Unveiling the Power of Source: Source-based Minimum Bayes Risk Decoding for Neural Machine Translation

Anonymous ACL submission

Abstract

001Maximum a posteriori decoding, a commonly002used method for neural machine translation003(NMT), aims to maximize the estimated poste-004rior probability. However, high estimated prob-005ability does not always lead to high translation006quality. Minimum Bayes Risk (MBR) decod-007ing (Kumar and Byrne, 2004) offers an alter-008native by seeking hypotheses with the highest009expected utility.

Inspired by Quality Estimation (QE) reranking which uses the QE model as a ranker (Fernandes et al., 2022), we propose source-based MBR (sMBR) decoding, a novel approach that utilizes quasi-sources (generated via paraphrasing or back-translation) as "support hypotheses" and a reference-free quality estimation metric as the utility function, marking the first work to solely use sources in MBR decoding. Experiments show that sMBR outperforms QE reranking and the standard MBR decoding. Our findings suggest that sMBR is a promising approach for NMT decoding.¹

1 Introduction

011

013

015

017

024

037

Neural Machine Translation (NMT) models typically aim to select a hypothesis with the highest estimated posterior probability during decoding, an approach known as Maximum A Posteriori (MAP) decoding. Beam search (Graves, 2012; Sutskever et al., 2014), which balances computational cost and search accuracy, has become the standard approximate decoding method for MAP.

However, the underlying assumption of beam search - that estimated probability is a good proxy for translation quality - has been challenged by evidence showing that estimated probability and quality do not always correlate positively (Ott et al., 2018; Freitag et al., 2021). For example, Fig 1 illustrates a case where a human reference translation



Figure 1: Example of $De \rightarrow En$, with source "*Kommt* einem Spitzel nahe". BS denotes beam search. The estimated log probability of a human reference is lower than that of the beam search output, and even lower than that of a bad translation.

has a lower estimated probability than the hypothesis generated by beam search, and even lower than that of a poor translation. Furthermore, the true MAP output is sometimes an empty string or overly brief translation (Koehn and Knowles, 2017; Murray and Chiang, 2018; Ott et al., 2018; Stahlberg and Byrne, 2019). These suggest that solely searching for high estimated probability hypotheses may not be an effective strategy for improving quality.

040

041

042

043

044

051

053

054

059

060

061

062

063

064

065

067

Given the limitations of using estimated probability as a proxy for quality, an attractive alternative is to directly target translation quality during decoding (Freitag et al., 2022). Minimum Bayes Risk (MBR) decoding, proposed in the era of statistical machine translation (Kumar and Byrne, 2004; Tromble et al., 2008), aims to find the hypothesis with the highest expected utility with respect to a set of hypotheses called "support hypotheses". Traditionally, surface-based evaluation metrics like BLEU (Papineni et al., 2002) were used as the utility function in MBR decoding (Eikema and Aziz, 2020; Eikema and Aziz, 2022). However, these metrics have shown limited correlation with human judgments (Mathur et al., 2020; Freitag et al., 2023b), hindering the widespread adoption of MBR decoding based on them. Recent work has explored using state-of-the-art neural metrics, such as COMET (Rei et al., 2022a), as utility functions for MBR decoding (Freitag et al., 2022; Fernan-

¹We will make the source code publicly available once the paper is accepted.

100

101

104

105

108

109

068

des et al., 2022; Freitag et al., 2023a), showing promising improvements in human evaluations.

Moreover, advances in reference-free evaluation metrics (Rei et al., 2021; Rei et al., 2022b; Rei et al., 2023) have enabled their direct application to hypothesis reranking, which we refer to as Quality Estimation (QE) reranking (Fernandes et al., 2022). QE reranking selects the hypothesis with the highest reference-free quality estimation score among the candidate hypotheses. However, QE reranking remains understudied compared to MBR decoding.

Inspired by QE reranking, which uses the QE model as a reranker, we propose a novel approach called source-based MBR (sMBR) decoding, which uses quasi-sources generated by paraphrasing or back-translation as "support hypotheses" and a QE metric as the utility function. This marks the first work to solely use sources as support hypotheses in MBR decoding, breaking the long-standing tradition of relying on using other hypotheses to approximate true utility for this purpose. See Fig 2 for a overview of our methodology. We provide empirical evidence through comprehensive experiments on three translation directions in both classic (large transformer models trained from scratch, including high-resource and low-resource sub-setups) and large language model (LLM) setups, demonstrating that sMBR outperforms QE reranking and the standard MBR decoding. These findings suggest that sMBR is a promising NMT decoding approach.

2 Decoding methods in NMT

Decoding can be viewed as two phases: *hypothesis generation* and *decision*. Specifically, in the hypothesis generation phase, a certain generation method, such as beam search, is used to generate N hypotheses from the model $\{h_0, h_1, \ldots, h_{N-1}\}$. Then, in the decision phase, N decision scores need to be computed for each of these N hypotheses $\{score_0, score_1, \ldots, score_{N-1}\}$. Finally, the hypothesis with the highest decision score is selected as the final output.

110 2.1 MAP decoding

111Given a source sentence x and a hypothesis space112 \mathcal{H} , the translation model $P_{mt}(\cdot \mid x)$ estimates the113probability of any hypothesis $h \in \mathcal{H}$. MAP decod-114ing aims to find the hypothesis y^{MAP} that maxi-

mizes the probability:

$$y^{MAP} = \operatorname*{argmax}_{h \in \mathcal{H}} P_{mt}(h \mid x). \tag{1}$$

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

147

148

149

150

151

152

153

154

155

156

157

158

160

In other words, MAP decoding simply takes the estimated probability as the decision score.

However, considering all possible hypotheses in \mathcal{H} is computationally intractable. Therefore, beam search is widely used as an efficient approximation of MAP decoding, balancing the trade-off between computational cost and search accuracy.

Increasing the beam size leads to searching for hypotheses with higher estimated probabilities. However, in practice, when the beam size exceeds 5 or 10, it often leads to a performance degradation instead (Tu et al., 2017; Koehn and Knowles, 2017). This phenomenon is known as the *beam search curse*, considered one of the six challenges of NMT (Koehn and Knowles, 2017).

2.2 MBR decoding

Unlike MAP decoding, which aims to find the highest estimated probability hypothesis, Minimum Bayes Risk (MBR) decoding seeks the hypothesis that minimizes the expected loss (or equivalently, maximizes the expected utility) with respect to a set of hypotheses, called "support hypotheses".

In practice, it is common to use a set of hypotheses from a model as support hypotheses. Formally, let $S \subseteq \mathcal{H}$ be a set of support hypotheses from model $P_{mt}(\cdot \mid x)$, for support hypothesis $h_s \in S$, MBR decoding selects the hypothesis y^{MBR} that has the least risk:

$$y^{MBR} = \operatorname*{argmin}_{h \in \mathcal{H}} \mathbb{E}_{h_s \in \mathcal{S}} \left[L(h_s, h) \mid x \right]$$
(2)

$$= \underset{h \in \mathcal{H}}{\operatorname{argmin}} \sum_{h_s \in \mathcal{S}} L(h_s, h) P_{mt}(h_s \mid x). \quad (3)$$

In practice, it is common to use a utility function, correlated to human evaluation results, such as BLEU or COMET, as an alternative to the loss function $L(\cdot, \cdot)$. Thus, the purpose of MBR decoding is actually to select the hypothesis with the maximum expected utility. In addition, the hypothesis space is usually too large to traverse all the hypotheses to find a translation that satisfies the above conditions. Therefore, a set of hypotheses C from the hypothesis space \mathcal{H} , called "candidate hypotheses", is often used as a representative of the whole hypothesis space. Combining these two points, for a given utility function $u(\cdot, \cdot)$, the MBR decoding objective can be reformulated as:



Figure 2: Overview of decoding methods in NMT. The diagram illustrates the process for MBR decoding, QE reranking, and the proposed sMBR decoding. It also shows two practices of sMBR: sMBR-BT and sMBR-PP. The figure demonstrates how the score used for selecting the final hypothesis is computed for each method.

$$y^{MBR} \approx \operatorname*{argmax}_{h \in \mathcal{C}} \sum_{h_s \in \mathcal{S}} u(h_s, h) P_{mt}(h_s \mid x).$$
 (4)

A widely used practice is to use the same set of hypotheses for both C and S, and to assume that all h_s have the same probability, instead of the estimated probability given by the model. That is, the expected utility for a chosen hypothesis is approximated by averaging its utilities to other hypotheses. Hypotheses can be obtained through, for example, beam search or ancestral sampling². Formally, the objective of MBR can be approximated as:

$$y^{MBR} \approx \operatorname*{argmax}_{h \in \mathcal{C}} score_h^{MBR},$$
 (5)

$$score_h^{MBR} = \frac{1}{|\mathcal{S}|} \sum_{h_s \in \mathcal{S}} u(h_s, h).$$
 (6)

2.3 QE reranking

161

162

163

166

168

170

171

172

173

174

175

176

177

179

180

181

182

184

Quality estimation (QE) is a task that aims to assess the quality of a translated sentence without reference translations but the original source sentence x. Recently, QE models have been employed to develop a new decoding method called QE reranking (Fernandes et al., 2022), which leverages QE models to rerank the candidate hypotheses. The main idea behind QE reranking is to select the hypothesis h with the highest estimated quality with the QE model, rather than relying on the estimated probability $P_{mt}(h|x)$. Formally, for a source x, a candidate hypothesis from the hypothesis space $h \in C$, and QE function f_{QE} , QE reranking aims at finding a y^{QE} that has the highest QE decision score $score_h^{QE}$:

$$y^{QE} = \operatorname*{argmax}_{h \in \mathcal{C}} score_{h}^{QE}$$
(7)

$$= \operatorname*{argmax}_{h \in \mathcal{C}} f_{QE}(x, h). \tag{8}$$

185

187

188

189

190

191

192

194

195

196

197

198

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

3 Method

In this section, we first introduce our proposed method, source-based MBR (sMBR) decoding. Then, we show that QE reranking is actually a special case of sMBR. Finally, we introduce two practices of sMBR: paraphrasing-based (sMBR-PP) and back-translation-based (sMBR-BT). See Figure 2 for a quick overview of how sMBR differs from standard MBR decoding and QE reranking.

3.1 sMBR

In this subsection, we introduce MBR decoding using a novel method to calculate the utility, sMBR decoding. Then, we will show that QE reranking is actually a special case of sMBR.

We hypothesize that for a hypothesis, decision scores calculated by QE using only the hypothesis and sources are a good proxy for translation quality, since QE reranking can achieve promising performance using only a reference-free reranker (QE model) (Fernandes et al., 2022). We are therefore interested in calculating the utility of MBR decoding using only sources and a QE model, and call MBR decoding based on this idea **sMBR decoding**. In other words, in the sMBR, sources are used as

²In the context of NMT, ancestral sampling refers to sampling from the entire vocabulary without any pruning.

217

218

221

228

229

239

240

241

242

244

245

246

247

248

249

250

259

"support hypotheses", and the QE model is used as a utility function.

Given that empirically better performance can be obtained by using more support hypotheses in the standard MBR decoding (Freitag et al., 2022; Fernandes et al., 2022; Freitag et al., 2023a), we would like to use more "support hypotheses" (sources) for sMBR as well. However, in the standard MBR decoding, there are usually multiple support hypotheses (i.e., |S| > 1), while we usually only have one source sentence. Therefore, we obtain additional "support hypotheses" for sMBR by considering other source language sentences that have the same meaning as the original source sentence.

Formally, let $P_{pp}(X'|X)$ be a paraphrasing distribution of source language sentences with the same meaning as the original source x, and $\widetilde{X'}$ be a finite sample from $P_{pp}(X'|X = x)$. Then, for a reference-free QE utility function $u(\cdot, \cdot)$, sMBR looks for y^{sMBR} that has the highest sMBR decision score $score_h^{sMBR}$ in the set of candidate hypotheses C:

$$y^{sMBR} = \operatorname*{argmax}_{h \in \mathcal{C}} \mathbb{E}_{x' \in \widetilde{X'}} \left[u(x', h) \mid x \right]$$
(9)

$$= \operatorname*{argmax}_{h \in \mathcal{C}} \sum_{x' \in \widetilde{X'}} u(x', h) P_{pp}(x'|x).$$
(10)

Similarly to eq. 6, it is also possible to approximate by assuming that all sources x' have the same probability, which simplifies the calculation. Formally, let $K = |\widetilde{X'}|$, then the objective of sMBR can be approximated as:

$$y^{sMBR} \approx \operatorname*{argmax}_{h \in \mathcal{C}} score_h^{sMBR},$$
 (11)

$$score_{h}^{sMBR} = \frac{1}{K} \sum_{x' \in \widetilde{X}'} u(x', h).$$
(12)

Here, QE reranking is a special case of sMBR when K = 1, i.e., when only the original source is used. Unlike QE reranking, sMBR considers multiple quasi-sources (K > 1), which are intended to serve a more diverse and representative utility.

3.2 sMBR-PP and sMBR-BT

We refer to the exact formulation presented in section 3.1 as sMBR-PP. In addition, we study an alternative approach, which indirectly generates quasi-sources by back-translation, sMBR-BT. Specifically, for an original source x, we first generat a translation h_0 using the forward translation model $P_{mt}(h|x)$ and then use h_0 as the input to a back-translation model to generate K quasi-sources $\{x'^1, x'^2, \dots, x'^k\}$. We then use the260set $\widetilde{X'} = \{x, x'^1, x'^2, \dots, x'^k\}$ of size K + 1 for261sMBR decoding. Note that in sMBR-BT, the input262to the back-translation model is a single hypothesis, where we simply use the one with the highest264estimated probability. Then, we obtain K quasisources by beam search with beam size K.260

267

268

270

271

272

273

274

275

276

277

278

279

280

282

284

285

286

288

289

290

291

292

293

294

295

296

297

298

299

301

302

303

304

305

4 **Experiments**

4.1 Setup

In this subsection, we present the details of NMT systems, decoding methods, and evaluation.

4.1.1 Data and models for NMT

We consider two setups: the classic (encoderdecoder based transformer model trained from scratch on parallel corpora) and the LLM setup.

Our experiments were performed on the English to German (En \rightarrow De), English to Russian (En \rightarrow Ru), and Chinese to English (Zh \rightarrow En), as COMET and COMET-QE on them proved to be highly correlated with human judgments at the segment level (Rei et al., 2022a). We use generaltest2023 (Kocmi et al., 2023) as the test set for each translation direction.

Classic setup The classic setup does not include the Zh \rightarrow En due to computing resource limitations. For both En \rightarrow De and En \rightarrow Ru in the classic setup, we use newstest2017-2019 as the development set.

We further consider both high-resource and lowresource sub-setups. For the high-resource setup, we use Facebook FAIR's WMT19 news translation task submission (Ng et al., 2019). For low resource setup, training data consists of 0.44M and 0.38M parallel sentences for En \rightarrow De and En \rightarrow Ru, respectively. These systems allow us to assess the performance in a data-constrained scenario. See Appendix D for more details.

LLM setup We use TowerInstruct-13B (Alves et al., 2024), a state-of-the-art LLM for translation-related tasks as NMT model. We prompt LLM to perform zero-shot translation. See Appendix E for more details.

4.1.2 Decoding

We employ four approaches for hypothesis generation: beam search, ancestral sampling, top-k sampling, and epsilon sampling. Beam search has

373

374

375

376

377

378

379

380

381

382

383

384

385

386

387

388

389

390

391

392

394

395

396

397

398

352

353

been widely used for MBR decoding in the past (Kumar and Byrne, 2004; Stahlberg et al., 2017;
Shu and Nakayama, 2017; Blain et al., 2017), while
ancestral sampling has gained popularity in recent
work (Eikema and Aziz, 2020; Freitag et al., 2022;
Eikema and Aziz, 2022). We include top-k sampling and epsilon sampling since we find that they
yield better performance than ancestral sampling.

315

316

317

320

321

324

325

330

350

351

Both for the classic and LLM setups, we evaluate the following decision rules:

- MAP: A widely used rule that selects the hypothesis with the highest estimated probability.
- MBR: MBR decoding based on COMET, using the Unbabel/wmt22-comet-da³ model. It calculates decision scores for candidate hypotheses using the approach described in 2.2, and then select the highest score one.
- QE reranking: A special form of MBR, which calculates decision scores with the quality estimation model $f_{QE}(\cdot, \cdot)$. Since COMET does not support reference-free quality estimation, we use COMET-QE (Unbabel/wmt22-cometkiwi-da⁴) as the utility function.
- sMBR: Source-based MBR decoding. We evaluate its two practices: sMBR-PP and sMBR-BT. The same QE model used for 333 QE reranking is employed as the utility func-334 tion. For sMBR-PP, we fine-tuned a unique 335 T5-large (Chung et al., 2024) model as a paraphrase generator for each source lan-337 guage, which was fine-tuned on paraphrase 338 data created through back-translation. We detail the sMBR-PP implementation in Ap-340 pendix A. As the back-translation model of 341 sMBR-BT, in the classic setup, we again use 342 Facebook FAIR's WMT19 news translation task submission in the high resource setup; in the low resource setup, we use the same model architecture and data as in the forwardtranslation model. For the LLM setup, we 347 reused the LLM (TowerInstruct-13B) as a 348 back-translation model.

Our baseline is a MAP based on beam search with a beam size of 5 both for the classic and LLM

setups, since we found that larger beams do not lead to better performance. Except for the baseline, we use **400** candidate hypotheses for the classic setup unless otherwise specified, as we find that more hypotheses have limited gains in performance but result in higher costs.⁵ For the LLM setup, we only use **128** hypotheses because generating more hypotheses leads to too much computing.

For MBR decoding, we tried two setups: (1) using the same set of hypotheses as candidate hypotheses and support hypotheses; and (2) using QE reranking to filter the set of support hypotheses to a smaller size that matches the size of the set of support hypotheses for sMBR. For sMBR-PP and sMBR-BT, we study case using 16 quasi-sources, as adding more did not yield further gains.

More details are provided in Appendix B.

4.1.3 Evaluation metrics

We use automatic evaluation metrics, including BLEU (Post, 2018), **XCOMET** (XCOMET-XXL) (Guerreiro et al., 2023), and MetricX (MetricX-24-Hybrid-XXL) (Juraska et al., 2024) to evaluate different methods. Our choice of XCOMET and MetricX is motivated by two key factors: 1) they are state-of-the-art neural metrics (Freitag et al., 2024), which correlate with human judgments even when evaluating NMT systems that use neural metric-based reranking (Kovacs et al., 2024); 2) they were trained on Multidimensional Ouality Metrics (MOM) data (Guerreiro et al., 2023; Juraska et al., 2024), while the COMET series, on which MBR and sMBR were directly optimized, were trained only on Direct Assessments (DA) data (Rei et al., 2022a). Given the limited correlation between MQM and DA (Freitag et al., 2021), we expect XCOMET and MetricX to provide a more independent assessment, as they are less likely to be biased towards the COMET-optimized MBR and sMBR methods. We perform significance tests using paired bootstrap resampling (Koehn, 2004).

4.2 Results

In this subsection, we analyze the performance by observing the results of the automatic evaluation metrics. Due to space constraints, we show in the main text only the results of beam search based results in the classic setup and epsilon sampling

³huggingface.co/Unbabel/wmt22-comet-da

⁴huggingface.co/Unbabel/wmt22-cometkiwi-da

⁵Appendix G shows the impact of candidate hypotheses number on the metrics.

Decoding method			$En \rightarrow De$ $En \rightarrow Ru$								
w/ beam search	$ \mathcal{C} $	$ \mathcal{S} $	BLEU ↑	XCOMET ↑	MetricX↓	BLEU ↑	XCOMET ↑	MetricX↓			
			High Resource (55.4M and 52.0M training data)								
MAP	5	_	34.80	84.89	3.63	27.82	82.90	5.36			
MBR	5	5	34.97	85.12	3.58	27.35	83.62	5.01			
MAP	$\bar{400}$		34.30	85.10	3.69	23.05	82.99	6.49			
QE reranking	400	1	35.23	86.48	3.22	22.80	86.20	4.27			
MBR	400	400	34.83	85.74	3.50	22.81	84.95	5.17			
MBR	400	17	34.93	85.88	3.34	23.05	85.00	5.37			
sMBR-PP	400	17	34.81	86.73 [†]	3.09 ^{††}	22.91	86.52 ^{††}	4.14 [†]			
sMBR-BT	400	17	33.80	86.17	3.33	22.36	84.99	4.65			
				Low Resour	ce (0.44M ar	nd 0.38M t	raining data)				
MAP	5	_	11.44	60.29	12.43	17.44	65.85	10.98			
MBR	5	5	12.56	60.33	12.16	17.35	66.95	10.51			
MAP	400	_	9.78	62.58	12.24	17.51	66.36	10.85			
QE reranking	400	1	13.63	65.63	10.34	17.71	74.66	7.81			
MBR	400	400	12.66	63.79	11.05	17.56	70.69	9.01			
MBR	400	17	11.75	64.53	11.11	18.17 ^{††}	71.10	9.07			
sMBR-PP	400	17	13.49	66.36 ^{††}	10.19 [†]	17.95	74.96 ^{††}	7 . 76 †			
sMBR-BT	400	17	9.66	63.77	11.68	16.87	69.69	9.35			

Table 1: Compares decision rules in the classic setup. |C| and |S| indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17 support hypotheses (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; † and †† indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

based results in the LLM setup; other results are included in the Appendix C.

400

Classic setup Table 1 highlights the effectiveness 401 of sMBR decoding in the classic setup with beam 402 search. Regarding XCOMET and MetricX, sMBR-403 PP significantly outperforms QE reranking, prov-404 ing the validity of our extension to QE reranking. 405 The results of the experiments based on sampling 406 407 methods are shown in Appendix C, where similar gains to those based on beam search can be ob-408 served. The sMBR-PP outperforms the standard 409 MBR on the neural metric, regardless of whether 410 the standard MBR can use the full 400 support 411 hypotheses or only 17. Thus, we conclude that 412 the sMBR-PP outperforms QE reranking and the 413 standard MBR in the classic setup. 414

415LLM setupTable 2 shows the results of epsilon416sampling in the LLM setup. We observe that417sMBR-PP can still significantly outperform QE418reranking regarding XCOMET and MetricX. The419XCOMET and MetricX of sMBR-PP are compara-420ble to the standard MBR and sometimes outperform421standard MBR. As with the results of the classic

setup, a gain similar to that based on epsilon sampling can be observed in Appendix C. Thus, we conclude that sMBR-PP still significantly outperforms QE reranking and is competitive with the standard MBR in the LLM setup. 422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

The performance of sMBR-PP relative to standard MBR differs between the two setups. The standard MBR shows performance similar to that of the sMBR-PP in the LLM setup. We believe this is due to the better quality of the support hypothesis generated by LLM, which leads to a higher approximation accuracy in eq. 6.

Compared to QE reranking, sMBR-BT shows gains regarding XCOMET and MetricX in the LLM setup, but even lower metrics than QE reranking in the classic setup. We investigate this reason in section I.

5 Discussion

In this section, we discuss the effectiveness of
sMBR-PP and the mechanism behind it through440some experiments. Due to space constraints, we
discuss the efficiency of sMBR-PP in Appendix F,
where we show that sMBR-PP is much slower than441

Decoding method				En→De		Zh→En		
	$ \mathcal{C} $	$ \mathcal{S} $	BLEU↑	XCOMET ↑	MetricX↓	BLEU↑	XCOMET ↑	MetricX↓
w/ epsilon sampling				Τον	verInstruct-1	13B		
MAP	5	_	30.24	86.06	3.32	20.73	88.15	2.41
MBR	5	5	28.10	87.30	2.95	19.74	88.96	2.20
MAP	128		32.64	86.43	3.22	23.12	89.14	2.20
QE reranking	128	1	29.40	88.76	2.56	19.88	90.64	1.89
MBR	128	128	29.84	89.19 [†]	$2.46^{\dagger\dagger}$	22.01**	90.39	1.90
MBR	128	17	31.93††	88.83	2.60	23.34††	90.43	1.87
sMBR-PP	128	17	27.19	89.47 ^{††}	2.44 ^{††}	19.87	90.70 †	1.87
sMBR-BT	128	17	28.73	89.04	2.50	19.77	90.38	1.98

Table 2: Compares decision rules in the LLM setup. The meaning of table elements is the same as Table 1.

$ \mathcal{S} $	1	6	11	17	33		Self-BLEU↓	Semantic Similarity↑
$\textbf{XCOMET} \uparrow$	86.48	86.65	86.73	86.73	86.74	sMBR-PP	41.68	94.32
MetricX \downarrow	3.22	3.12	3.10	3.09	3.09	sMBR-BT	48.25	94.53

Table 3: Impact of increasing sources for sMBR-PP: The number of sources is positively correlated with the evaluation metrics. |S| = 1 + K, i.e., (# of the original source) + (# of quasi-sources). Candidate hypotheses were generated by beam search. When |S| = 1, sMBR-PP is QE reranking.

QE reranking and MBR with simple optimization.

5.1 Effects of increasing sources

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

Since sMBR is an extension to QE reranking by increasing the number of sources, we first investigate the impact of increasing the number of sources on the performance of sMBR-PP. We focus on the En \rightarrow De high resource setup of the classic setup and evaluate with neural metrics. Table 3 presents the results, demonstrating a positive correlation between source number and evaluation metrics. This observation again shows that our extensions to QE reranking are effective. In addition, the increase in synthesis sources from 16 to 32 does not result in further gains, which we hypothesize is due to the inability of the paraphrase generator to achieve the generation of up to 32 generative high-quality synthesis sources.

5.2 Qualities of quasi-sources

463To understand the properties of the quasi-sources464of sMBR, we analyzed them in terms of surface465diversity and semantic similarity with the original466source. Surface diversity was measured using Self-467BLEU (Zhu et al., 2018), while semantic similarity468was assessed by cosine distance between sentence

Table 4: Analyzing of quasi-sources: analyzed on the $En \rightarrow De$ generaltest2023, high resource. Lower Self-BLEU means richer surface diversity; higher semantic similarity means closer semantics to the original source.

$ \mathcal{S} $	1	6	11	17
Ave. QE	81.28	80.58	80.57	80.54

Table 5: Average QE scores with the original source: QE scores are negatively correlated with source number. The analysis was performed in the En \rightarrow De high resource setup on sMBR-PP. Candidate hypotheses were generated by beam search.

embeddings⁶.

The results presented in Table 4 reveal that the quasi-sources generated by sMBR-PP exhibit much lower Self-BLEU scores compared to those produced by sMBR-BT, indicating greater surface diversity. On the other hand, the scores of the two in terms of semantic similarity are close, implying that both generated quasi-sources do not deviate too much from the original source's semantics. We hypothesize that the poor performance of sMBR-BT in the classic setup can be attributed to the limited surface diversity of its quasi-sources.

5.3 Why sMBR-PP works

The results shown in Section 4.2 demonstrate that sMBR-PP is significantly better than QE rerank-

477 478 479

469

470

471

472

473

474

475

- 480 481
- 482 483

⁶We use huggingface.co/sentence-transformers/ all-mpnet-base-v2

ing in terms of neural metrics. We discuss the 484 mechanism behind producing these gains. We hy-485 pothesize that the QE model used in QE rerank-486 ing is overly sensitive to the specific phrasing and 487 structure of the original source, leading to an over-488 reliance on a single source that could negatively 489 impact performance. In contrast, aggregating OE 490 scores across multiple sources in sMBR decoding 491 is expected to provide more robust QE. We ob-492 served that as the number of sources increases in 493 sMBR decoding, the average QE score between the 494 selected translations and the original source sen-495 tence decreases (Table 5). This suggests that sMBR 496 decoding no longer relies solely on QE with respect 497 to the original source. We conjecture that this is 498 because sMBR decoding tends to select hypothe-499 ses that are more generally applicable to different source variants.

6 Related Work

502

505

506

507

510

511

512

513

514

515

517

518

521

523

525

527

529

533

MBR decoding has been used in speech recognition (Stolcke et al., 1997; Goel and Byrne, 2000), word alignment (Kumar and Byrne, 2002), and statistical machine translation (Kumar and Byrne, 2004; Tromble et al., 2008). Recently, some works have re-explored the application of MBR decoding in NMT and demonstrated promising results (Stahlberg et al., 2017; Shu and Nakayama, 2017; Eikema and Aziz, 2020; Eikema and Aziz, 2022). These works have shown that MBR decoding can help overcome some of the limitations of MAP.

In past work, MBR decoding is usually based on beam search to generate candidate hypotheses (Stahlberg et al., 2017; Shu and Nakayama, 2017). Recently, Eikema and Aziz (2020) proposed sampling-based MBR decoding and found that the samples from the model are faithful to the training data statistics, while the beam search is not. Freitag et al. (2022) further explored the impact of the generation method on the performance.

In terms of utility functions, past work has primarily used surface-based metrics such as BLEU and BEER (Stanojević and Sima'an, 2014). However, these metrics have limited correlation with human judgments (Mathur et al., 2020; Freitag et al., 2023b). Recently, a trend has been to combine advanced neural metrics with MBR decoding, such as COMET and BLEURT. These works demonstrate that neural metrics-based MBR can improve performance in human evaluations (Freitag et al., 2022; Fernandes et al., 2022; Freitag et al., 2023a). However, they are also limited by the high cost, as MBR decoding has a secondary cost for the number of candidate hypotheses. Eikema and Aziz (2022) investigated decoupling candidate and support hypotheses, enabling the exploration of more potential candidate hypotheses within a limited computational budget. In addition, some recent work has focused on improving the efficiency of the MBR decoding (Cheng and Vlachos, 2023; Vamvas and Sennrich, 2024; Deguchi et al., 2024).

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

578

579

580

On the other hand, some work has found that models used to assess the quality of NMT systems (i.e., quality estimation) can perform well even in the absence of a reference (Rei et al., 2021; Rei et al., 2022b; Rei et al., 2023). Fernandes et al. (2022) explored the direct use of quality estimation models as rerankers for NMT and showed promising results, referred to as QE reranking.

7 Conclusions and Future Work

In this work, inspired by QE reranking, we propose sMBR decoding, which uses sources and QE model to calculate the utility, the first practical method to solely rely on sources as "support hypotheses" in MBR decoding. Experimental results (Section 4.2) show that sMBR decoding outperforms QE reranking and the standard MBR decoding.

Despite its limitations, such as the challenge of generating quasi-sources, sMBR represents a significant step forward in MBR decoding. By breaking with the tradition of approximating true utility using only the average of utilities with respect to other hypotheses, sMBR opens up new possibilities for future research.

Our analysis in Appendix I indicates that using a more powerful paraphrase generator, such as GPT4 (Achiam et al., 2023), for sMBR-PP shows promise for further performance improvements. The analysis in Appendix I suggests that using Diverse Beam Search (Vijayakumar et al., 2017) for sMBR-BT has the potential to enhance performance. Therefore, we plan to explore these methods for generating quasi-sources in our future work. In addition to methods for generating quasi-sources, in our future work we will continue to investigate broadening the boundaries of "support hypotheses" to include sentences in languages other than the source and target.

583

586

587

591

593

594

595

596

597

606

610

612

613

615

616

617

618

619

621

622

625 626

627

628

630

8 Limitations

While our proposed sMBR decoding approach shows promising results, it has some limitations.

Firstly, reranking methods that directly optimize evaluation metrics may "overfit" to those metrics, causing the optimized metrics to become unreliable (Kovacs et al., 2024). We mitigate this problem by using automatic metrics that are likely to have low correlation with the metrics that are directly optimized. However, since they still have common parts in the training data, this measure may not completely avoid the problem of unreliable metrics. Moreover, extending QE reranking with a quasisource may exacerbate this problem, as it may result in the final selection being favored by the QE metric itself but translations that are unrelated to the source (referred to as "universal translations" by Yan et al. (2023)). Therefore, our conclusion that sMBR-PP outperforms QE reranking and standard MBR decoding may be questioned. Human evaluation can mitigate this issue but is costly and time-consuming.

Secondly, generating high-quality quasi-source sentences remains a challenge. We explored two methods based on paraphrasing and backtranslation, but the back-translation approach did not consistently improve reranking performance. This suggests that further research is needed to identify more effective techniques for generating diverse, representative quasi-sources.

Finally, we have only tested the proposed method in a limited number of translation directions and domains. However, not all language pairs have well-performing quality estimation models available. In the case of some language pairs, this may lead to a questioning of one of our basic hypothesis, i.e., the quality estimation model is a good proxy for the true utility. Therefore, the effectiveness of sMBR in a wider range of settings remains an open question.

Acknowledgements

References

- OpenAI Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mo Bavarian, Jeff Belgum, Irwan Bello, et al. 2023. Gpt-4 technical report.
- Duarte M. Alves, José P. Pombal, Nuno M. Guerreiro,

Pedro H. Martins, Joao Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, Pierre Colombo, Jos'e G. C. de Souza, and André Martins. 2024. Tower: An open multilingual large language model for translation-related tasks. *ArXiv*, abs/2402.17733. 631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

- Frédéric Blain, Lucia Specia, and Pranava Madhyastha. 2017. Exploring hypotheses spaces in neural machine translation. In *Proceedings of Machine Translation Summit XVI: Research Track*, pages 282–298, Nagoya Japan.
- Julius Cheng and Andreas Vlachos. 2023. Faster minimum Bayes risk decoding with confidence-based pruning. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12473–12480, Singapore. Association for Computational Linguistics.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. 2024. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53.
- Hiroyuki Deguchi, Yusuke Sakai, Hidetaka Kamigaito, Taro Watanabe, Hideki Tanaka, and Masao Utiyama. 2024. Centroid-based efficient minimum Bayes risk decoding. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 11009–11018, Bangkok, Thailand and virtual meeting. 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.
- Patrick Fernandes, António Farinhas, Ricardo Rei, José G. C. de Souza, Perez Ogayo, Graham Neubig, and Andre Martins. 2022. Quality-aware decoding for neural machine translation. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1396–1412, Seattle, United States. Association for Computational Linguistics.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.

- 6 6 6 6
- 6 6
- 69
- 69
- 700 701
- 703 704 705 706 707
- 708 709 710 711

713

- 714 715 716 717 718 719 720
- 720 721 722 723 723 724
- 725 726 727 728
- 730
- 732 733
- 734 735
- 736

737 738 739

- 740 741
- 742
- 743 744

- Markus Freitag, Behrooz Ghorbani, and Patrick Fernandes. 2023a. 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.
- Markus Freitag, David Grangier, Qijun Tan, and Bowen Liang. 2022. High quality rather than high model probability: Minimum Bayes risk decoding with neural metrics. *Transactions of the Association for Computational Linguistics*, 10:811–825.
- Markus Freitag, Nitika Mathur, Daniel Deutsch, Chi-Kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Frederic Blain, Tom Kocmi, Jiayi Wang, David Ifeoluwa Adelani, Marianna Buchicchio, Chrysoula Zerva, and Alon Lavie. 2024. Are LLMs breaking MT metrics? results of the WMT24 metrics shared task. In *Proceedings of the Ninth Conference on Machine Translation*, pages 47–81, Miami, Florida, USA. Association for Computational Linguistics.
- Markus Freitag, Nitika Mathur, Chi-kiu Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frederic Blain, Daniel Deutsch, Craig Stewart, Chrysoula Zerva, Sheila Castilho, Alon Lavie, and George Foster. 2023b. Results of WMT23 metrics shared task: Metrics might be guilty but references are not innocent. In *Proceedings of the Eighth Conference on Machine Translation*, pages 578–628, Singapore. Association for Computational Linguistics.
- Vaibhava Goel and William J Byrne. 2000. Minimum bayes-risk automatic speech recognition. *Computer Speech & Language*, 2(14):115–135.
- Alex Graves. 2012. Sequence transduction with recurrent neural networks. *CoRR*, abs/1211.3711.
- Nuno M. Guerreiro, Ricardo Rei, Daan van Stigt, Luísa Coheur, Pierre Colombo, and André Martins. 2023. xCOMET: Transparent machine translation evaluation through fine-grained error detection. *ArXiv*, abs/2310.10482.
- John Hewitt, Christopher Manning, and Percy Liang. 2022. Truncation sampling as language model desmoothing. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 3414– 3427, Abu Dhabi, United Arab Emirates. 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.
- Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In 3rd International Conference on Learning Representations,

ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings. 745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

768

769

770

771

772

773

774

775

776

777

778

779

781

782

783

784

787

791

793

794

795

796

797

798

- 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, Makoto Nagata, Toshiaki Nakazawa, Martin Popel, Maja Popović, and Mariya Shmatova. 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.
- Philipp Koehn. 2004. Statistical significance tests for machine translation evaluation. In Proceedings of the 2004 Conference on Empirical Methods in Natural Language Processing, pages 388–395, Barcelona, Spain. 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.
- Geza Kovacs, Daniel Deutsch, and Markus Freitag. 2024. Mitigating metric bias in minimum Bayes risk decoding. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1063–1094, Miami, Florida, USA. Association for Computational Linguistics.
- Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Shankar Kumar and Bill Byrne. 2004. Minimum bayesrisk 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.
- Shankar Kumar and William Byrne. 2002. Minimum Bayes-risk word alignments of bilingual texts. In *Proceedings of the 2002 Conference on Empirical Methods in Natural Language Processing (EMNLP 2002)*, pages 140–147. Association for Computational Linguistics.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net.

- 810
- 811
- 813
- 816 817

- 823 824 825
- 826

- 829
- 830
- 835
- 837 838

836

839

- 847

851

855

- Nitika Mathur, Timothy Baldwin, and Trevor Cohn. 2020. Tangled up in BLEU: Reevaluating the evaluation of automatic machine translation evaluation metrics. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 4984–4997, Online. Association for Computational Linguistics.
 - Kenton Murray and David Chiang. 2018. Correcting length bias in neural machine translation. In Proceedings of the Third Conference on Machine Translation: Research Papers, pages 212-223, Brussels, Belgium. Association for Computational Linguistics.
 - Nathan Ng, Kyra Yee, Alexei Baevski, Myle Ott, Michael Auli, and Sergey Edunov. 2019. Facebook FAIR's WMT19 news translation task submission. In Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1), pages 314-319, Florence, Italy. Association for Computational Linguistics.
 - Myle Ott, Michael Auli, David Grangier, and Marc'Aurelio Ranzato. 2018. Analyzing uncertainty in neural machine translation. In International Conference on Machine Learning.
 - Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In Proceedings of NAACL-HLT 2019: Demonstrations.
 - 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.
 - 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, Belgium, Brussels. Association for Computational Linguistics.
 - Ricardo Rei, José G. C. de Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André F. T. Martins. 2022a. COMET-22: Unbabel-IST 2022 submission for the metrics shared task. In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 578–585, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Ricardo Rei, Ana C Farinha, Chrysoula Zerva, Daan van Stigt, Craig Stewart, Pedro Ramos, Taisiya Glushkova, André F. T. Martins, and Alon Lavie. 2021. Are references really needed? unbabel-IST 2021 submission for the metrics shared task. In Proceedings of the Sixth Conference on Machine Translation, pages 1030-1040, Online. Association for Computational Linguistics.

Ricardo Rei, Nuno M. Guerreiro, José Pombal, Daan van Stigt, Marcos Treviso, Luisa Coheur, José G. C. de Souza, and André Martins. 2023. Scaling up CometKiwi: Unbabel-IST 2023 submission for the quality estimation shared task. In *Proceedings of the* Eighth Conference on Machine Translation, pages 841-848, Singapore. Association for Computational Linguistics.

856

857

858

859

860

861

862

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

- Ricardo Rei, Marcos Treviso, Nuno M. Guerreiro, Chrysoula Zerva, Ana C Farinha, Christine Maroti, José G. C. de Souza, Taisiya Glushkova, Duarte Alves, Luisa Coheur, Alon Lavie, and André F. T. Martins. 2022b. CometKiwi: IST-unbabel 2022 submission for the quality estimation shared task. In Proceedings of the Seventh Conference on Machine Translation (WMT), pages 634-645, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Neural machine translation of rare words with subword units. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Raphael Shu and Hideki Nakavama. 2017. Later-stage minimum bayes-risk decoding for neural machine translation. CoRR, abs/1704.03169.
- Felix Stahlberg and Bill Byrne. 2019. On NMT search errors and model errors: Cat got your tongue? In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pages 3356-3362, Hong Kong, China. Association for Computational Linguistics.
- Felix Stahlberg, Adrià de Gispert, Eva Hasler, and Bill Byrne. 2017. Neural machine translation by minimising the Bayes-risk with respect to syntactic translation lattices. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers, pages 362-368, Valencia, Spain. Association for Computational Linguistics.
- Miloš Stanojević and Khalil Sima'an. 2014. Fitting sentence level translation evaluation with many dense features. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pages 202-206, Doha, Qatar. Association for Computational Linguistics.
- Andreas Stolcke, Yochai Konig, and Mitch Weintraub. 1997. Explicit word error minimization in n-best list rescoring. In EUROSPEECH.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. Advances in neural information processing systems, 27.

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

992

993

994

995

996

997

998

999

1000

1002

1003

1004

1005

1006

1007

1008

1009

1010

1011

Roy Tromble, Shankar Kumar, Franz Och, and Wolfgang Macherey. 2008. Lattice Minimum Bayes-Risk decoding for statistical machine translation. In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 620–629, Honolulu, Hawaii. Association for Computational Linguistics.

913

914

915

916

917

918

919

922 923

924

925

926

927

928

929

932

933

934

937

939

941

942

943

945

947

949

950

951

953

957

958

959

960

961

962

963

964

- Zhaopeng Tu, Yang Liu, Lifeng Shang, Xiaohua Liu, and Hang Li. 2017. Neural machine translation with reconstruction. In *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, pages 3097–3103. AAAI Press.
- 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.
- Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Neural Information Processing Systems*.
- Ashwin K Vijayakumar, Michael Cogswell, Ramprasaath R. Selvaraju, Qing Sun, Stefan Lee, David Crandall, and Dhruv Batra. 2017. Diverse beam search: Decoding diverse solutions from neural sequence models.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. In North American Chapter of the Association for Computational Linguistics.
- Yiming Yan, Tao Wang, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Mingxuan Wang. 2023. BLEURT has universal translations: An analysis of automatic metrics by minimum risk training. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 5428–5443, Toronto, Canada. Association for Computational Linguistics.
- Yuan Zhang, Jason Baldridge, and Luheng He. 2019. PAWS: Paraphrase Adversaries from Word Scrambling. In *Proc. of NAACL*.
- Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, SIGIR '18, page 1097–1100, New York, NY, USA. Association for Computing Machinery.

A Additional details of sMBR-PP

For sMBR-PP, we generally use the same paraphrase generator for both the classic and LLM setups for the same source language.

When the source is English, for the paraphrase generator used in sMBR-PP, its specific model is google/flan-t5-large⁷. This model is trained for instruction following and thus works out-of-thebox for paraphrase generation. However, we found its performance to be rather poor, thus we chose to fine-tune it.

The fine-tuning training data consists of a publicly available paraphrase generation dataset, PAWS (Zhang et al., 2019), concatenated with a dataset we created. The dataset we created is based on En-De's News-Commentary parallel corpus⁸ and uses machine translation to create paraphrased sentences. Specifically, we first input German sentences from the parallel corpus into the $De \rightarrow En NMT$ model, and then paired its output with English sentences from the original parallel corpus to compose the samples in the paraphrase generation dataset. We use the De→En model from Facebook FAIR's WMT19 news translation task submission (Ng et al., 2019). We use a semantic similarity-based approach to estimate the quality of the dataset we created, and then filter out sentence pairs with low similarity. We use the sentence-transformers/all-mpnet-base-v2⁹ model to compute the similarity between paraphrase pairs and filter out sentence pairs with a similarity of 0.88 or less. In the end, the training data for the model consisted of a total of about 339.2K paraphrased sentence pairs, of which about 317.4K came from the data we created and about 21.8K came from PAWS.

For the training of this model, we used the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of 3e-4, weight decay of 0.0, and a batch size of 1536 examples, trained with fp32 full precision (This is because we found that the flan-t5 series is prone to training failure at fp16 precision). We set the maximum number of training epochs to 10. We randomly separate 3K sentence pairs from the dataset as the development set, and then select the checkpoints with the lowest loss on the development set.

⁷huggingface.co/google/flan-t5-large

⁸data.statmt.org/news-commentary/v18.1/

⁹huggingface.co/sentence-transformers/

all-mpnet-base-v2

When the source is Chinese, we follow a simi-1012 lar procedure as above. The difference is that the 1013 base model is mT5-large (Xue et al., 2020) and the 1014 training data only includes the dataset created by us-1015 ing TowerInstruct-13B to perform reverse trans-1016 lation on the News-Commentary parallel corpus. 1017 We use lier007/xiaobu-embedding-v2¹⁰ to cal-1018 culate the cosine similarity between the rephrase 1019 and the original sentence, and filter out samples 1020 with a similarity below 0.925. 1021

1022

1023

1024

1025

1026

1028

1029

1030

1031

1032

1033

1035

1036

1037

1038

1039

1040

1041

1042

1044

1045

1046

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

In the inference phase of generating paraphrases, we used epsilon sampling (Hewitt et al., 2022) (epsilon = 0.02), as we found that this setup balances the diversity and quality of the synthesized sources well. Training and inference were done on a single NVIDIA H100.

We used the following prompts during training and inference:

Source (input to the encoder):

paraphrase: {source_text}

Target (input to the decoder):

{target_text}

B Additional details on decoding and training of low-resource NMT models

We completed training of the NMT low-resource model, and all decoding experiments on a machine with 4 NVIDIA RTX A6000. In the hypothesis generation phase, we used CTranslate2¹¹ to generate hypotheses because of its efficiency.

For the training of low-resource NMT models, we use the fairseq (Ott et al., 2019) tool. We use base size transformer (Vaswani et al., 2017) architectures with a dropout rate of 0.3. And train for a maximum of 100 epochs at full fp16 precision. we select the checkpoints with the highest BLEU on the development set. We use adam (Kingma and Ba, 2015) optimizer with an initial learning rate of 1e-3, weight decay of 1e-4, a warm-up step of 4000, and batch size is 1e5 tokens. We build vocabulary of size 32000 with Byte-Pair Encoding (Sennrich et al., 2016) using the sentencepiece (Kudo and Richardson, 2018) tool. The vocabulary is shared between source and target languages.

C Additional experimental results

In this section, we present additional experimental results that could not be included in the main text

xiaobu-embedding-v2

due to space constraints. For the classic setup, this section includes the results of the generation methods other than beam search; for the LLM setup, this section includes not only the results of the generation methods other than Epsilon sampling, but also the complete experimental results of En->Ru.

In addition to beam search and epsilon sampling, we attempted to use top-k sampling and ancestral sampling to generate hypotheses. Unlike top-k sampling, ancestral sampling the entire vocabulary for each time step in autoregressive decoding without any pruning. The results of the experiments are shown in Table 6, 8, 7, 9, and 10. Compared to other generation methods, ancestral sampling performs poorly on both surface-based and neural metrics. Among the sampling-based methods, epsilon sampling performs best, which is consistent with the findings of Freitag et al., 2023a.

We used k = 10 for top-k sampling and epsilon = 0.02 for epsilon sampling (Freitag et al., 2023a). Due to implementation issues with some CUDA programming, we do not consider epsilon sampling with low resource setup.

In conclusion, similar to the experimental results based on beam search and epsilon sampling, the significantly boosted neural metrics demonstrate that sMBR-PP significantly outperforms QE reranking. However, improvements in neural metrics do not always lead to gains in surface-based metrics and even lead to deterioration compared to the baseline, especially when using sampling-based hypothesis generation. One possible explanation is that sampling leads to more diverse hypotheses, making it easier to generate candidates hypotheses that would lead to higher neural metrics but not favored by BLEU. Unfortunately, sMBR decoding does not consistently mitigate this issue compared to QE reranking, suggesting potential limitations in the utility functions.

D Details of the classic setup

For the high-resource setup, training data consists of 27.7M and 26.0M parallel sentences for $En \rightarrow De$ and $En \rightarrow Ru$, respectively, augmented with an equal amount of back-translation sentences. We use a single model without ensembling or language model reranking to focus on the impact of the proposed methods. For the low resource setup, we train two base Transformer models using the News-Commentary dataset¹². 1059 1060 1061

1058

1062 1063 1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1091

1092

1093

1094

1095

1097

1098

1099

1100

1101

1102

1103

1104

1105

¹⁰https://huggingface.co/lier007/

¹¹https://github.com/OpenNMT/CTranslate2

¹²data.statmt.org/news-commentary/v18.1/

Decoding method			$En \rightarrow De$ $En \rightarrow Ru$					
w/ top-k sampling	$ \mathcal{C} $	$ \mathcal{S} $	BLEU↑	XCOMET ↑	MetricX↓	BLEU ↑	XCOMET ↑	MetricX↓
				High Resour	ce (55.4M ar	nd 52.0M t	raining data)	
MAP	5	_	24.36	70.80	7.81	17.98	66.34	8.69
MBR	5	5	23.19	71.56	7.46	16.26	68.18	7.78
MAP	$\bar{400}$		24.42	69.55	8.11	21.31**	77.83	7.99
QE reranking	400	1	26.69	81.99	4.46	17.94	84.34	4.38
MBR	400	400	25.13	78.51	5.57	19.26††	80.86	5.18
MBR	400	17	27.99 ^{††}	80.28	4.89	22.21 ^{††}	81.55	5.70
sMBR-PP	400	17	25.74	82.14	4.37 [†]	17.30	84.77 ^{††}	4.20 ^{††}
sMBR-BT	400	17	25.01	77.29	6.00	16.75	79.54	5.80
				Low Resour	ce (0.44M ar	nd 0.38M t	raining data)	
MAP	5	_	7.91	53.11	14.39	14.36	61.89	12.32
MBR	5	5	8.48	53.62	14.09	13.76	62.77	11.68
MAP	$\bar{400}$		6.05	58.64	13.48	16.93**	65.55	11.12
QE reranking	400	1	10.76	62.47	11.01	15.06	73.50	7.94
MBR	400	400	10.02	60.80	11.99	15.48	69.61	9.10
MBR	400	17	10.13	61.44	11.74	17 .36 ††	69.95	9.32
sMBR-PP	400	17	10.36	63.35 ^{††}	10.90 ^{††}	15.11	73.71 [†]	7.88 [†]
sMBR-BT	400	17	5.64	59.96	12.86	14.70	68.22	9.78

Table 6: Compares sMBR with other decision rules for En \rightarrow De and En \rightarrow Ru in the **classic setup**. |C| and |S| indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17 "support hypotheses" (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; † and †† indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

1110

1111

1112

1113

1114

1115

1116

1117

1118

1120

1121

1122

1123

1124

1125

1126

1127

1107

E Details of the LLM setup

We use the target language to prompt the model to perform zero-shot translation. We used the following prompts during inference:

 $En \rightarrow De:$

Übersetzen Sie den folgenden Text vom Englischen ins Deutsche.\n Englischen:\n{source_text}\nDeutsche:

 $En \rightarrow Ru$:

Переведите следующий текст с английского на русский.\n Английский:\n{source text}\nРусский:

1119 Zh→En:

Translate the following text from
 Chinese to English.\nChinese:
 \n{source_text}\nEnglish:

We completed all decoding on a server with four NVIDIA H100s with bfloat16 precision. For sMBR-PP on Zh \rightarrow En, we trained a mT5-based (Xue et al., 2020) model for the paraphrase generator. See Appendix A for details.

F Efficiency comparison

In this section, we discuss the efficiency of sMBR-PP and compare it with QE reranking and MBR. For the latter, we compare the average decision time required for MBR, an optimized implementation of MBR (MBR-fast), and sMBR-PP to translate a single sentence. We ran each method five times on a single NVIDIA H100 with batch size 256 examples and then report the means. 1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

As expected, the decision time required for sMBR-PP to translate a sentence is much larger than that of QE reranking. Specifically, the decision time of sMBR-PP consists of two parts: the time required to generate the quasi-source and the time required to calculate the quality-aware utility function. In fact, we find that the time required to generate the quasi-source is only a small part of the overall decision time, which is about 0.13 seconds for each sentence, while the large number of quality-aware utility functions requires 3.56 seconds. In contrast, the decision time for QE reranking is 0.21 seconds per sentence, which is much faster than sMBR-PP.

Obviously, compared to QE reranking, sMBR-

Decoding method				$En \rightarrow De$ $En \rightarrow Ru$						
w/ ancestral sampling	$ \mathcal{C} $	$ \mathcal{S} $	BLEU↑	XCOMET ↑	MetricX↓	BLEU ↑	XCOMET ↑	MetricX↓		
		High Resource (55.4M and 52.0M training data)								
MAP	5	_	10.00	35.78	18.68	14.05	61.37	12.22		
MBR	5	5	6.80	30.50	20.09	12.40	58.33	12.60		
MAP	400		11.60	51.08	14.71	21.83	79.61	7.06		
QE reranking	400	1	16.98	61.07	11.98	19.12	82.12	5.12		
MBR	400	400	10.48	45.41	15.88	17.64	78.47	6.05		
MBR	400	17	18.23 ^{††}	58.37	12.49	21.44	81.47	5.29		
sMBR-PP	400	17	17.09	61.29	11.97	18.55	82.22	5.11		
sMBR-BT	400	17	14.73	56.94	12.95	19.40	80.42	5.74		
				Low Resour	ce (0.44M ar	nd 0.38M t	raining data)			
MAP	5	_	5.29	39.33	18.35	11.13	49.98	16.13		
MBR	5	5	5.85	36.64	18.92	8.47	46.38	16.73		
MAP	400		3.72	51.78	15.35	15.23 ^{††}	62.82	12.00		
QE reranking	400	1	6.67	54.18	14.14	14.22	68.13	9.96		
MBR	400	400	7.67 ^{††}	47.81	15.70	12.47	62.45	11.44		
MBR	400	17	6.99	53.16	14.49	15.13††	66.23	10.47		
sMBR-PP	400	17	6.67	54.85 ^{††}	14.11	14.08	68.41 [†]	9.88 †		
sMBR-BT	400	17	3.62	51.60	15.19	13.90	64.30	11.33		

Table 7: Compares sMBR with other decision rules for En \rightarrow De and En \rightarrow Ru in the classic setup. $|\mathcal{C}|$ and $|\mathcal{S}|$ indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17"support hypotheses" (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; † and †† indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

Decoding method				En→De		En→Ru		
w/ epsilon sampling	$ \mathcal{C} $	$ \mathcal{S} $	BLEU↑	XCOMET ↑	MetricX↓	BLEU↑	XCOMET ↑	MetricX↓
				High Resour	ce (55.4M ar	nd 52.0M t	raining data)	
MAP	5	_	24.41	81.26	5.18	19.34	73.56	7.57
MBR	5	5	25.85	82.44	4.19	16.88	74.14	7.17
MAP	400	_	13.62	76.83	7.97	22.01	77.65	6.58
QE reranking	400	1	28.64	86.12	3.12	18.94	83.12	4.57
MBR	400	400	28.23	86.20	3.17	19.48	80.09	5.43
MBR	400	17	26.72	85.18	3.84	23.05 ^{††}	81.75	5.10
sMBR-PP	400	17	27.55	86.47 [†]	3.00 [†]	18.89	83.41 [†]	4.47 [†]
sMBR-BT	400	17	12.94	78.26	7.53	17.36	79.69	5.81

Table 8: Compares sMBR with other decision rules for $En \rightarrow De$ and $En \rightarrow Ru$ in the classic setup with epsilon sampling. $|\mathcal{C}|$ and $|\mathcal{S}|$ indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17 "support hypotheses" (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; \dagger and \dagger indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

PP uses the quality perception utility function times (number of quasi-sources +1).

Next, we compare the decision time of sMBR-PP with that of the standard MBR. In MBR, the 1155 COMET model can be decomposed into a sen-1156 tence encoder f_{emb} for computing sentence em-1157

1152

1153

1154

beddings, and a simple estimator $f_{est}(\cdot, \cdot, \cdot)$ based 1158 on a multilayer perceptron. For a source x, a candi-1159 date hypothesis h, and support hypotheses $h_s \in S$, 1160 COMET-based MBR first computes the source em-1161 bedding x^{emb} , the candidate hypothesis embedding 1162 h^{emb} , and a set of support hypotheses embeddings 1163

Decoding method				En→De		Zh→En		
	$ \mathcal{C} $	$ \mathcal{S} $	BLEU↑	XCOMET [↑]	MetricX↓	BLEU ↑	XCOMET [↑]	MetricX↓
w/ beam search				To	werInstruct-1	3B		
MAP	5	_	39.75	87.11	2.47	24.94	89.13	2.13
MBR	5	5	39.84	87.29	2.46	25.00	89.14	2.13
MAP	128		40.04	87.07	2.44	24.96	89.12	2.14
QE reranking	128	1	40.11	87.35	2.40	24.96	89.22	2.11
MBR	128	128	40.07	87.29	2.41	25.00	89.17	2.13
MBR	128	17	40.14	87.19	2.41	24.96	89.13	2.14
sMBR-PP	128	17	40.15	87.45 [†]	2.36	24.93	89.28 [†]	2.10
sMBR-BT	128	17	40.18	87.39	2.37	24.90	89.21	2.11
w/ top-k sampling								
MAP	5	_	28.12	84.85	3.31	19.94	88.14	2.42
MBR	5	5	25.52	86.12	3.12	18.52	88.82	2.31
MAP	128		32.13	86.17	3.30	21.90	88.56	$\bar{2}.\bar{2}2$
QE reranking	128	1	27.56	88.09	2.73	19.07	90.10	1.93
MBR	128	128	28.28^{\dagger}	$88.68^{\dagger \dagger}$	2.56^{+}	21.19 ^{††}	90.23	1.94
MBR	128	17	30.66††	$88.76^{\dagger \dagger}$	2.59	22.40 ^{††}	90.47 ^{††}	1.87
sMBR-PP	128	17	25.20	89.05 ^{††}	2.47 ^{††}	18.85	90.19	1.91
sMBR-BT	128	17	27.13	88.48	2.65	18.87	89.82	2.03
w/ ancestral sampling								
MAP	5	_	26.29	83.20	4.00	19.07	87.21	2.62
MBR	5	5	24.20	84.67	3.62	17.97	87.79	2.52
MAP	128		30.81**	85.76	3.42	21.86	88.57	2.32
QE reranking	128	1	27.10	87.57	2.89	18.67	89.69	2.07
MBR	128	128	26.88	87.59	2.75	20.65**	89.65	2.03
MBR	128	17	$28.64^{\dagger\dagger}$	87.50	2.78	$21.72^{\dagger\dagger}$	90.02 [†]	1.97
sMBR-PP	128	17	25.78	87.91 [†]	2.86	18.29	89.70	2.08
sMBR-BT	128	17	26.78	87.79	2.84	18.40	89.41	2.13
w/ epsilon sampling								
MAP	5	_	30.24	86.06	3.32	20.73	88.15	2.41
MBR	5	5	28.10	87.30	2.95	19.74	88.96	2.20
MAP	128		32.64**	86.43	3.22	23.12 ^{††}	89.14	$-\bar{2}.\bar{2}0$
QE reranking	128	1	29.40	88.76	2.56	19.88	90.64	1.89
MBR	128	128	29.84	89.19 [†]	$2.46^{\dagger\dagger}$	22.01**	90.39	1.90
MBR	128	17	31.93††	88.83	2.60	23.34**	90.43	1.87
sMBR-PP	128	17	27.19	89.47 ^{††}	2.44 ^{††}	19.87	90.70 [†]	1.87
sMBR-BT	128	17	28.73	89.04	2.50	19.77	90.38	1.98

Table 9: Compares sMBR with other decision rules for En \rightarrow De and Zh \rightarrow En in the LLM setup. |C| and |S|indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17"support hypotheses" (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; † and †† indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

1164 1165

1166

,

$$S^{emb}$$
, using f_{emb} . Then, the MBR score of h
score_h^{MBR} can then be computed as:

$$score_{h}^{MBR} = \frac{1}{|\mathcal{S}|} \sum_{\substack{h_{s}^{emb} \in \mathcal{S}^{emb}}} f_{est}(x^{emb}, h_{s}^{emb}, h^{emb})$$

When S = C, the cost of computing utility for all candidate hypotheses in a naive MBR implementa-1168 tion is $\mathcal{O}(|\mathcal{C}|^2)$, implying a quadratic cost for both 1169 f_{emb} and $f_{est}(\cdot, \cdot, \cdot)$. 1170

However, MBR-fast optimizes embedding com-1171 putation by recognizing that the embedding any 1172

(13)

Decoding method				En→Ru	
	$ \mathcal{C} $	$ \mathcal{S} $	BLEU ↑	XCOMET [↑]	MetricX↓
w/ beam search		To	werInstruc	t-13B	
MAP	5	_	29.46	89.51	3.00
MBR	5	5	29.33	89.61	2.96
MAP	128		$2\bar{9}.\bar{5}0^{\dagger}$	89.58	3.00
QE reranking	128	1	29.33	89.65	2.97
MBR	128	128	29.38	89.69	2.97
MBR	128	17	29.52 ^{††}	89.57	3.00
sMBR-PP	128	17	29.36	89.77 [†]	2.96
sMBR-BT	128	17	29.39	89.69	2.96
w/ top-k sampling					
MAP	5	_	23.22	86.79	3.69
MBR	5	5	21.23	87.02	3.53
MAP	128		28.07**	89.30	2.97
QE reranking	128	1	22.02	91.37	2.45
MBR	128	128	24.24 ^{††}	91.26	2.48
MBR	128	17	27.47 ^{††}	91.36	2.41
sMBR-PP	128	17	20.24	91.49 ^{††}	2.40
sMBR-BT	128	17	21.27	90.80	2.51
w/ ancestral sampling					
MAP	5	_	23.28	86.79	3.69
MBR	5	5	21.12	87.18	3.51
MAP	128		28.32**	89.27	3.01
QE reranking	128	1	21.35	90.96	2.55
MBR	128	128	23.77 ^{††}	90.98	2.53
MBR	128	17	26.55 ^{††}	91.23 ^{††}	2.47
sMBR-PP	128	17	20.52	90.99	2.50
sMBR-BT	128	17	21.63	90.78	2.57
w/ epsilon sampling					
MAP	5	_	26.19	88.53	3.20
MBR	5	5	23.92	89.56	2.87
MAP	128		29.61	89.68	2.97
QE reranking	128	1	22.35	91.67	2.34
MBR	128	128	$25.90^{\dagger\dagger}$	91.33	2.35
MBR	128	17	$28.94^{\dagger\dagger}$	91.54	2.40
sMBR-PP	128	17	21.26	92.12 ^{††}	2.20 ^{††}
sMBR-BT	128	17	22.47	91.68	2.38

Table 10: Compares sMBR with other decision rules for $En \rightarrow Ru$, the NMT model is **TowerInstruct-13B**. |C| and |S| indicate the number of candidate hypotheses and supportive hypotheses, respectively. For sMBR, we used |S| = 17 support hypotheses (1 original source + 16 quasi-sources). We performed paired bootstrap resampling; † and †† indicate significantly better than QE reranking within groups (p < 0.05 and p < 0.01, respectively; Multiple testing correction is not applied). The best in each group is marked in bold.

1205

1206

1207

1208

1210

1173sentence in a triple (x, h, h_s) is independent of the1174other elements. By pre-computing sentence embed-1175dings independently for all sources and hypotheses,1176MBR-fast avoids duplicate f_{emb} computations and1177reduces its cost to $\mathcal{O}(|\mathcal{C}|)$ when $\mathcal{S} = \mathcal{C}$. The esti-1178mator $f_{est}(\cdot, \cdot, \cdot)$ still has a quadratic cost $\mathcal{O}(|\mathcal{C}|^2)$ 1179since the order of elements within the triple affects1180the output, but it is computationally cheaper com-1181pared to a f_{emb} consisting of multiple transformer1182blocks. Note that this optimization is not univer-1183sal due to the fact that it takes advantage of the1184particular architecture of COMET.

In contrast, the COMET-QE model used in sMBR consists of an encoder f_{emb}^{QE} that takes the concatenated source and hypothesis as input and outputs their joint embedding, and an estimator f_{est}^{QE} . The joint embeddings must be computed separately for each source-hypothesis pair, resulting in a cost of $\mathcal{O}(K \times |\mathcal{C}|)$ for both f_{emb}^{QE} and f_{est}^{QE} .

Table 11 shows the measurement results. sMBR is faster than the naive implementation of MBR because it uses a smaller number of support hypotheses. However, it is much slower than the optimized implementation of MBR due to the difficulty of further optimizing its utility function itself.

In summary, while sMBR-PP significantly improves translation quality compared to QE reranking and has competitive performance to MBR, there is still room for improving its efficiency to match optimized COMET-based MBR decoding.

	MBR	MBR-fast	sMBR-PP
Decision time	135.29 s	0.32 s	3.56 (+ 0.13) s

Table 11: Decision time for translating a sentence: measured on newstest2020 in $En \rightarrow De$. For sMBR-PP, the number in parentheses is the quasi-source generation time. The batch size is 256.

G Impact of the number of candidate hypotheses

We explored the impact of the number of candidate hypotheses on the evaluation metrics in an $En \rightarrow De$ high resource setting. Figure 3 shows the results. We find that 400 is an appropriate number, as more candidate hypotheses bring small performance gains and lead to higher costs.

H Using LLM to generate quasi-sources for sMBR-PP

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1237

1238

1239

1240

1241

1242

We show in the results of investigating the nature of quasi-sources in sMBR-PP and sMBR-BT based on Self-BLEU and semantic similarity. These results suggest that quasi-sources with richer surface forms and greater semantic similarity to the original source may lead to better translation quality. In fact, during the early stages of this research, we tried using GPT4-0125 (Achiam et al., 2023), the state-of-the-art LLM at the time, to generate quasi-sources for sMBR-PP.

	XCOMET ↑	MetricX↓
sMBR-PP (T5)	86.52	4.14
sMBR-PP (GPT4)	87.10	4.09

Table 12: Comparison of sMBR-PP performance based on different paraphrase generators: Experiments conducted on high-resource sub-setup, $En \rightarrow Ru$ language pair, using beam search to generate candidate hypotheses.

Table 12 shows our results on the classic setup, En \rightarrow Ru language pair, high-resource sub-setup, using beam search to generate candidate hypotheses.

We found that sMBR-PP based on GPT4-0125 achieved better performance on both XCOMET and MetricX.

	$Self\text{-}BLEU{\downarrow}$	Semantic Similarity \uparrow
sMBR-PP (T5)	45.95	92.83
sMBR-PP (GPT4)	18.67	93.46

Table 13: Analyzing of quasi-sources: analyzed on the $En \rightarrow Ru$ generaltest2023, high resource. Lower Self-BLEU means richer surface diversity; higher semantic similarity means closer semantics to the original source.

We investigated the properties of quasi-sources generated by GPT4-0125 using the same method as in , and the results are presented in Table 13.

We observed that quasi-sources generated by GPT4-0125 had lower Self-BLEU than those generated by T5, while maintaining similar semantic similarity. This indicates that quasi-sources generated by GPT4-0125 have richer surface forms.

Therefore, we believe that using paraphrases with richer surface forms as quasi-sources can indeed improve the performance of sMBR-PP. However, considering that research based on non-open models like GPT4-0125 would make our work difficult to reproduce—after all, GPT4-0125 produces



Figure 3: Impact of the number of candidate hypotheses on the evaluation metrics: in the $En \rightarrow De$ high resource setup. The horizontal axis indicates the number of candidate hypotheses and the vertical axis indicates the evaluation indicators.

different outputs for identical inputs even with a specified random seed or temperature of 0, and GPT4-0125 may not be accessible in the future. On the other hand, although sMBR-PP based on T5 doesn't perform as well on XCOMET and MetricX as sMBR-PP based on GPT4-0125, T5 is an open model that allows us to ensure the reproducibility of our research and support our conclusions. Therefore, we chose T5-based sMBR-PP as our main experimental setup. Nevertheless, recent research on open LLMs has made significant progress, and we will explore using open LLMs as alternative experimental setups for sMBR-PP in future work.

1243

1244

1245

1246

1247

1248 1249

1250

1251

1253

1254

1255

1257

1258

1259

1261

1263

1265

I Generating quasi-sources for sMBR-PP with diverse beam search

The analysis results in suggest that the poor performance of sMBR-BT may be due to the lack of diversity in the surface form of the quasi-source. This could be because we used simple beam search to generate the quasi-source.

Diverse Beam Search (DBS) (Vijayakumar et al., 2017) is an improved beam search method that can generate more diverse text. We tried using different

beam search methods to generate quasi-sources for sMBR-BT in the classic setup.

	XCOMET ↑	MetricX↓
QE ranking	86.48	3.22
sMBR-BT (BS)	86.17	3.33
sMBR-BT (DBS)	86.19	3.29

Table 14: Comparison of sMBR-BT performance based on different search methods: Experiments conducted on high-resource sub-setup, $En \rightarrow De$ language pair, using beam search to generate candidate hypotheses.

Table 14 shows the experimental results for the $En \rightarrow De$ task (classic setting, high-resource subset, using beam search to generate candidate hypotheses), where sMBR-BT (BS) indicates using simple beam search to generate the quasi-source, while sMBR-BT (DBS) indicates using DBS. We found that using DBS to generate quasi-sources only slightly improved the performance of sMBR-BT, but it was still worse than QE reranking.

We will attempt to use more improved methods (such as sampling-based generation methods) to generate quasi-sources to enhance the performance

1268

1276 1277 1278

1279

of sMBR-BT as future work.