

Agreement-Constrained Probabilistic Minimum Bayes Risk Decoding

Koki Natsumi[†], Hiroyuki Deguchi[‡],

Yusuke Sakai[†], Hidetaka Kamigaito[†], Taro Watanabe[†]

[†]Nara Institute of Science and Technology (NAIST) [‡]NTT, Inc.

natsumi.koki.ng5@naist.ac.jp, hiroyuki.deguchi@ntt.com

{sakai.yusuke.sr9, kamigaito.h, taro}@is.naist.jp

Abstract

Minimum Bayes risk (MBR) decoding generates high-quality translations by maximizing the expected utility of output candidates, but it evaluates all pairwise scores over the candidate set; hence, it takes quadratic time with respect to the number of candidates. To reduce the number of utility function calls, probabilistic MBR (PMBR) decoding partially evaluates quality scores using sampled pairs of candidates and completes the missing scores with a matrix completion algorithm. Nevertheless, it degrades the translation quality as the number of utility function calls is reduced. Therefore, to improve the trade-off between quality and cost, we propose agreement-constrained PMBR (AC-PMBR) decoding, which leverages a knowledge distilled model to guide the completion of the score matrix. Our AC-PMBR decoding improved approximation errors of matrix completion by up to 3 times and achieved higher translation quality compared with PMBR decoding at a comparable computational cost on the WMT’23 En↔De translation tasks.

1 Introduction

Maximum a posteriori (MAP) decoding, which finds the most probable candidate, is the standard inference strategy in translation tasks, while such high-probability translations do not always align with human assessment (Koehn and Knowles, 2017; Eikema and Aziz, 2020). To overcome the limitation, minimum Bayes risk (MBR) decoding selects a high-quality translation rather than a high-probability one by maximizing expected utility (Goel and Byrne, 2000; Kumar and Byrne, 2004). For estimating expected utility, it calculates the utility score matrix, evaluating all candidates against multiple pseudo-references, which are sample translations drawn from the output distribution. Thus, it requires utility function calls proportional to the square of the number of candidates and is computationally expensive, especially when using

neural metrics that highly correlate with human assessment, e.g., BLEURT (Sellam et al., 2020).

Recent studies improve the efficiency of MBR decoding (Cheng and Vlachos, 2023; Jinnai and Ariu, 2024; Deguchi et al., 2024b; Vamvas and Sennrich, 2024; Trabelsi et al., 2024). Among them, probabilistic MBR (PMBR) decoding (Trabelsi et al., 2024) drastically reduces the number of utility function calls by completing the score matrix using partially observed scores. Nevertheless, as the number of utility function calls is reduced, the approximation error of matrix completion increases, and the translation quality deteriorates. That is, there exists a trade-off between completion accuracy and computational cost.

To relax this trade-off, we propose *agreement-constrained PMBR (AC-PMBR)* decoding, which facilitates score matrix completion by leveraging a knowledge distilled metric. Our agreement constraint minimizes the difference between the target and distilled low-rank matrices, thereby reducing the approximation error in matrix completion.

Experiments demonstrated that our AC-PMBR decoding improved the matrix completion accuracy by up to 3 times in mean squared error (MSE) against the full score matrix and achieved higher translation quality compared with PMBR decoding at comparable costs in the WMT’23 En↔De translation tasks (Kocmi et al., 2023).

2 Background

MBR decoding MBR decoding finds higher-quality translations than widely used MAP decoding, such as N -best beam search, by maximizing the expected utility of output candidates (Kumar and Byrne, 2004; Eikema and Aziz, 2020). Let \mathcal{T} be all possible translations. The goal of MBR decoding is to find the translation that maximizes the expected utility, i.e., $\operatorname{argmax}_{y \in \mathcal{T}} \mathbb{E}_{\hat{y} \sim \Pr(y|x)}[u(y, \hat{y})]$, where x is an in-

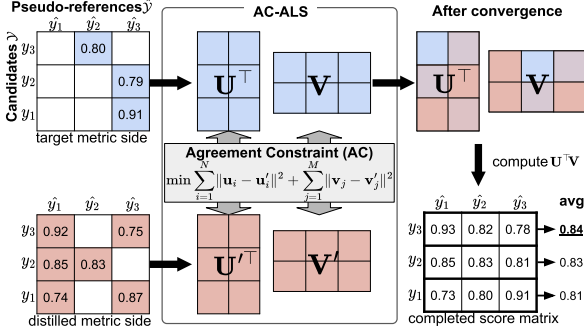


Figure 1: Overview of our proposed Agreement-Constrained PMBR (AC-PMBR) decoding.

put text, $\Pr(y|x)$ is the true translation probability, and $u: \mathcal{T} \times \mathcal{T} \rightarrow \mathbb{R}$ denotes a utility function that satisfies $y \succeq y' \iff u(y, \hat{y}) \geq u(y', \hat{y})$ under the preference relation \succeq . Since enumerating all $y \in \mathcal{T}$ is infeasible and calculating the true probability \Pr is unknown, MBR decoding estimates the expected utility using sample translations drawn from the model output probability, called pseudo-references $\hat{\mathcal{Y}} := \{\hat{y}_1, \dots, \hat{y}_M\} \subset \mathcal{T}$, and selects the translation from a candidate set $\mathcal{Y} := \{y_1, \dots, y_N\} \subset \mathcal{T}$. The expected utility is typically estimated by the Monte Carlo method (Eikema and Aziz, 2022) with a score matrix $\mathbf{O} := [O_{ij} = u(y_i, \hat{y}_j)] \in \mathbb{R}^{N \times M}$, and then, the best candidate is selected, i.e., $y_{\text{MBR}} := \arg\max_{y_i \in \mathcal{Y}} \frac{1}{M} \sum_{j=1}^M O_{ij}$.

MBR decoding generates high-quality translations, while its time complexity is $\mathcal{O}(NM)$. Recent studies often employ $N \geq 1,000$ (Freitag et al., 2023), making it extremely slow.

PMBR decoding Probabilistic MBR (PMBR) decoding accelerates MBR decoding by reducing the number of utility function calls (Trabelsi et al., 2024). It does not evaluate scores for all O_{ij} ; instead, it partially evaluates only sampled pairs of hypotheses and pseudo-references. The other missing scores are completed using a low-rank matrix factorization from a partially observed score matrix $\tilde{\mathbf{O}} \in \mathbb{R}^{N \times M}$. Specifically, the incomplete matrix $\tilde{\mathbf{O}}$ is approximated by the matrix multiplication of two d -dimensional low-rank matrices $\mathbf{U} \in \mathbb{R}^{d \times N}$ and $\mathbf{V} \in \mathbb{R}^{d \times M}$, i.e., $\tilde{\mathbf{O}} \approx \mathbf{U}^\top \mathbf{V}$. Here, $\mathbf{u}_i, \mathbf{v}_j \in \mathbb{R}^d$ are d -dimensional column vectors, and $\mathbf{U} = [\mathbf{u}_1; \dots; \mathbf{u}_N]$ and $\mathbf{V} = [\mathbf{v}_1; \dots; \mathbf{v}_M]$ stack the rank reduced vectors for the row and column directions of $\tilde{\mathbf{O}}$, respectively. Let $\text{Obs}(\tilde{\mathbf{O}}) := \{(i, j) \mid \tilde{O}_{ij} \text{ is observed}\}$ be the set of observed indices in $\tilde{\mathbf{O}}$. We obtain \mathbf{U} and \mathbf{V} that minimize

the following objective:

$$\mathcal{L}_{\text{MF}}(\mathbf{U}, \mathbf{V}; \tilde{\mathbf{O}}) = \sum_{(i,j) \in \text{Obs}(\tilde{\mathbf{O}})} (\tilde{O}_{ij} - \mathbf{u}_i^\top \mathbf{v}_j)^2 + \lambda \left(\sum_{i=1}^N \|\mathbf{u}_i\|^2 + \sum_{j=1}^M \|\mathbf{v}_j\|^2 \right), \quad (1)$$

where $\lambda \in \mathbb{R}_+$ is a weight of the regularization term. This optimization is solved by the alternating least-squares (ALS) algorithm (Zachariah et al., 2012), and the complete score matrix is obtained by $\mathbf{U}^\top \mathbf{V}$. PMBR decoding successfully reduces the number of utility function calls. Nevertheless, there is still a cost-quality trade-off, because further reductions in utility function calls significantly degrade approximation accuracy.

3 Agreement-Constrained PMBR decoding

We propose agreement-constrained PMBR (AC-PMBR) decoding, which alleviates the PMBR decoding cost-quality trade-off without increasing total cost by reallocating a fixed evaluation budget. Instead of adding the distilled metric on top of PMBR decoding, AC-PMBR decoding reduces target metric calls and spends the saved budget on many distilled metric calls. Thereby enabling more total utility function calls at the same computational cost as PMBR decoding and yielding higher matrix-completion accuracy of the MBR's score matrix. AC-PMBR decoding proceeds in two steps: (1) score matrix construction, and (2) agreement-constrained matrix completion, as illustrated in Figure 1.

Score matrix construction We compute the score matrices, $\tilde{\mathbf{O}} \in \mathbb{R}^{N \times M}$ and $\tilde{\mathbf{O}}' \in \mathbb{R}^{N \times M}$, with the target and its distilled metrics, respectively. Hereafter, we denote a prime $'$ for the distilled metric side. Let r and r' denote the reduction rates; we observe only a $1/r$ and $1/r'$ fraction of the $N \times M$ entries in $\tilde{\mathbf{O}}$ and $\tilde{\mathbf{O}}'$, respectively. The time complexity of evaluating the partially observed samples of hypotheses and pseudo-references is $\mathcal{O}(\frac{NM}{r})$. As the reduction rate r or r' increases, the number of observed samples decreases. To alleviate the cost-quality trade-off, we set $r > r'$, i.e., we call a distilled metric more frequently than an expensive target metric for denser guidance at almost the same cost.

Agreement-constrained matrix completion We factorize the matrices $\tilde{\mathbf{O}}$ and $\tilde{\mathbf{O}}'$ with the alternating least squares (ALS) algorithm (Zachariah

Algorithm 1 Agreement-constrained ALS

Require: Regularization weight $\lambda \in \mathbb{R}_+$, agreement weight $\gamma \in \mathbb{R}_+$, rank $d \in \mathbb{N}$, and identity matrix $\mathbf{I} \in \mathbb{R}^{d \times d}$

Ensure: $\mathbf{U} \in \mathbb{R}^{d \times N}$ and $\mathbf{V} \in \mathbb{R}^{d \times M}$

```
1: repeat
2:   Initialize  $\mathbf{U}, \mathbf{U}' \in \mathbb{R}^{d \times N}$  and  $\mathbf{V}, \mathbf{V}' \in \mathbb{R}^{d \times M}$ 
3:   for  $i = 1 \dots N$  do
4:      $\mathbf{M}' = \text{diag}(\mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}}')[(i, 1)]}, \dots, \mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}}')[(i, M)]})$ 
5:      $\mathbf{M} = \text{diag}(\mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}})[(i, 1)]}, \dots, \mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}})[(i, M)]})$ 
6:      $\mathbf{u}'_i = (\mathbf{V}'\mathbf{M}'\mathbf{V}'^\top + (\lambda + \gamma)\mathbf{I})^{-1}(\mathbf{V}'\mathbf{M}'\tilde{\mathbf{O}}'_{i*} + \gamma\mathbf{u}_i)$ 
7:      $\mathbf{u}_i = (\mathbf{V}\mathbf{M}\mathbf{V}^\top + (\lambda + \gamma)\mathbf{I})^{-1}(\mathbf{V}\mathbf{M}\tilde{\mathbf{O}}_{i*} + \gamma\mathbf{u}'_i)$ 
8:   end for
9:   for  $j = 1 \dots M$  do
10:     $\mathbf{N}' = \text{diag}(\mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}}')[(1, j)]}, \dots, \mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}}')[(N, j)]})$ 
11:     $\mathbf{N} = \text{diag}(\mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}})[(1, j)]}, \dots, \mathbb{1}_{\text{Obs}(\tilde{\mathbf{O}})[(N, j)]})$ 
12:     $\mathbf{v}'_j = (\mathbf{U}'\mathbf{N}'\mathbf{U}'^\top + (\lambda + \gamma)\mathbf{I})^{-1}(\mathbf{U}'\mathbf{N}'\tilde{\mathbf{O}}'_{*j} + \gamma\mathbf{v}_j)$ 
13:     $\mathbf{v}_j = (\mathbf{U}\mathbf{N}\mathbf{U}^\top + (\lambda + \gamma)\mathbf{I})^{-1}(\mathbf{U}\mathbf{N}\tilde{\mathbf{O}}_{*j} + \gamma\mathbf{v}'_j)$ 
14:   end for
15: until convergence
16: return  $\mathbf{U}, \mathbf{V}$ 
```

et al., 2012), i.e., we minimize $\mathcal{L}_{\text{MF}}(\mathbf{U}, \mathbf{V}, \tilde{\mathbf{O}})$ and $\mathcal{L}_{\text{MF}}(\mathbf{U}', \mathbf{V}', \tilde{\mathbf{O}}')$ with our proposed agreement constraint. The constraint encourages the rank reduced representation on the target metric to be closer to that of its knowledge distilled metric:

$$\begin{aligned} \mathcal{L}_{\text{AC}}(\mathbf{U}, \mathbf{V}, \mathbf{U}', \mathbf{V}') \\ = \sum_{i=1}^N \|\mathbf{u}_i - \mathbf{u}'_i\|^2 + \sum_{j=1}^M \|\mathbf{v}_j - \mathbf{v}'_j\|^2. \end{aligned} \quad (2)$$

Formally, our AC-PMBR decoding minimizes the following objective:

$$\begin{aligned} \arg\min_{\mathbf{U}, \mathbf{V}} \mathcal{L}_{\text{MF}}(\mathbf{U}, \mathbf{V}; \tilde{\mathbf{O}}) + \mathcal{L}_{\text{MF}}(\mathbf{U}', \mathbf{V}'; \tilde{\mathbf{O}}') \\ + \gamma \mathcal{L}_{\text{AC}}(\mathbf{U}, \mathbf{V}, \mathbf{U}', \mathbf{V}'), \end{aligned} \quad (3)$$

where $\gamma \in \mathbb{R}_+$ controls the strength of the agreement constraint. This constrained optimization problem can be solved by extending the ALS algorithm, as shown in Algorithm 1; a detailed derivation is provided in Appendix C. The rank reduced representations of the distilled metric side are first updated, as shown in Algorithm 1 and 12, so that those of the target metric side are not affected by the unupdated matrices that do not have meaningful information. Now, we complete the incomplete score matrix by multiplying matrices \mathbf{U} and \mathbf{V} , calculated in Equation 3, i.e., we use $\mathbf{U}^\top \mathbf{V}$ as the completed score matrix and estimate the expected utility in the same way as MBR decoding.

4 Experimental Settings

We provide the experimental settings below; further details are available in Appendix A.

Evaluation We conduct experiments on the WMT'23 En \leftrightarrow De translation tasks (Kocmi et al., 2023). We evaluate translation quality using BLEURT (Sellam et al., 2020), XCOMET (Guerreiro et al., 2024), BLEU (Papineni et al., 2002), chrF (Popović, 2015), and MetricX (Juraska et al., 2024). To assess the performance of each method in completing the score matrix, we compute the mean squared error (MSE) between the matrix completed by ALS or Agreement-constrained ALS and the ground-truth matrix.

Methods We compare AC-PMBR decoding with PMBR decoding, and also evaluate MAP and MBR decoding. We also evaluate the upper bound of translation quality in candidate sets by selecting translations that maximize the target metric using references (Oracle). Following Deguchi et al. (2024a), we sampled 1,024 translation candidates via ε -sampling with $\varepsilon = 0.02$ (Freitag et al., 2023) from M2M100 (Fan et al., 2020) and used the same set as pseudo-references.

Utility function We employ BLEURT-20 as the target metric, along with its three distilled versions, BLEURT-20-{D3, D6, D12} (Pu et al., 2021), as the distilled metrics. To ensure comparable computational costs between PMBR and AC-PMBR decoding, we choose the reduction rates r and r' in two settings: a low-reduction (**Low**) setting, where we use $r = 16$ for PMBR and $r = 32$ for AC-PMBR decoding, with $r' = \{2, 4, 8\}$ for BLEURT-20-{D3, D6, D12}, respectively; and a high-reduction (**High**) setting, where we use $r = 512$ for PMBR and $r = 1,024$ for AC-PMBR decoding, with $r' = \{64, 128, 256\}$, respectively.

Hyperparameter The agreement weight γ and regularization weight λ were tuned on the WMT'22 En \rightarrow De translation task (Kocmi et al., 2022) by minimizing the MSE of the score matrices. ALS and agreement-constrained ALS were run for up to 30 iterations or until the loss difference fell below 10^{-4} . For both PMBR and AC-PMBR decoding, we fixed the rank $d = 8$ and $\lambda = 0.1$. In AC-PMBR decoding, we also tuned $\gamma \in \{0.1, \dots, 1.0\}$ for each reduction rate r , and used the optimal values $\gamma = 0.1$ for $r = 32$ and $\gamma = 1.0$ for $r = 1,024$.

5 Experimental Results and Discussions

Table 1 shows the main results.

Decoding	Dist.	En→De						De→En					
		BLRT↑	XCT↑	BLEU↑	chrF↑	MX↓	MSE↓	BLRT↑	XCT↑	BLEU↑	chrF↑	MX↓	MSE↓
MAP	–	45.27	59.61	10.74	30.38	12.94	–	56.27	65.79	16.56	37.39	11.68	–
MBR	–	57.42	67.83	18.97	46.19	8.87	–	65.19	77.49	23.86	50.99	8.46	–
	D3	46.68	57.20	18.12	46.69	10.79	10.38	60.65	70.67	23.37	51.89	9.66	8.53
	D6	48.89	59.54	18.31	46.49	10.33	9.43	61.80	73.42	24.28	51.37	9.38	6.97
	D12	51.33	63.54	17.73	44.53	9.79	9.58	62.21	74.28	23.92	50.55	9.45	6.53
Reduction rate: Low PMBR: $r = 16$, AC-PMBR: $r = 32$													
PMBR	–	57.19	67.79	19.05	46.21	8.81	3.04	64.87	77.27	23.47	50.72	8.46	2.54
AC-PMBR	D3	57.01	67.57	19.19	46.60	8.80	3.23	64.87	77.01	23.82	51.08	8.51	2.80
	D6	57.26	68.00	19.14	46.27	8.74	3.12	64.87	77.26	23.69	50.95	8.51	2.68
	D12	57.29	67.94	19.09	46.23	8.77	3.10	64.95	77.47	24.05	51.10	8.46	2.65
Reduction rate: High PMBR: $r = 512$, AC-PMBR: $r = 1,024$													
PMBR	–	50.42	60.92	16.76	44.05	10.80	26.06	60.02	70.74	21.34	48.69	10.37	33.80
AC-PMBR	D3	50.49	60.56	18.15	46.17	10.42	10.29	61.08	71.98	22.96	50.50	9.71	11.89
	D6	51.21	61.92	18.23	45.75	10.12	11.16	61.02	72.10	22.83	49.94	9.90	12.57
	D12	51.81	63.41	17.54	44.48	10.12	15.78	61.26	72.66	22.31	49.12	9.97	16.80
Oracle	–	57.43	69.10	18.60	44.54	8.52	–	70.93	79.28	26.14	52.51	7.21	–

Table 1: Results of WMT’23 En↔De. BLEURT is abbreviated as BLRT, XCOMET as XCT, and MetricX as MX. “Dist.” indicates distilled metrics. The best scores within each reduction rate setting are highlighted in **bold**.

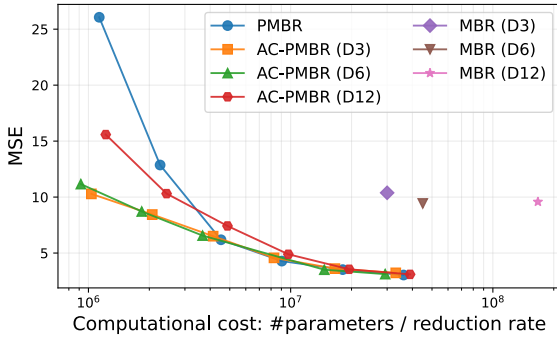


Figure 2: MSE of score matrices when varying computational costs in WMT’23 En→De. #parameters refers to the number of parameters in the evaluation metric model, serving as an indicator of model scale, while the computational cost is described by Equation 4 and 5.

5.1 Translation quality

Under the low-reduction setting, AC-PMBR decoding retains an advantage over baseline PMBR in both translation directions. It delivers gains of up to 0.6% BLEU and 0.4% chrF on surface-form metrics, and up to 0.2% on the semantic metric XCOMET. Under the high-reduction setting, PMBR’s quality fell sharply, whereas AC-PMBR decoding curbed that decline in both directions and stayed ahead on every metric. Across the two directions, AC-PMBR decoding delivered roughly 2.5% higher XCOMET and up to 2% higher chrF, clearly outperforming baseline PMBR decoding in the most demanding scenario. The results of the significance tests are reported in Appendix B.

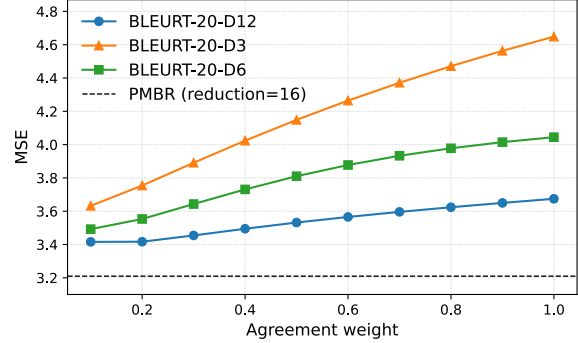


Figure 3: Agreement weight tuning in the low reduction rate setting on WMT’22 En→De.

5.2 Matrix completion accuracy

In the low-reduction setting, AC-PMBR decoding achieved performance comparable to PMBR decoding in MSE evaluation. In addition, the MSE was significantly improved in the high-reduction setting, with up to a 3 times improvement. This setting has the highest number of utility function calls, and therefore, it is more effective for MSE to evaluate a large number of pairs with low accuracy than to calculate a small number of pairs with a high accuracy model at low computational cost. Figure 2 also shows that AC-PMBR decoding suppresses MSE, which worsens as computational cost decreases. Moreover, AC-PMBR decoding achieved higher accuracy at a lower cost than MBR decoding with distilled metrics. This suggests that distilled metrics are effective when used to assist matrix completion via our agreement constraint,

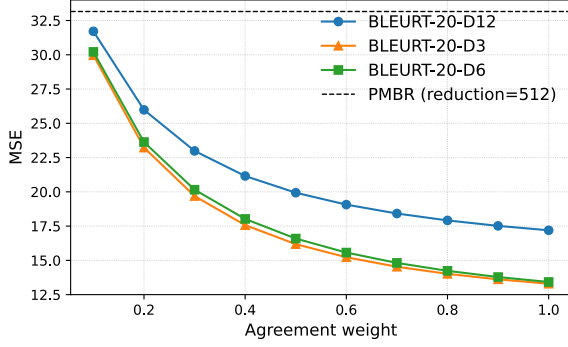


Figure 4: Agreement weight tuning in the high-reduction setting on WMT’22 En→De.

rather than being used for the utility function.

5.3 Effect of Reduction Rate on Translation

Table 1 shows that with the proposed AC-PMBR decoding, the decline in the approximation accuracy of the score matrix is more gradual compared to conventional PMBR decoding, especially as the reduction rate r increases. This suggests that AC-PMBR decoding can perform matrix completion more accurately than PMBR decoding at a comparable computational cost. Furthermore, as shown in Figures 3 and 4, when the reduction rate r is small, a low agreement weight is sufficient because the information from the target model alone is adequate for matrix completion. Conversely, at high reduction rates where the target model’s information is insufficient, the information from the knowledge distilled model effectively contributes to the completion process. The fact that its translation quality surpasses that of MBR decoding using only the knowledge distilled model also indicates that our method effectively incorporates information from the target model, even under high reduction rates.

5.4 Approximation Accuracy of the Score Matrix

As shown in Figure 2, the MSE evaluation of the score-matrix approximation clearly demonstrates the robustness of the proposed AC-PMBR decoding, especially at high reduction rates, where conventional PMBR decoding collapses. Under high reduction, the number of observed scores becomes critically small, turning matrix completion in PMBR decoding into an ill-posed problem and causing a sharp drop in approximation accuracy. This is reflected in a dramatic increase in MSE and a notable decline in translation quality. In contrast, AC-PMBR decoding maintains significantly lower

MSE under the same conditions. This robustness is attributed to the information from the dense score matrix provided by the knowledge-distilled model. We presume that this matrix, which captures the general distribution of the true score matrix, acts as a guide that prevents the approximation from failing. Given that the performance of both methods is comparable at low reduction rates, where observed scores are relatively abundant, AC-PMBR decoding is expected to perform effectively in more difficult, information-scarce situations.

6 Conclusion

In this study, we proposed AC-PMBR decoding, which assists score-matrix completion by aligning a target metric with its distilled metrics. We evaluated it on the WMT’23 En↔De translation tasks. AC-PMBR decoding mitigated the quality degradation observed in PMBR decoding, particularly under high-reduction settings, improved evaluation scores across metrics, and reduced approximation error by up to three times. Our study focused on distilled metrics to achieve a better cost–quality trade-off, but the framework inherently supports multi-metric ensembles, suggesting potential for multi-aspect decoding in future work.

Limitations

Evaluation This study primarily uses BLEURT as the target metric. While applying AC-PMBR decoding to other neural metrics, such as XCOMET, would be ideal for exploring broader robustness, this short paper emphasizes the proposal of AC-PMBR decoding and its metric aggregation framework, using BLEURT as a case study. Prior work on PMBR decoding reports consistent trends and strong scores across multiple language pairs, suggesting that the underlying phenomenon is largely language-independent (Trabelsi et al., 2024). Accordingly, we report results on the representative WMT’23 En↔De directions. We believe this already provides sufficient evidence to support our motivation. Similarly, although evaluating other tasks or language pairs would offer more comprehensive validation, the primary goal of this paper is to demonstrate the effectiveness of AC-PMBR decoding. Thus, we consider the current experimental scope sufficient and reserve broader extensions for future work.

Hardware-level Optimization In our experiments, we used a single GPU, but it is also possible

to use two GPUs to compute the score matrices for the target metric and the distilled metric in parallel. While further hardware-level optimizations could improve efficiency, we did not pursue them as they fall outside the core focus of this study.

Distillation Metric This study assumes the availability of a distilled metric and focuses on improving the trade-off between computational cost and translation quality under that assumption. Furthermore, the performance of AC-PMBR may vary depending on the quality of the distilled metric, as we have already demonstrated experimentally in this paper. While we use BLEURT and its existing distilled metrics as a case study, the results consistently show the effectiveness of AC-PMBR across multiple settings. Since the main focus lies in the methodological contribution, we do not explore the availability or development of better distilled metrics. Nonetheless, we expect that as research in metric distillation advances, the benefits of AC-PMBR will become even more pronounced.

Ethical Considerations

This study fully complies with the ACL Ethics Policy and addresses all required items in the Responsible Research Checklist. All resources used in this work are publicly available and appropriately licensed, with no license-related issues. The study does not involve or generate any harmful content. While AI assistants were used for minor writing support, such as rephrasing and spell-checking, all original content was manually created by the authors. Based on the above, we confirm that this work raises no ethical concerns.

References

- Julius Cheng and Andreas Vlachos. 2023. [Faster minimum Bayes risk decoding with confidence-based pruning](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12473–12480, Singapore. Association for Computational Linguistics.
- Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, and Taro Watanabe. 2024a. [mbrs: A library for minimum Bayes risk decoding](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 351–362, Miami, Florida, USA. Association for Computational Linguistics.
- Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, Taro Watanabe, Hideki Tanaka, and Masao Utiyama. 2024b. [Centroid-based efficient minimum Bayes risk decoding](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11009–11018, Bangkok, Thailand. Association for Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2020. [Is MAP decoding all you need? the inadequacy of the mode in neural machine translation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 4506–4520, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Bryan Eikema and Wilker Aziz. 2022. [Sampling-based approximations to minimum Bayes risk decoding for neural machine translation](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 10978–10993, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Man-deep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, Naman Goyal, Tom Birch, Vitaliy Liptchinsky, Sergey Edunov, Edouard Grave, Michael Auli, and Armand Joulin. 2020. [Beyond english-centric multilingual machine translation](#). *CoRR*, abs/2010.11125.
- Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023. [Epsilon sampling rocks: Investigating sampling strategies for minimum Bayes risk decoding for machine translation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9198–9209, Singapore. Association for Computational Linguistics.
- Vaibhava Goel and William J Byrne. 2000. [Minimum bayes-risk automatic speech recognition](#). *Computer Speech & Language*, 14(2):115–135.
- Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luisa Coheur, Pierre Colombo, and André F. T. Martins. 2024. [xcomet: Transparent machine translation evaluation through fine-grained error detection](#). *Transactions of the Association for Computational Linguistics*, 12:979–995.
- Yuu Jinnai and Kaito Ariu. 2024. [Hyperparameter-free approach for faster minimum Bayes risk decoding](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 8547–8566, Bangkok, Thailand. Association for Computational Linguistics.
- Juraj Juraska, Daniel Deutsch, Mara Finkelstein, and Markus Freitag. 2024. [MetricX-24: The Google submission to the WMT 2024 metrics shared task](#). In *Proceedings of the Ninth Conference on Machine Translation*, pages 492–504, Miami, Florida, USA. Association for Computational Linguistics.
- Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, Barry Haddow,

- Philipp Koehn, Benjamin Marie, Christof Monz, Makoto Morishita, Kenton Murray, Masaaki Nagata, Toshiaki Nakazawa, Martin Popel, and 3 others. 2023. [Findings of the 2023 conference on machine translation \(WMT23\): LLMs are here but not quite there yet](#). In *Proceedings of the Eighth Conference on Machine Translation*, pages 1–42, Singapore. Association for Computational Linguistics.
- Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, Rebecca Knowles, Philipp Koehn, Christof Monz, Makoto Morishita, Masaaki Nagata, Toshiaki Nakazawa, Michal Novák, Martin Popel, and Maja Popović. 2022. [Findings of the 2022 conference on machine translation \(WMT22\)](#). In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 1–45, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Philipp Koehn and Rebecca Knowles. 2017. [Six challenges for neural machine translation](#). In *Proceedings of the First Workshop on Neural Machine Translation*, pages 28–39, Vancouver. Association for Computational Linguistics.
- Shankar Kumar and William Byrne. 2004. [Minimum Bayes-risk decoding for statistical machine translation](#). In *Proceedings of the Human Language Technology Conference of the North American Chapter of the Association for Computational Linguistics: HLT-NAACL 2004*, pages 169–176, Boston, Massachusetts, USA. Association for Computational Linguistics.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Maja Popović. 2015. [chrF: character n-gram F-score for automatic MT evaluation](#). In *Proceedings of the Tenth Workshop on Statistical Machine Translation*, pages 392–395, Lisbon, Portugal. Association for Computational Linguistics.
- Matt Post. 2018. [A call for clarity in reporting BLEU scores](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 186–191, Brussels, Belgium. Association for Computational Linguistics.
- Amy Pu, Hyung Won Chung, Ankur Parikh, Sebastian Gehrmann, and Thibault Sellam. 2021. [Learning compact metrics for MT](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 751–762, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. [BLEURT: Learning robust metrics for text generation](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7881–7892, Online. Association for Computational Linguistics.
- Firas Trabelsi, David Vilar, Mara Finkelstein, and Markus Freitag. 2024. [Efficient minimum bayes risk decoding using low-rank matrix completion algorithms](#). In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Jannis Vamvas and Rico Sennrich. 2024. [Linear-time minimum Bayes risk decoding with reference aggregation](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 790–801, Bangkok, Thailand. Association for Computational Linguistics.
- Dave Zachariah, Martin Sundin, M. Jansson, and Saikat Chatterjee. 2012. [Alternating least-squares for low-rank matrix reconstruction](#). *IEEE Signal Processing Letters*, 19:231–234.

A Detailed Experimental Settings

Metrics	#parameters
BLEURT-20	579M
BLEURT-20-D12	167M
BLEURT-20-D6	45M
BLEURT-20-D3	30M

Table 2: Number of parameters of BLEURT-20 (Sellam et al., 2020) and its distilled metrics (Pu et al., 2021).

Datasets The datasets we used in our experiments and the number of sentences they contain are listed in Table 3, and we tuned hyperparameters only on WMT’22 En→De, using WMT’23 En→De and WMT’23 De→En exclusively as test sets.

Dataset	En→De	De→En
WMT’22	2,037	1,984
WMT’23	557	549

Table 3: Number of sentences for each dataset.

Computational costs The major bottleneck of AC-PMBR decoding is utility score calculation using target and distilled metrics. In the low-reduction setting of AC-PMBR decoding, utility calculation took more than 1,000 times longer than the agreement-constrained ALS algorithm. Therefore, Algorithm 1 can be ignored from the overall cost, and the utility score calculation is dominant in the computational costs of both PMBR and AC-PMBR decoding.

We used BLEURT-20 (Sellam et al., 2020) as the utility metric, alongside its distilled metrics BLEURT-20-D3,6,12 (Pu et al., 2021), which shrink the parameter count from 579M to 30M, 45M, and 167M, respectively, as listed in Table 2. Since the wall-clock computation time highly depends on the hardware environment, we evaluated the computation cost based on the time complexity. In the PMBR decoding, the time complexity is formally defined as $\mathcal{O}(\frac{NMC}{r})$, where C is the cost of utility function calls. We fixed the number of candidates and pseudo-references, i.e., N and M are constant; thus, the computational costs in our settings now depend on $\mathcal{O}(\frac{C}{r})$. Here, the cost of the utility function call C is sublinearly proportional

to the number of parameters in a metric model¹. Therefore, we defined the total computational cost of PMBR decoding $\text{Cost}_{\text{PMBR}}$ as follows:

$$\text{Cost}_{\text{PMBR}} := \frac{\text{\#parameters}}{r}. \quad (4)$$

Similarly, the cost of AC-PMBR decoding $\text{Cost}_{\text{AC-PMBR}}$ is defined as follows:

$$\text{Cost}_{\text{AC-PMBR}} := \frac{\text{\#parameters of target metric}}{r} + \frac{\text{\#parameters of distilled metric}}{r'}. \quad (5)$$

In all experiments, to compare the performance of PMBR and AC-PMBR decoding at a comparable cost, we set roughly the same costs for $\text{Cost}_{\text{PMBR}}$ and $\text{Cost}_{\text{AC-PMBR}}$, i.e., we set $\text{Cost}_{\text{PMBR}} \approx \text{Cost}_{\text{AC-PMBR}}$ in both high and low reduction rate settings.

Hyperparameter tuning In our AC-PMBR decoding, we tuned $\gamma \in \mathbb{R}_+$, which is a weight of the agreement term in Equation (3), with a fixed random seed. As shown in Figures 3 and 4, we varied $\gamma \in \{0.1, 0.2, \dots 1.0\}$ in each reduction setting, i.e., low-reduction setting and high-reduction setting, and selected $\gamma = 0.1$ and $\gamma = 1.0$, respectively. These tuning experiments revealed that the optimal agreement weight tends to increase with the reduction rate. This is because, under high-reduction settings, utility scores are sparsely observed, and the information from the distilled metric becomes more beneficial for completing the score matrix.

Computational environment All experiments were conducted on a single NVIDIA RTX A6000 GPU with an Intel® Xeon® Gold 6426Y processor, and our method was implemented with MBRS (Deguchi et al., 2024a).

¹<https://github.com/google-research/bleurt/blob/master/checkpoints.md>

B Statistical Significance Tests

Decoding	En→De		De→En	
	BLEU	chrF2	BLEU	chrF2
Reduction rate: Low	PMBR: $r = 16$, AC-PMBR: $r = 32$			
PMBR (baseline)	19.1 ± 1.5	46.2 ± 1.8	23.5 ± 1.7	50.7 ± 2.2
AC-PMBR	19.2 ± 1.5	46.6 ± 1.7	23.8 ± 1.7	51.1 ± 2.1
Reduction rate: High	PMBR: $r = 512$, AC-PMBR: $r = 1,024$			
PMBR (baseline)	16.8 ± 1.4	44.1 ± 1.8	21.3 ± 1.6	48.7 ± 2.1
AC-PMBR	18.1 ± 1.4	46.2 ± 1.6	23.0 ± 1.6	50.5 ± 2.0

Table 4: Results of statistical significance tests on the WMT’23 En↔De translation tasks comparing AC-PMBR and PMBR decoding. All scores are reported as mean \pm 95% confidence intervals. Entries with $p < 0.05$ and higher scores are highlighted in **bold**.

We conduct statistical significance tests using sacreBLEU (Post, 2018) for BLEU and chrF metrics, with each result based on 1,000 bootstrap resampling iterations drawn from the WMT’23 translation tasks. As shown in Table 4, the tests confirmed that the improvements achieved by AC-PMBR decoding over PMBR decoding were statistically significant ($p < 0.05$) in all high-reduction settings and for chrF under the low-reduction setting, supporting the robustness of the observed gains.

C Detailed Derivation

In our AC-PMBR decoding, we minimize the following loss function \mathcal{L} :

$$\mathcal{L} := \mathcal{L}_{\text{MF}}(\mathbf{U}, \mathbf{V}; \tilde{\mathbf{O}}) + \mathcal{L}_{\text{MF}}(\mathbf{U}', \mathbf{V}'; \tilde{\mathbf{O}}') + \gamma \mathcal{L}_{\text{AC}}(\mathbf{U}, \mathbf{V}, \mathbf{U}', \mathbf{V}'). \quad (6)$$

Thus, the rank reduced representations in the target utility, \mathbf{u}_i and \mathbf{v}_j , are updated as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}_i} = -2 \sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}) \\ n=i}} \left(O_{ij} - \mathbf{u}_i^\top \mathbf{v}_j \right) \mathbf{v}_j + 2\lambda \mathbf{u}_i + 2\gamma (\mathbf{u}_i - \mathbf{u}'_i) \quad (7)$$

$$\left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}) \\ n=i}} \mathbf{v}_j \mathbf{v}_j^\top + (\lambda + \gamma) \mathbf{I} \right) \mathbf{u}_i = \sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}) \\ n=i}} O_{ij} \mathbf{v}_j + \gamma \mathbf{u}'_i \quad (8)$$

$$\mathbf{u}_i = \left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}) \\ n=i}} \mathbf{v}_j \mathbf{v}_j^\top + (\lambda + \gamma) \mathbf{I} \right)^{-1} \left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}) \\ n=i}} O_{ij} \mathbf{v}_j + \gamma \mathbf{u}'_i \right). \quad (9)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}_j} = -2 \sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}) \\ m=j}} \left(O_{ij} - \mathbf{u}_i^\top \mathbf{v}_j \right) \mathbf{u}_i + 2\lambda \mathbf{v}_j + 2\gamma (\mathbf{v}_j - \mathbf{v}'_j) \quad (10)$$

$$\left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}) \\ m=j}} \mathbf{u}_i \mathbf{u}_i^\top + (\lambda + \gamma) \mathbf{I} \right) \mathbf{v}_j = \sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}) \\ m=j}} O_{ij} \mathbf{u}_i + \gamma \mathbf{v}'_j \quad (11)$$

$$\mathbf{v}_j = \left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}) \\ m=j}} \mathbf{u}_i \mathbf{u}_i^\top + (\lambda + \gamma) \mathbf{I} \right)^{-1} \left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}) \\ m=j}} O_{ij} \mathbf{u}_i + \gamma \mathbf{v}'_j \right). \quad (12)$$

Likewise, the rank reduced representations in the distilled utility, \mathbf{u}'_i and \mathbf{v}'_j , are updated as follows:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}'_i} = -2 \sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}') \\ n=i}} \left(O'_{ij} - \mathbf{u}'_i{}^\top \mathbf{v}'_j \right) \mathbf{v}'_j + 2\lambda \mathbf{u}'_i + 2\gamma (\mathbf{u}'_i - \mathbf{u}_i) \quad (13)$$

$$\left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}') \\ n=i}} \mathbf{v}'_j \mathbf{v}'_j{}^\top + (\lambda + \gamma) \mathbf{I} \right) \mathbf{u}'_i = \sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}') \\ n=i}} O'_{ij} \mathbf{v}'_j + \gamma \mathbf{u}_i, \quad (14)$$

$$\mathbf{u}'_i = \left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}') \\ n=i}} \mathbf{v}'_j \mathbf{v}'_j{}^\top + (\lambda + \gamma) \mathbf{I} \right)^{-1} \left(\sum_{\substack{(n,j) \in \text{Obs}(\tilde{\mathbf{O}}') \\ n=i}} O'_{ij} \mathbf{v}'_j + \gamma \mathbf{u}_i \right). \quad (15)$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{v}'_j} = -2 \sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}') \\ m=j}} \left(O'_{ij} - \mathbf{u}'_i{}^\top \mathbf{v}'_j \right) \mathbf{u}'_i + 2\lambda \mathbf{v}'_j + 2\gamma (\mathbf{v}'_j - \mathbf{v}_j), \quad (16)$$

$$\left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}') \\ m=j}} \mathbf{u}'_i \mathbf{u}'_i{}^\top + (\lambda + \gamma) \mathbf{I} \right) \mathbf{v}'_j = \sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}') \\ m=j}} O'_{ij} \mathbf{u}'_i + \gamma \mathbf{v}_j, \quad (17)$$

$$\mathbf{v}'_j = \left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}') \\ m=j}} \mathbf{u}'_i \mathbf{u}'_i{}^\top + (\lambda + \gamma) \mathbf{I} \right)^{-1} \left(\sum_{\substack{(i,m) \in \text{Obs}(\tilde{\mathbf{O}}') \\ m=j}} O'_{ij} \mathbf{u}'_i + \gamma \mathbf{v}_j \right). \quad (18)$$

We obtained \mathbf{U} and \mathbf{V} for the matrix completion using Algorithm 1 with these derived update rules.