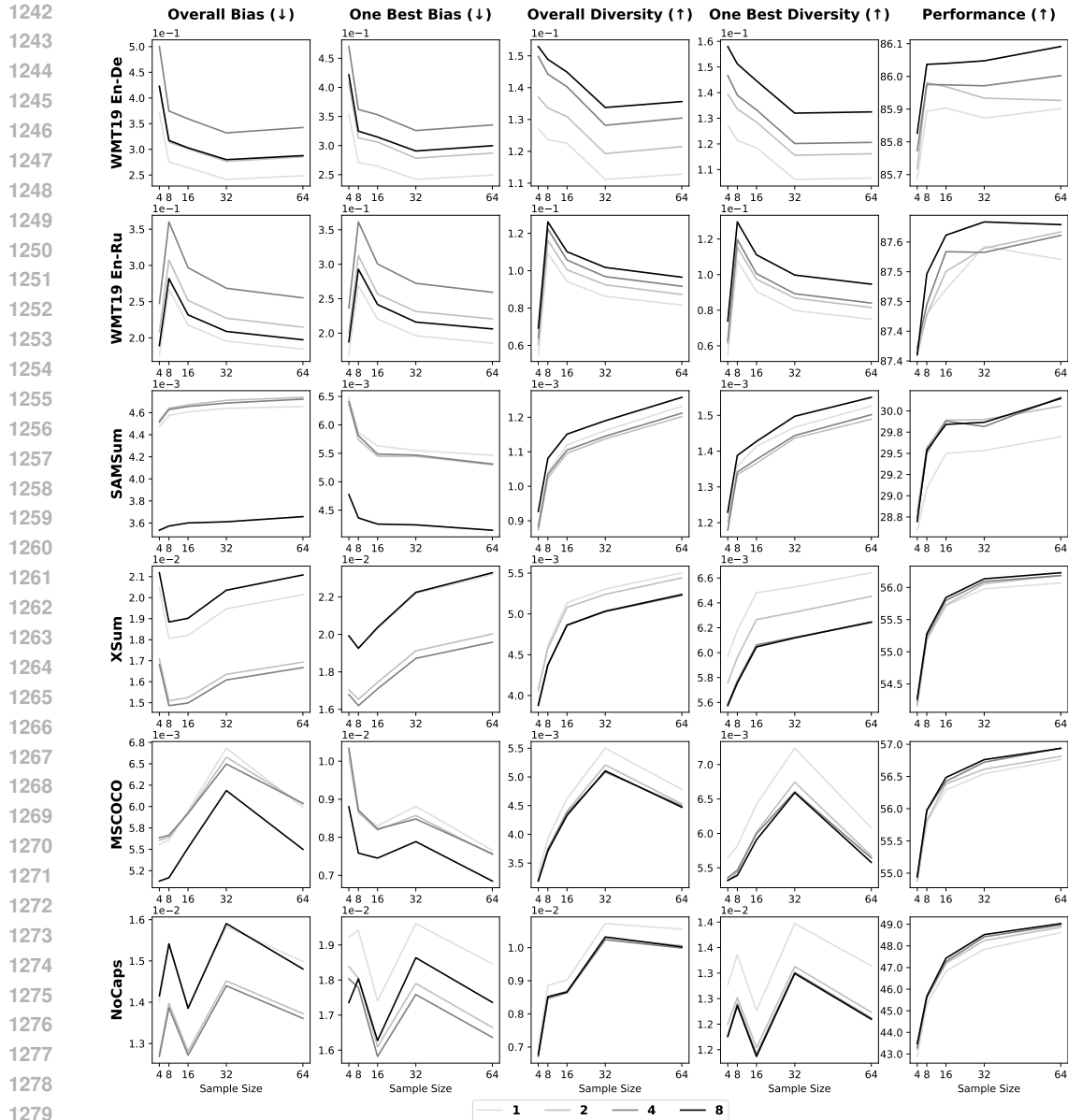


Figure 5: The relationship between bias, diversity, and performance in MAMBR decoding with pseudo-references generated by ancestral sampling. The notations are the same as Figure 6.

H LIMITATION

Although our bias-diversity decomposition for MBR decoding can explain the behavior of pseudo-references and utility functions, a theoretical explanation for the effectiveness of the used hypotheses is a model-side behavior and thus beyond the scope of our analysis. Therefore, corresponding to this limitation, we conduct a limited empirical analysis presented in Appendix G.1, similar to previous works (Eikema & Aziz, 2020; Fernandes et al., 2022; Freitag et al., 2023a).

1296
 1297
 1298
 1299
 1300
 1301
 1302
 1303
 1304
 1305
 1306
 1307
 1308
 1309
 1310
 1311
 1312
 1313
 1314
 1315
 1316
 1317
 1318
 1319
 1320
 1321
 1322
 1323
 1324
 1325
 1326
 1327
 1328
 1329
 1330
 1331
 1332
 1333
 1334
 1335
 1336
 1337
 1338
 1339
 1340
 1341
 1342
 1343
 1344
 1345
 1346
 1347
 1348
 1349

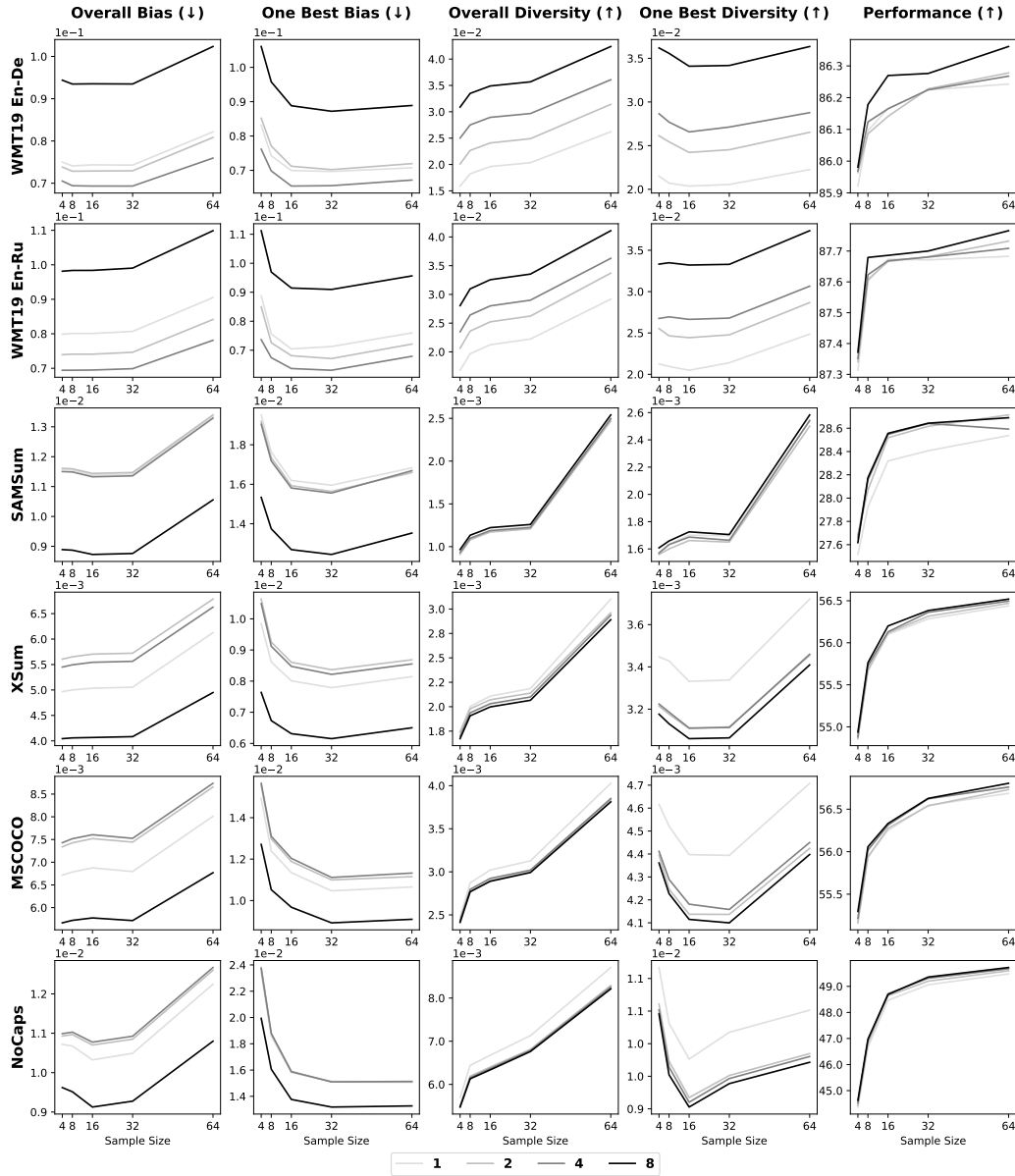


Figure 6: The relationship between bias, diversity, and performance in MAMBR decoding with pseudo-references generated by epsilon sampling. The lines indicate the score for each number of used metric models. Other notations are the same as Figure 6.

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

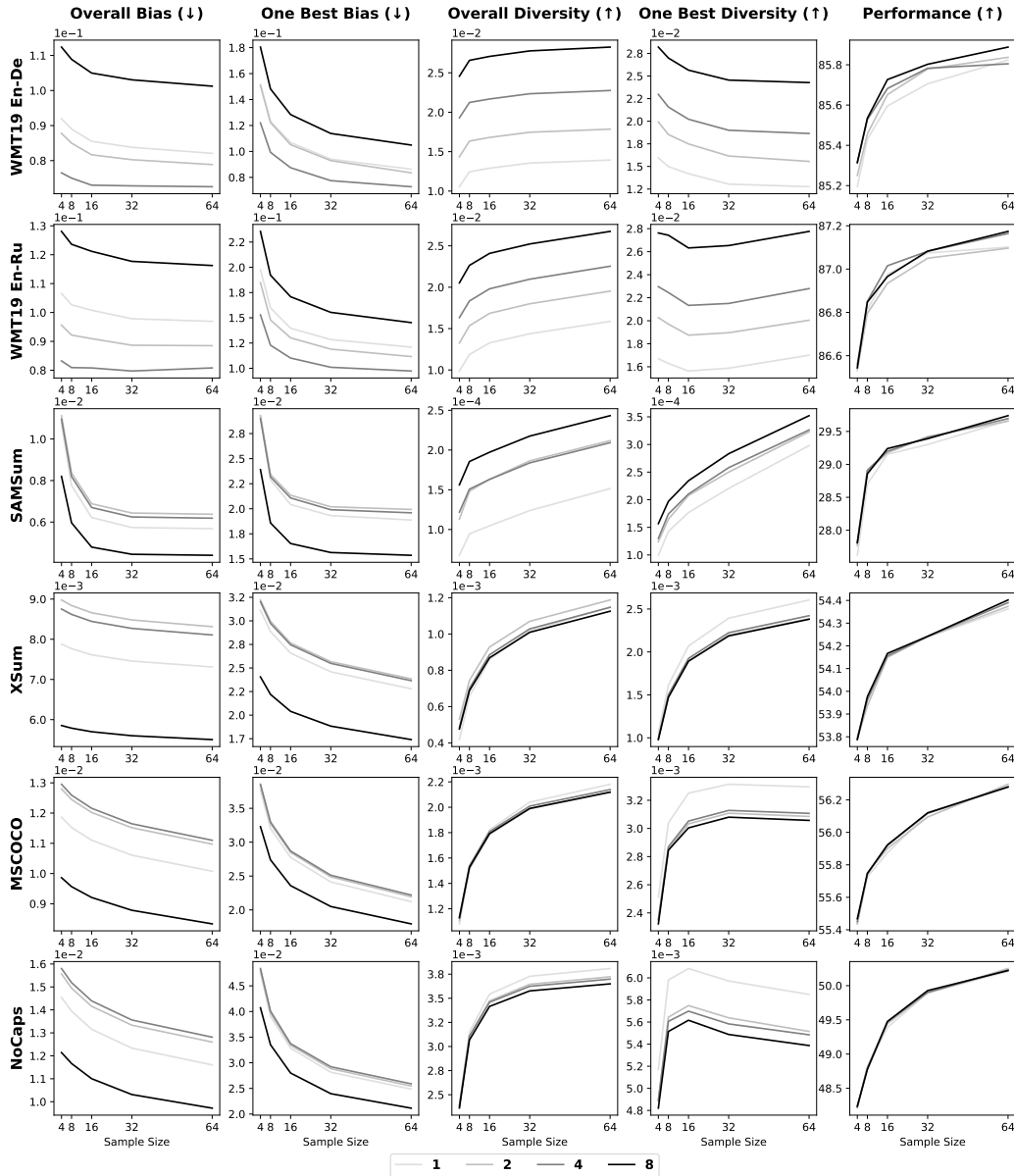


Figure 7: The relationship between bias, diversity, and performance in MAMBR decoding with pseudo-references generated by beam decoding. The notations are the same as Figure 6.

1404
 1405
 1406
 1407
 1408
 1409
 1410
 1411
 1412
 1413
 1414
 1415
 1416
 1417
 1418
 1419
 1420
 1421
 1422
 1423
 1424
 1425
 1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

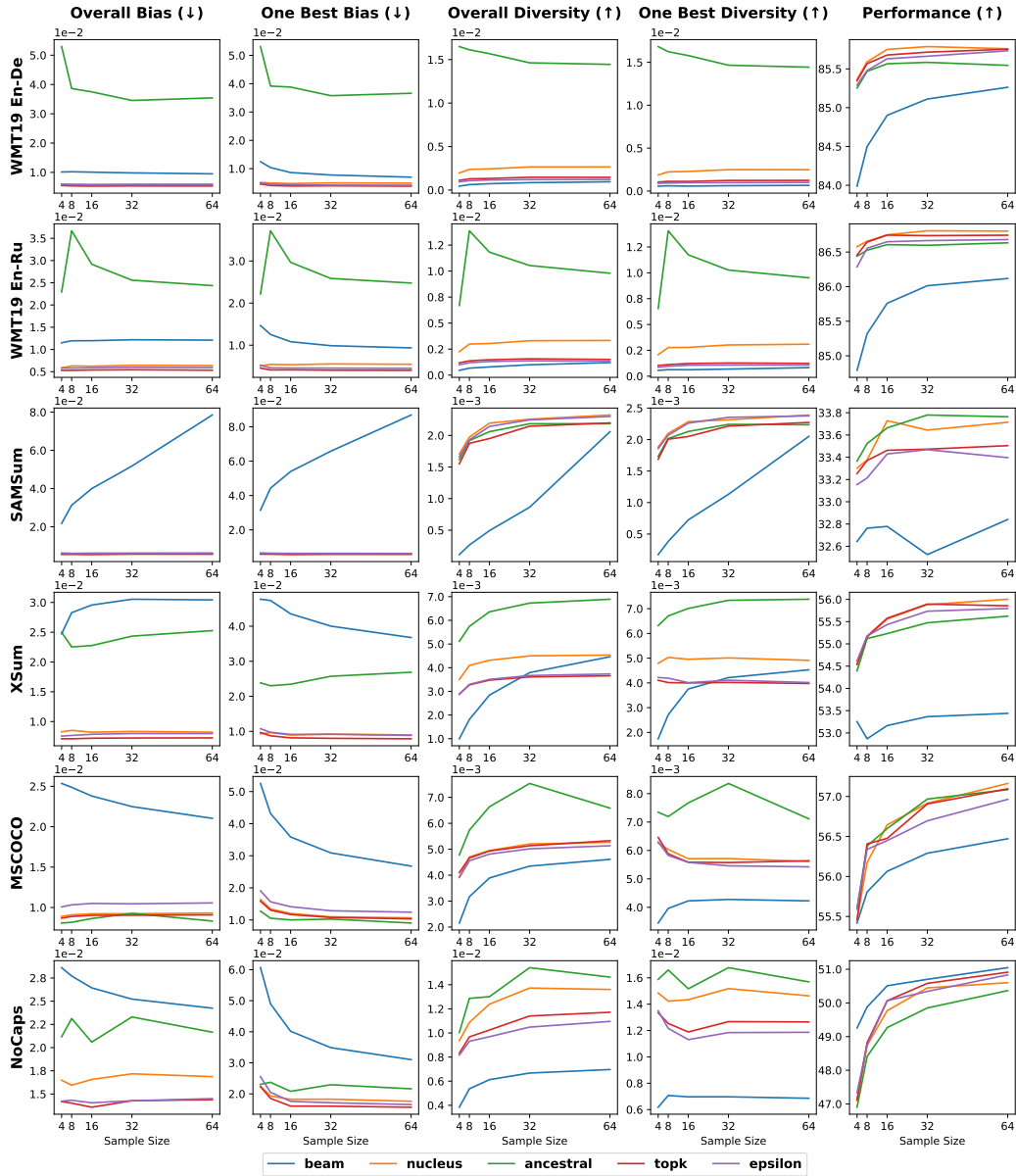


Figure 8: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with hypotheses generated by beam decoding. The notations are the same as Figure 2.

1458
 1459
 1460
 1461
 1462
 1463
 1464
 1465
 1466
 1467
 1468
 1469
 1470
 1471
 1472
 1473
 1474
 1475
 1476
 1477
 1478
 1479
 1480
 1481
 1482
 1483
 1484
 1485
 1486
 1487
 1488
 1489
 1490
 1491
 1492
 1493
 1494
 1495
 1496
 1497
 1498
 1499
 1500
 1501
 1502
 1503
 1504
 1505
 1506
 1507
 1508
 1509
 1510
 1511

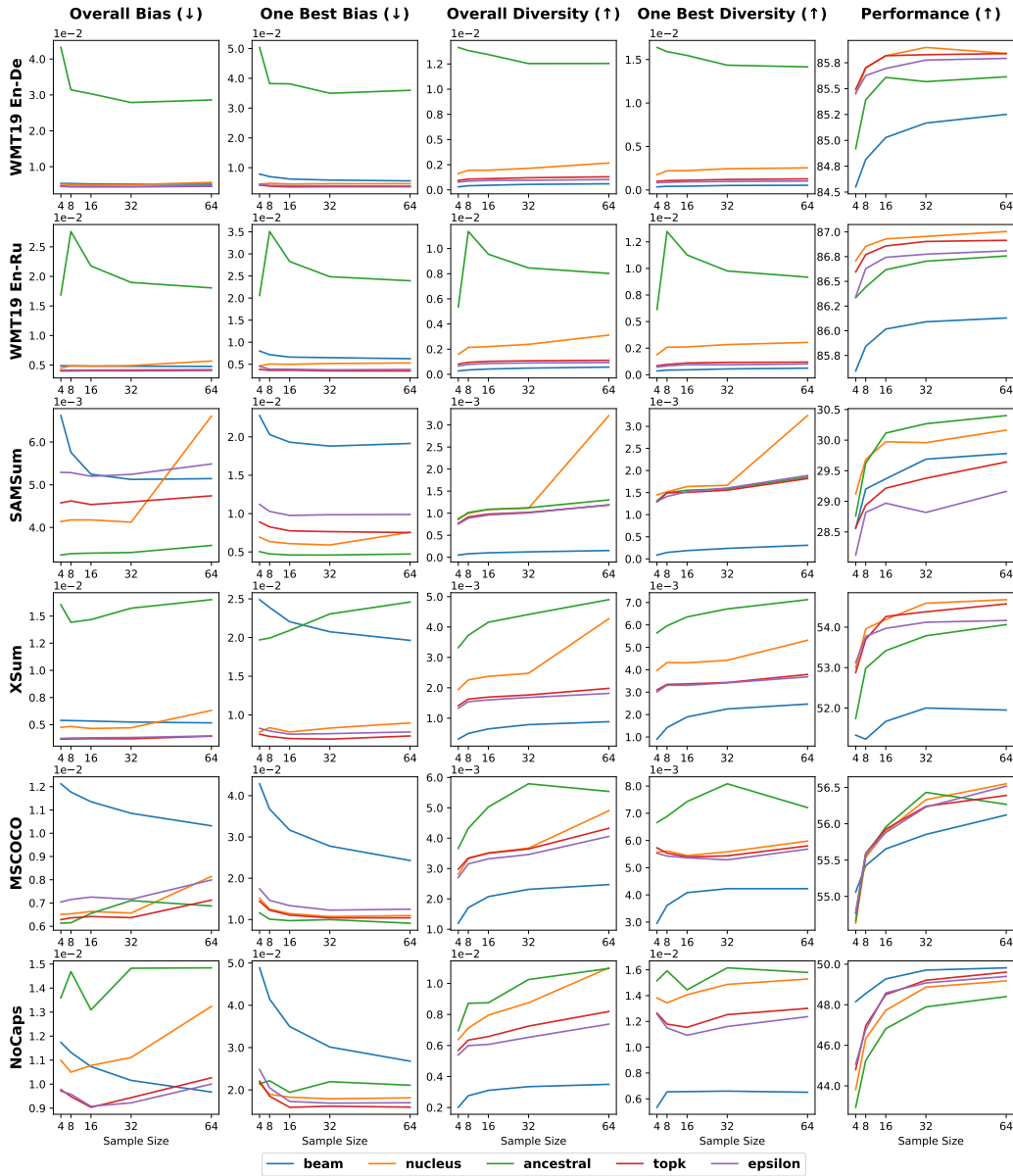


Figure 9: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with hypotheses generated by nucleus sampling. The notations are the same as Figure 2.

1512
 1513
 1514
 1515
 1516
 1517
 1518
 1519
 1520
 1521
 1522
 1523
 1524
 1525
 1526
 1527
 1528
 1529
 1530
 1531
 1532
 1533
 1534
 1535
 1536
 1537
 1538
 1539
 1540
 1541
 1542
 1543
 1544
 1545
 1546
 1547
 1548
 1549
 1550
 1551
 1552
 1553
 1554
 1555
 1556
 1557
 1558
 1559
 1560
 1561
 1562
 1563
 1564
 1565

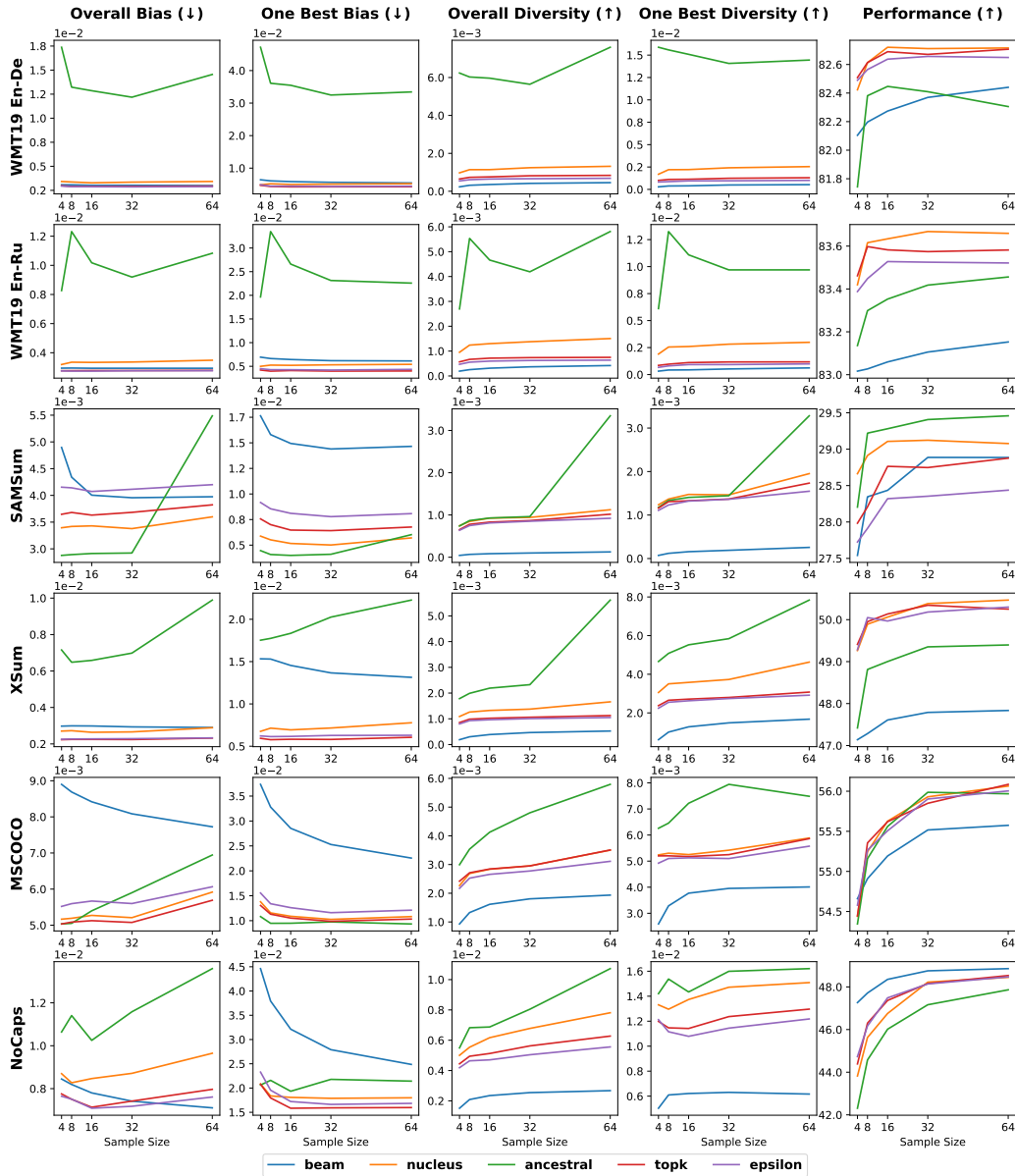


Figure 10: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with hypotheses generated by ancestral sampling. The notations are the same as Figure 2.

1566
 1567
 1568
 1569
 1570
 1571
 1572
 1573
 1574
 1575
 1576
 1577
 1578
 1579
 1580
 1581
 1582
 1583
 1584
 1585
 1586
 1587
 1588
 1589
 1590
 1591
 1592
 1593
 1594
 1595
 1596
 1597
 1598
 1599
 1600
 1601
 1602
 1603
 1604
 1605
 1606
 1607
 1608
 1609
 1610
 1611
 1612
 1613
 1614
 1615
 1616
 1617
 1618
 1619

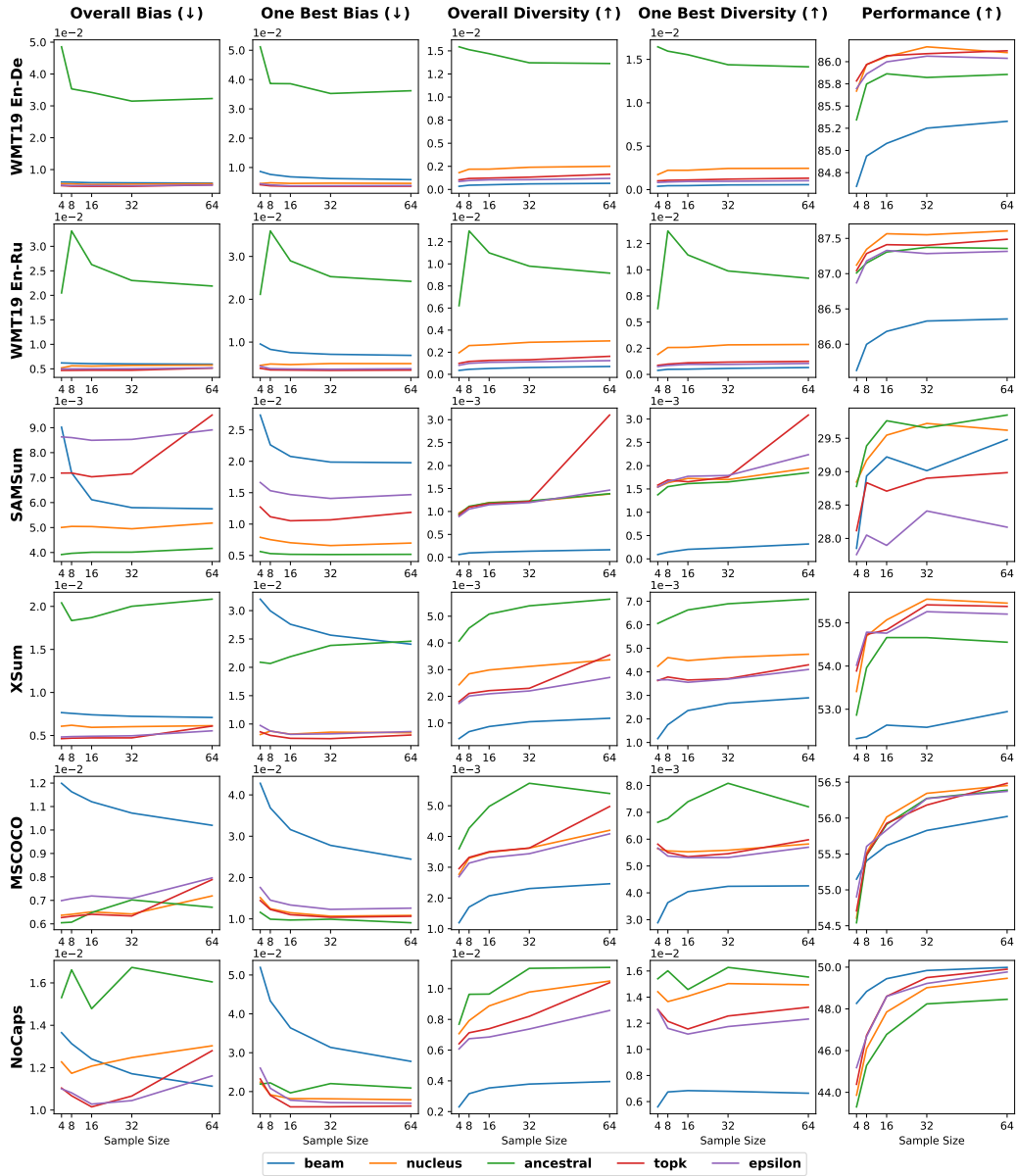


Figure 11: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with hypotheses generated by top-k sampling. The notations are the same as Figure 2.

1620
 1621
 1622
 1623
 1624
 1625
 1626
 1627
 1628
 1629
 1630
 1631
 1632
 1633
 1634
 1635
 1636
 1637
 1638
 1639
 1640
 1641
 1642
 1643
 1644
 1645
 1646
 1647
 1648
 1649
 1650
 1651
 1652
 1653
 1654
 1655
 1656
 1657
 1658
 1659
 1660
 1661
 1662
 1663
 1664
 1665
 1666
 1667
 1668
 1669
 1670
 1671
 1672
 1673

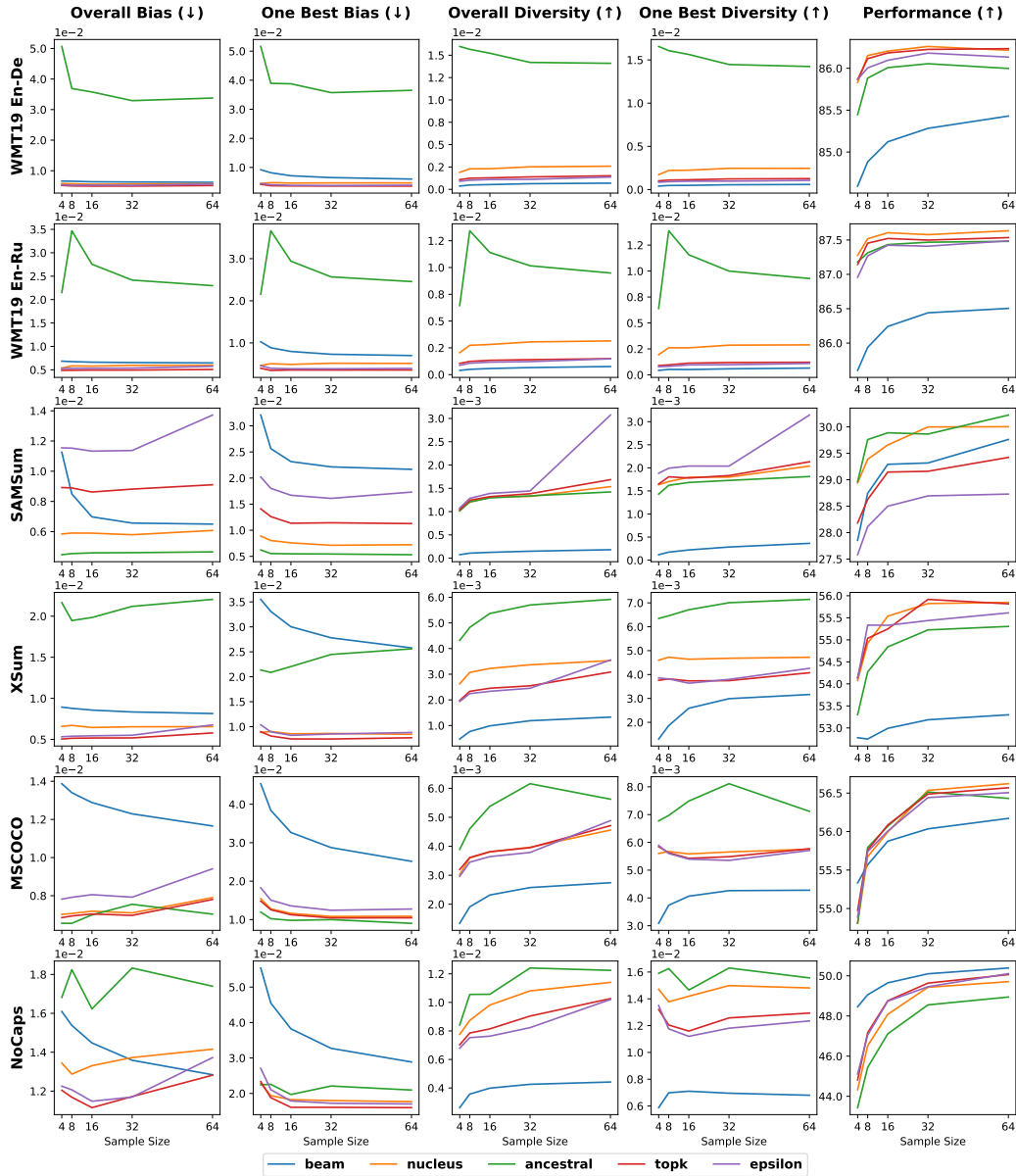


Figure 12: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with hypotheses generated by epsilon sampling. The notations are the same as Figure 2.

Table 4: Results of MAMBR with samples generated by ancestral sampling. Notations are the same as Table 1.

		WMT19 En-De					WMT19 En-Ru				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	84.5	84.7	84.8	85.0	85.1	85.8	86.1	86.3	86.5	86.5
Num. of Models	2	84.5	84.8	85.0	85.1	85.2	85.8	86.1	86.3	86.4	86.5
	4	84.6	84.8	85.0	85.1	85.2	85.8	86.1	86.4	86.5	86.6
	8	84.7	84.8	85.0	85.1	85.3	85.8	86.1	86.3	86.5	86.6
		SAMSum					XSum				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	28.6	29.1	29.5	29.5	29.7	53.3	54.1	55.0	55.1	55.2
Num. of Models	2	28.8	29.6	29.9	29.9	30.1	53.2	54.2	55.0	55.3	55.3
	4	28.7	29.5	29.9	29.8	30.2	53.3	54.2	55.0	55.3	55.3
	8	28.7	29.5	29.8	29.9	30.1	53.4	54.3	55.0	55.2	55.3
		MSCOCO					NoCaps				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	54.8	55.8	56.1	56.4	56.6	43.2	45.3	46.9	48.2	49.0
Num. of Models	2	54.8	55.8	56.2	56.4	56.6	43.4	45.7	47.2	48.7	49.1
	4	54.8	55.9	56.4	56.5	56.8	43.8	45.8	47.5	48.8	49.5
	8	54.9	55.9	56.3	56.6	56.8	43.9	45.9	47.4	49.0	49.5

Table 5: Results of MAMBR with samples generated by epsilon sampling. Notations are the same as Table 1.

		WMT19 En-De					WMT19 En-Ru				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	85.2	85.4	85.6	85.6	85.6	86.7	87.0	87.1	87.0	87.1
Num. of Models	2	85.3	85.5	85.6	85.7	85.7	86.7	87.0	87.1	87.1	87.1
	4	85.3	85.5	85.6	85.7	85.7	86.7	87.0	87.1	87.0	87.1
	8	85.3	85.6	85.7	85.7	85.8	86.7	87.1	87.1	87.0	87.2
		SAMSum					XSum				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	27.5	27.9	28.3	28.4	28.5	54.1	55.1	55.1	55.4	55.5
Num. of Models	2	27.7	28.1	28.5	28.6	28.7	54.1	55.2	55.1	55.4	55.4
	4	27.7	28.2	28.5	28.6	28.6	54.1	55.2	55.2	55.5	55.6
	8	27.6	28.2	28.6	28.6	28.7	54.2	55.2	55.3	55.4	55.6
		MSCOCO					NoCaps				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
	1	55.0	55.8	56.0	56.4	56.6	45.0	46.9	48.6	49.5	49.8
Num. of Models	2	55.0	55.7	56.1	56.4	56.6	45.0	47.1	48.9	49.5	50.1
	4	55.0	55.9	56.2	56.6	56.7	45.2	47.2	48.9	49.8	50.2
	8	55.1	55.8	56.2	56.5	56.8	45.2	47.2	49.0	49.8	50.3

1728
 1729
 1730
 1731
 1732
 1733
 1734
 1735
 1736
 1737
 1738
 1739
 1740
 1741
 1742
 1743
 1744
 1745
 1746
 1747
 1748
 1749
 1750
 1751
 1752
 1753
 1754
 1755
 1756
 1757
 1758
 1759
 1760
 1761
 1762
 1763
 1764
 1765
 1766
 1767
 1768
 1769
 1770
 1771
 1772
 1773
 1774
 1775
 1776
 1777
 1778
 1779
 1780
 1781

Table 6: Results of MAMBR with samples generated by beam decoding. Notations are the same as Table 1.

		WMT19 En-De					WMT19 En-Ru				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
1		85.1	85.3	85.4	85.4	85.3	86.7	86.9	86.9	86.9	86.9
Num. of Models	2	85.1	85.4	85.5	85.4	85.4	86.7	86.8	86.9	86.9	87.0
	4	85.2	85.4	85.5	85.4	85.4	86.7	86.8	86.9	86.9	86.9
	8	85.2	85.5	85.5	85.5	85.5	86.6	86.9	86.9	86.9	86.9
		SAMSum					XSum				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
1		27.6	28.7	29.2	29.3	29.7	52.8	52.8	53.1	53.2	53.3
Num. of Models	2	27.8	28.9	29.2	29.4	29.7	52.7	52.9	53.0	53.2	53.4
	4	27.8	28.9	29.2	29.4	29.7	52.8	52.8	53.0	53.2	53.3
	8	27.8	28.9	29.2	29.4	29.7	52.7	52.9	53.1	53.2	53.3
		MSCOCO					NoCaps				
Num. of Samples		4	8	16	32	64	4	8	16	32	64
1		55.3	55.5	55.9	56.0	56.2	48.5	49.0	49.6	50.1	50.5
Num. of Models	2	55.3	55.5	55.8	56.1	56.2	48.5	49.0	49.8	50.2	50.5
	4	55.3	55.5	55.8	56.0	56.3	48.5	49.1	49.7	50.2	50.5
	8	55.3	55.5	55.8	56.1	56.2	48.5	49.2	49.7	50.3	50.5

1782
 1783
 1784
 1785
 1786
 1787
 1788
 1789
 1790
 1791
 1792
 1793
 1794
 1795
 1796
 1797
 1798
 1799
 1800
 1801
 1802
 1803
 1804
 1805
 1806
 1807
 1808
 1809
 1810
 1811
 1812
 1813
 1814
 1815
 1816
 1817
 1818
 1819
 1820
 1821
 1822
 1823
 1824
 1825
 1826
 1827
 1828
 1829
 1830
 1831
 1832
 1833
 1834
 1835

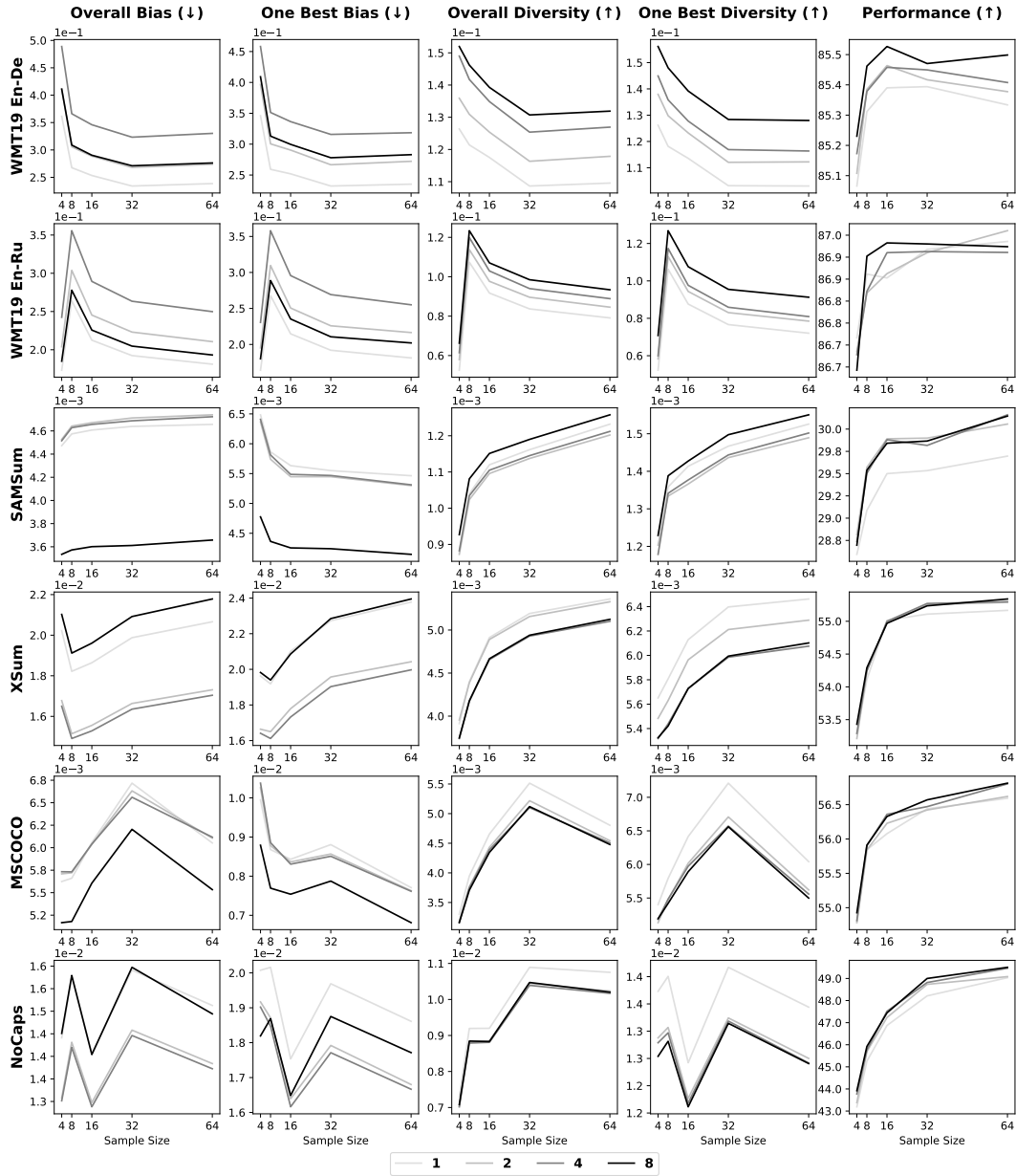


Figure 13: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with pseudo-references generated by ancestral sampling. The notations are the same as Figure 6.

1836
 1837
 1838
 1839
 1840
 1841
 1842
 1843
 1844
 1845
 1846
 1847
 1848
 1849
 1850
 1851
 1852
 1853
 1854
 1855
 1856
 1857
 1858
 1859
 1860
 1861
 1862
 1863
 1864
 1865
 1866
 1867
 1868
 1869
 1870
 1871
 1872
 1873
 1874
 1875
 1876
 1877
 1878
 1879
 1880
 1881
 1882
 1883
 1884
 1885
 1886
 1887
 1888
 1889

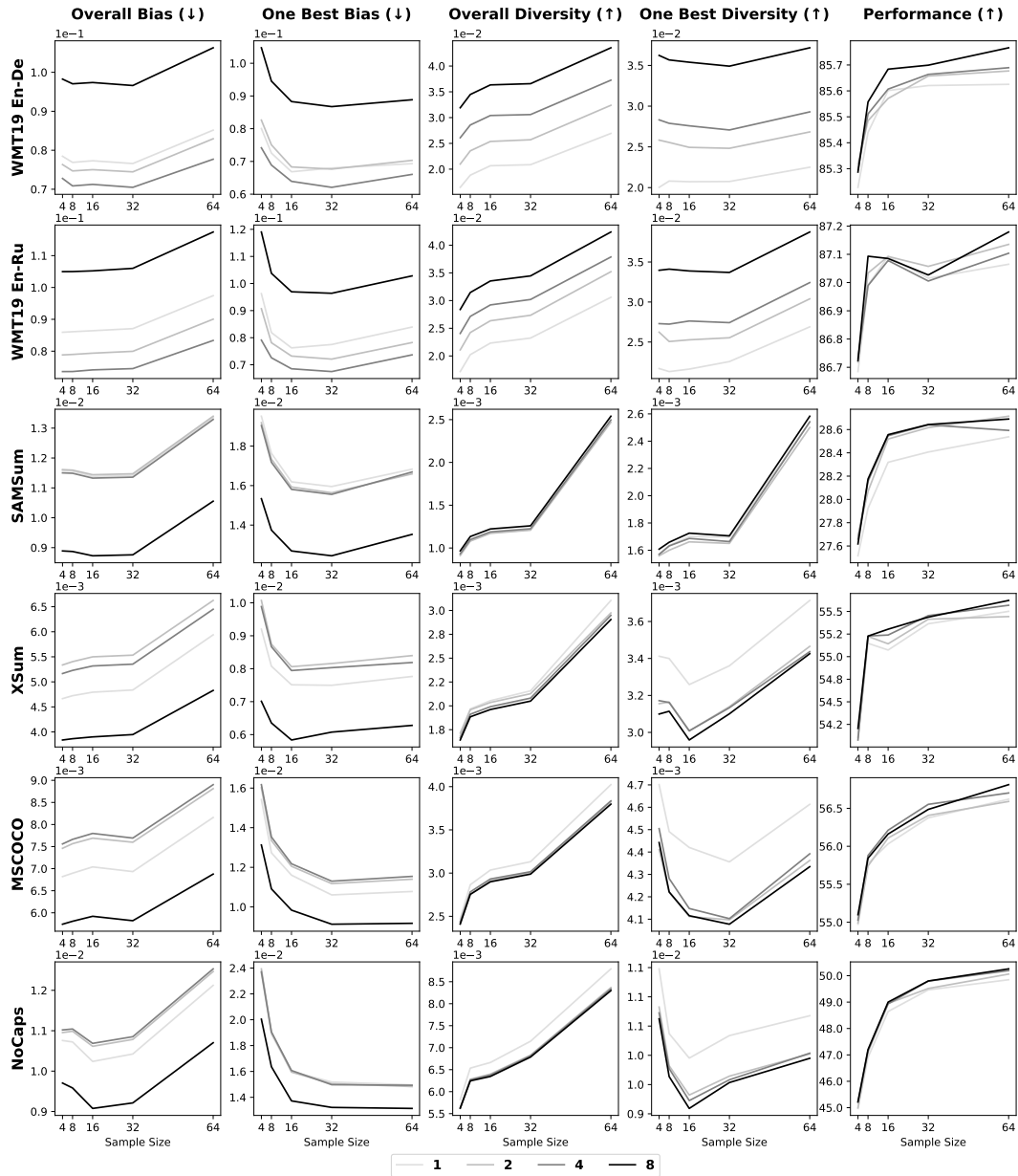


Figure 14: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with pseudo-references generated by epsilon sampling. The notations are the same as Figure 6.

1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943

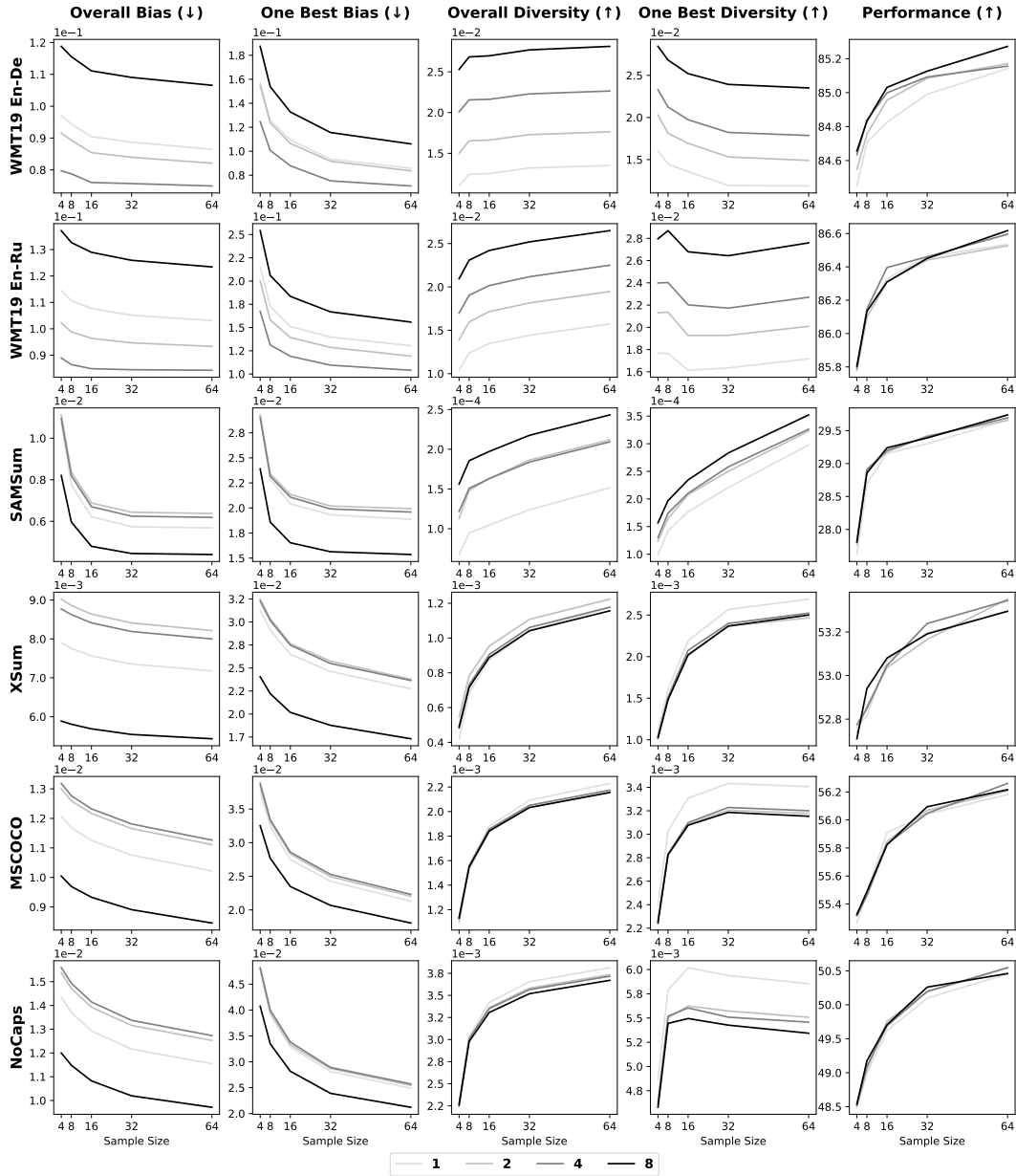


Figure 15: The relationship between bias, diversity, and performance on the first 1000 lines of each dataset in MBR decoding with pseudo-references generated by beam decoding. The notations are the same as Figure 6.