Has Machine Translation Evaluation Achieved Human Parity? The Human Reference and the Limits of Progress

Anonymous ACL submission

Abstract

In Machine Translation (MT) Evaluation, metric performance is assessed based on agreement with human judgments. In recent years, automatic metrics have demonstrated increasingly high levels of agreement with humans. To gain a clearer understanding of metric performance and establish an upper bound, we incorporate human baselines in the MT Meta-Evaluation, that is, the assessment of MT metrics capabilities. Our results show that human annotators are not consistently superior to automatic metrics, with state-of-the-art metrics often ranking on par with or higher than human baselines. Despite these findings suggesting human parity, we discuss several reasons for caution. Finally, we explore the broader implications of our re-017 sults for the research field, asking: Can we still 018 reliably measure improvements in MT Evalua-019 tion? With this work, we aim to shed light on the limits of our ability to measure progress in the field, fostering discussion on an issue that we believe is crucial to the entire MT Evaluation community.

011

024

033

037

041

1 **Introduction and Related Work**

Machine Translation (MT) Evaluation is the task of assessing the quality of the translated text, while MT Meta-Evaluation estimates the capabilities of automatic evaluation techniques, i.e., MT metrics. Historically, automatic metrics have been employed for their low cost and fast experimentation time, whereas human evaluation is still considered the gold standard, necessary for validating automatically derived findings. However, in recent years the MT Evaluation field has seen significant advancements. Neural-based metrics have demonstrated strong correlations with human judgments, largely replacing traditional heuristic-based metrics, and becoming the de facto standard in MT evaluation (Freitag et al., 2022, 2023, 2024). More recently, LLM-based approaches to MT Evaluation have emerged (Kocmi and Federmann, 2023b,a;

Fernandes et al., 2023; Bavaresco et al., 2024), offering not only high correlation with human judgments but also improved interpretability. Thus, we ask: What is missing for automatic techniques to achieve human parity – if they have not already? Indeed, unlike other Natural Language Processing (NLP) tasks, MT Evaluation lacks a human performance reference, making it difficult to gauge the true capabilities of MT metrics. For instance, in MT, human performance is measured by evaluating human references alongside system translations (Läubli et al., 2018; Kocmi et al., 2023, 2024a). Similarly, popular NLP benchmarks such as HellaSwag (Zellers et al., 2019), MMLU (Hendrycks et al., 2021), and MT-bench (Zheng et al., 2023) report the performance of human baselines.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

078

079

081

082

Since in MT Evaluation metric performance is measured based on agreement with human annotators, we posit that the agreement between different annotators can serve as a reference for human performance. Previous studies reported the Inter-Annotator Agreement (IAA) in MT Evaluation: Lommel et al. (2014b) used Cohen's kappa to measure the pairwise agreement between raters; Freitag et al. (2021a) grouped raters' assessments into seven score bins before computing pairwise agreement; and Kocmi et al. (2024b) used Kendall τ correlation coefficient. However, they have used different measures, making direct comparisons difficult, and none of them contextualized IAA in relation to the performance of automatic metrics. To the best of our knowledge, Perrella et al. (2024a) were the first to assess metric and human performance jointly. However, since comparing humans and metrics was not their primary focus, they used a single human annotation protocol that exhibited very poor performance - likely due to low annotation quality - rendering it ineffective as a human performance reference for MT metrics.

In this work, we address this gap by incorporating human baselines into the metric rankings from

	20	20	20	22	2023	2024
	\rightarrow DE	$ $ ZH \rightarrow $ $ \rightarrow DE $ $ \rightarrow ZH $ $		\rightarrow DE	\rightarrow es	
MQM	3	3	3	3	2	1
ESA	X	X	X	X	2	1
SQM	3	3	X	X	×	X
DA+SQM	×	×	1	1	1	×
#Seg	681	895	583	1065	156	449
#Sys	9	9	10	13	12	12

Table 1: The four top rows indicate the number of available and distinct annotations for each annotation protocol and test set. We list the studies that released these annotations in Appendix A. '2020' refers to the data released by Freitag et al. (2021a), while other years correspond to the test sets from the respective WMT editions. The notation $\rightarrow x$ indicates that the test set contains translations from English to X, whereas $x \rightarrow$ denotes translations from X to English. The two bottom rows contain the number of source segment present in the intersection of annotations from each protocol, for each test set.

various editions of the Metrics Shared Task of the Conference on Machine Translation (WMT). By using Meta-Evaluation strategies from WMT 2024 we derive a single, comprehensive ranking of MT *evaluators* – both human and automatic – establishing a human performance reference for MT metrics across several test sets, translation directions, and human annotation protocols, and offering a clearer understanding of the capabilities of current MT evaluation techniques. Then, given that our results suggest that automatic metrics may have reached human parity, we critically examine this claim and discuss its implications for future research in MT Evaluation.

084

095

100

101

102

2 Preliminaries and Experimental Setup

In this section, we describe the human annotations, the annotation protocols, the test sets selected for our work, the Meta-Evaluation strategies employed, and the automatic metrics included.

2.1 The Human Annotations

103Each year WMT conducts new manual annotation104campaigns to collect human judgments about trans-105lation quality. First, several test sets are created by106selecting N_t source segments per test set t, from107various sources. Source segments can be single108sentences or entire paragraphs. These segments are109then translated into the target language using M_t 110MT systems. As a result, each test set t contains

 $N_t \times M_t$ translations. Finally, human raters are hired to assess the quality of the collected translations (Kocmi et al., 2023; Freitag et al., 2023; Kocmi et al., 2024a; Freitag et al., 2024). In order to annotate such a large volume of translations, non-overlapping portions of the data are typically assigned to multiple raters. Consequently, an annotated test set is formed by combining annotations from multiple raters. Unless explicitly stated otherwise, the annotations used in this work follow this approach. In this context, we use the term evaluator to refer to any entity that performs the MT Evaluation task. An evaluator can be an MT metric, a human rater, an ensemble of MT Metrics, or an entity that selects annotations from different raters depending on the source segment.

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

2.2 Test Sets and Annotation Protocols

To estimate a human performance reference in MT Evaluation we require multiple human annotations about translation quality for the same test set. This way, we can designate one annotation as ground truth and the others as human baselines. Consequently, we use the test sets released by Freitag et al. (2021a) and at WMT editions from 2022 to 2024, which contain human annotations from at least two of the following protocols: Multidimensional Quality Metrics (Lommel et al., 2014a, MQM), Error Span Annotation (Kocmi et al., 2024b, ESA), Professional Scalar Quality Metrics (Freitag et al., 2021a, pSQM), and Direct Assessments + Scalar Quality Metrics (Kocmi et al., 2022a, DA+SQM). Table 1 summarizes the test sets and language directions used in our work, along with the availability of human annotations. We describe the aforementioned annotation protocols in Appendix A.

Following standard practice in the literature (Freitag et al., 2021a,b, 2022, 2023, 2024), we designate the MQM annotations released annually at WMT as the ground truth, and employ the others as human baselines. Indeed, the MQM protocol relies on experienced annotators, providing a more finegrained (and more expensive) evaluation compared to other protocols.

2.3 The MT Meta-Evaluation

We compute metric rankings using the Meta-Evaluation strategies employed at the WMT 2024 Metrics Shared Task:

 Soft Pairwise Accuracy (SPA) estimates evaluator performance based on ability to rank MT

Test set 2020		1	en→de		1	ZH→EN		Test set 2023	EN→		→DE			
		SPA		acc_{ea}^{*}		SP	A	$\operatorname{acc}_{ea}^{*}$			SI	PA	ac	c_{eq}^*
Metric	Rank	Acc	. Ra	nk 7	Acc.	Rank	Acc.	Rank	Acc.	Metric	Rank	Acc.	Rank	Acc.
MOM-2020-2	1	96 4	15	1 5	68.86	1	88 10	1	55 70	GEMBA-MQM*	1	96.07	3	57.50
nSOM-1	1	95 !	i9	6 4	19 41	1	79.16	13	43.89	CometKiwi-XXL*	1	95.51	3	57.56
MOM 2020 3			20	0 5	6 84	1	02.06	20	52.80	MetricX-23-QE-XXL*	1	94.82	1	61.44
NIQNI-2020-5	4			4 5	0.04	1	52.00	2	52.80	DA+SQM	2	93.00	13	46.41
BLEURI-0.2	2	86.8	51	4 5	0.81	2	72.59	3	50.57	ESA-1	2	92.21	13	46.49
pSQM-2	2	2 85.8	37	9 4	6.97	1	89.33	9	46.77	MQM-2023-2	3	89.45	14	43.08
BLEURT-20	2	2 85.5	52	3 - 5	51.68	3	67.46	4	50.12	ESA-2	3	89.14	11	49.77
Test set 2022			EN-	→DE			EN	→ZH		Test set 2024		EN-	→ES	
		SI	PA	a	cc_{eq}^*		SPA	a	cc_{eq}^*		SF	PA	acc_{eq}^{*}	
Metric		Rank	Acc.	Rank	Acc.	Rank	Acc.	Rank	Acc.	Metric	Rank	Acc.	Rank	Acc.
MetricX-23-QE-X	XL*	1	94.89	3	57.64	2	83.92	2	47.43	CometKiwi-XXL*	1	86.12	4	67.24
MQM-2022-2		1	94.49	6	55.55	2	80.82	3	47.05	gemba_esa*	1	85.72	3	67.68
MQM-2022-3		1	92.59	1	61.06	1	87.22	2	47.56	ESA	2	80.12	8	63.84
MetricX-23-XXL		2	92.34	2	59.27	1	87.69	1	48.43	metametrics_mt_mqm	2	80.10	1	68.95
DA+SQM		6	66.61	16	46.03	2	82.95	12	36.26	MetricX-24-Hybrid	2	79.75	1	69.20

Table 2: Results obtained by applying the WMT 2024 Meta-Evaluation strategies to the test sets illustrated in Section 2.2. The 'Acc.' column contains the Meta-Evaluation accuracy, while 'Rank' reports clusters of statistical significance computed following Freitag et al. (2024), using the PERM-BOTH hypothesis test introduced by Deutsch et al. (2021). Starred metrics are reference-less metrics, and rows highlighted in gray are human evaluators.

systems in the same order as in the ranking derived from ground truth annotations (Thompson et al., 2024).

• Pairwise Accuracy with Tie Calibration (acc^{*}_{eq}) estimates evaluator performance based on ability to rank individual translations in the same order as in the ranking derived from the human annotations selected as ground truth (Deutsch et al., 2023).

We describe these measures in more detail in Appendix D.

2.4 Metrics

160

161

162

163

164

165

166

167

168

169

170

172

173

174

175

176

178

179

180

181

183

185

186

189

The automatic evaluators considered – i.e., the MT metrics – are those submitted to the WMT Metrics Shared Task in the 2020, 2022, 2023, and 2024 editions. Additionally, we include several stateof-the-art metrics from recent WMT editions in rankings from previous years, provided they were not trained on the corresponding test sets. Table 3 in Appendix B lists all considered metrics.

3 Results

Table 2 presents the evaluator rankings. Due to space constraints, each table includes only a subset of evaluators. A complete set of results, including all the evaluators, is provided in Appendix C.

Results vary substantially across years and translation directions. Notably, human evaluators do not consistently rank higher than automatic metrics. Under SPA, human evaluators often share clusters of statistical significance with automatic metrics, whereas, under $\operatorname{acc}_{eq}^*$, they are frequently outperformed. For example, in 2020 EN \rightarrow DE, BLEURT-0.2 and BLEURT-20 fall within the same statistical significance cluster as MQM-2020-3 and pSQM-2 under SPA, with pSQM-2 ranking 9th under $\operatorname{acc}_{eq}^*$. Similarly, in 2022 EN \rightarrow DE, MQM-2022-2 and MQM-2022-3 share the top cluster with MetricX-23-QE-XXL* under SPA, with MQM-2022-2 ranking 6th under $\operatorname{acc}_{eq}^*$. Finally, in 2023 and 2024, human evaluators rank consistently below various automatic metrics under both SPA and $\operatorname{acc}_{eq}^*$. Even when restricted to the human evaluators who follow the same protocol as the annotations employed as gold – i.e., MQM – they rank consistently in the top positions solely in 2020.

190

191

192

193

194

195

196

197

198

199

200

201

202

203

205

207

209

210

211

212

213

214

215

216

217

218

219

220

221

These results may indicate that MT Evaluation has reached human parity. Nonetheless, we argue that our findings are insufficient to establish equivalence between human and automatic MT Evaluation and discuss our reasons in the next section.

4 Discussion

First, we outline several factors to consider before claiming human parity in MT Evaluation. Then, we discuss the broader implications of our findings, warning that measuring progress in the field may become increasingly challenging.

Meta-Evaluation measures Certain Meta-Evaluation measures may be inadequate for comparing human and automatic evaluators. In particular, our results consistently rank human evaluators higher under SPA than under $\operatorname{acc}_{eq}^{*}$. This discrepancy may be related to the findings of Perrella et al. (2024b), who show that acc_{eq}^* favors evaluators whose assessments fall within a continuous interval, whereas, as detailed in Appendix A, human evaluators produce discrete assessments.

223

230

231

233

235

239

241

242

243

245

246

247

251

255

257

258

261

262

263

265

266

Annotation quality Certain annotation campaigns might have produced low-quality annotations, either due to a lack of clarity in the annotation guidelines or to the ability of the involved raters. This is particularly concerning in the 2023 data, where, even if restricted to SPA, all human evaluators fall within the second and third clusters of statistical significance, alongside heuristic-based metrics such as BLEU.¹

Benchmarks difficulty Current test sets might be too easy for the MT systems whose translations are being evaluated. Supporting this hypothesis, we observe that sentinel-cand-mqm, a metric that assesses only translation fluency, ranks on par with the human evaluator ESA under SPA, and even higher under acc_{eq}^* (Table 7). This suggests that the evaluated translations may differ only in minor fluency-related nuances. Additionally, previous research has shown that metrics may struggle in unseen domains (Zouhar et al., 2024) and lack sensitivity to specific translation errors such as incorrect number, gender (Karpinska et al., 2022), or word sense disambiguation (Martelli et al., 2024). Thus, before claiming human parity, future studies should compare metrics and humans in more demanding contexts rather than relying solely on standard benchmarks.

4.1 Can We Still Measure Improvements in MT Evaluation?

As discussed, we believe claiming human parity is premature without first addressing the issues outlined above. Nonetheless, with automatic metrics ranking the same as, or higher than, human evaluators in standard benchmarks, our results raise a critical concern about our ability to measure progress in MT Evaluation: What does a higher or lower ranking truly mean?

If a metric ranks higher than a human evaluator using a non-MQM protocol, is the metric a better evaluator, or does it merely align more closely with the score distribution of the MQM protocol? More concerningly, if a metric ranks higher than an MQM evaluator, does this suggest superior evaluation capabilities, or does it simply reflect better alignment with the specific raters who produced the gold annotations? Indeed, Finkelstein et al. (2024) achieved an exceptionally high agreement with gold annotations by explicitly optimizing their metric to align with the raters themselves. More generally, we argue that in current benchmarks it is unclear whether a higher ranking – relative to either a human or an automatic evaluator – reflects genuine improvements in evaluation quality or merely closer alignment with a particular annotation protocol or rater style. 268

269

270

271

272

273

274

275

276

277

278

279

281

282

283

284

285

286

287

289

291

292

293

294

295

296

297

To ensure the reliability of Meta-Evaluation, future research should focus on exploring whether the gap between human and automatic evaluators can be restored. This could be pursued in several ways, including (but not limited to) selecting more challenging test sets, using test sets adversarial to MT metrics (e.g., from domains different from their training data), producing higher-quality human annotations, or designing new annotation protocols that yield stronger IAA. Additionally, greater resources could be allocated to human annotation campaigns – either by collecting multiple annotations per translation to reach a consensus among annotators or by increasing the number of segments in test sets, as suggested by Riley et al. (2024).

5 Conclusions

We incorporate human baselines into the metric 298 rankings from previous editions of the WMT Met-299 rics Shared Task. Our results show that MT met-300 rics frequently rank higher than human evaluators, 301 particularly when the latter follow annotation pro-302 tocols different from MQM - the protocol used as 303 the gold standard. While our findings may indicate 304 human parity, we recommend caution and highlight several issues the research community should 306 address before making such claims. Finally, we 307 discuss a critical concern arising from our findings: 308 the limits of measuring progress in MT Evaluation 309 as automatic metrics approach human baselines. 310 In this respect, we propose research directions to 311 ensure that progress remains measurable or, at the 312 very least, to extend the period during which it can 313 be reliably tracked. 314

¹However, we wish to highlight that our 2023 test set features only 156 segments annotated by all human evaluators (as reported in Table 1), which might have resulted in unreliable estimates of SPA and acc_{eq}^* in this test set.

Limitations

References

- 310
- 31
- 319
- 320
- 32
- 32
- 324
- 52

327

331

332

333

335 336

337

338

340

341

347

349

362

366

326

 David Anugraha, Garry Kuwanto, Lucky Susanto, Derry Tanti Wijaya, and Genta Winata. 2024.
 MetaMetrics-MT: Tuning meta-metrics for machine translation via human preference calibration. In *Proceedings of the Ninth Conference on Machine Translation*, pages 459–469, Miami, Florida, USA. Association for Computational Linguistics.

This study requires test sets annotated by multiple

human evaluators. Consequently, our analysis is

limited to six test sets and four language directions.

ious human evaluators required restricting our anal-

ysis to segments annotated by all of them. As a

result, some test sets contain only a small number

of segments, which might reduce the reliability of

the results. To mitigate this issue, our findings are

supported by statistical significance analyses.

Moreover, assessing the agreement between var-

- Anna Bavaresco, Raffaella Bernardi, Leonardo Bertolazzi, Desmond Elliott, Raquel Fernández, Albert Gatt, Esam Ghaleb, Mario Giulianelli, Michael Hanna, Alexander Koller, André F. T. Martins, Philipp Mondorf, Vera Neplenbroek, Sandro Pezzelle, Barbara Plank, David Schlangen, Alessandro Suglia, Aditya K Surikuchi, Ece Takmaz, and Alberto Testoni. 2024. Llms instead of human judges? a large scale empirical study across 20 nlp evaluation tasks. *Preprint*, arXiv:2406.18403.
- Rachel Bawden, Biao Zhang, Andre Tättar, and Matt Post. 2020. ParBLEU: Augmenting metrics with automatic paraphrases for the WMT'20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 887–894, Online. Association for Computational Linguistics.
- Daniel Deutsch, Rotem Dror, and Dan Roth. 2021. A statistical analysis of summarization evaluation metrics using resampling methods. *Transactions of the Association for Computational Linguistics*, 9:1132–1146.
- Daniel Deutsch, George Foster, and Markus Freitag. 2023. Ties matter: Meta-evaluating modern metrics with pairwise accuracy and tie calibration. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 12914– 12929, Singapore. Association for Computational Linguistics.
- Sören Dreano, Derek Molloy, and Noel Murphy. 2023a. Embed_Llama: Using LLM embeddings for the metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 738–745, Singapore. Association for Computational Linguistics.

Sören Dreano, Derek Molloy, and Noel Murphy. 2023b. Tokengram_F, a fast and accurate tokenbased chrF++ derivative. In *Proceedings of the Eighth Conference on Machine Translation*, pages 730–737, Singapore. Association for Computational Linguistics.

367

368

370

371

372

373

374

375

376

378

379

381

382

383

384

385

386

387

388

389

390

391

392

393

394

395

396

397

398

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

- Muhammad ElNokrashy and Tom Kocmi. 2023. eBLEU: Unexpectedly good machine translation evaluation using simple word embeddings. In *Proceedings of the Eighth Conference on Machine Translation*, pages 746–750, Singapore. Association for Computational Linguistics.
- Patrick Fernandes, Daniel Deutsch, Mara Finkelstein, Parker Riley, André Martins, Graham Neubig, Ankush Garg, Jonathan Clark, Markus Freitag, and Orhan Firat. 2023. The devil is in the errors: Leveraging large language models for fine-grained machine translation evaluation. In *Proceedings of the Eighth Conference on Machine Translation*, pages 1066– 1083, Singapore. Association for Computational Linguistics.
- Mara Finkelstein, Dan Deutsch, Parker Riley, Juraj Juraska, Geza Kovacs, and Markus Freitag. 2024. From jack of all trades to master of one: Specializing llm-based autoraters to a test set. *Preprint*, arXiv:2411.15387.
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021a. Experts, errors, and context: A large-scale study of human evaluation for machine translation. *Transactions of the Association for Computational Linguistics*, 9:1460–1474.
- 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. 2023. 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.
- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, Eleftherios Avramidis, Tom Kocmi, George Foster, Alon Lavie, and André F. T. Martins. 2022. Results of WMT22 metrics shared task: Stop using BLEU – neural metrics are better and more robust. In *Proceedings of the Seventh Conference* on Machine Translation (WMT), pages 46–68, Abu
- 5

482

425Dhabi, United Arab Emirates (Hybrid). Association426for Computational Linguistics.

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

469

470

471

472

473

474

475

476

477

478

479

480

481

- Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-kiu Lo, Craig Stewart, George Foster, Alon Lavie, and Ondřej Bojar. 2021b. Results of the WMT21 metrics shared task: Evaluating metrics with expert-based human evaluations on TED and news domain. In *Proceedings of the Sixth Conference on Machine Translation*, pages 733–774, Online. Association for Computational Linguistics.
- Thamme Gowda, Tom Kocmi, and Marcin Junczys-Dowmunt. 2023. Cometoid: Distilling strong reference-based machine translation metrics into Even stronger quality estimation metrics. In *Proceedings of the Eighth Conference on Machine Translation*, pages 751–755, Singapore. Association for Computational Linguistics.
 - 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.
 - Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
 2021. Measuring massive multitask language understanding. *Preprint*, arXiv:2009.03300.
 - 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.
 - Juraj Juraska, Mara Finkelstein, Daniel Deutsch, Aditya Siddhant, Mehdi Mirzazadeh, and Markus Freitag. 2023. MetricX-23: The Google submission to the WMT 2023 metrics shared task. In *Proceedings* of the Eighth Conference on Machine Translation, pages 756–767, Singapore. Association for Computational Linguistics.
 - Marzena Karpinska, Nishant Raj, Katherine Thai, Yixiao Song, Ankita Gupta, and Mohit Iyyer. 2022.
 DEMETR: Diagnosing evaluation metrics for translation. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 9540–9561, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
 - Fabio Kepler, Jonay Trénous, Marcos Treviso, Miguel Vera, and André F. T. Martins. 2019. OpenKiwi: An open source framework for quality estimation. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pages 117–122, Florence, Italy. 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, Marzena Karpinska, Philipp Koehn, Benjamin Marie, Christof Monz, Kenton Murray, Masaaki Nagata, Martin Popel, Maja Popović, Mariya Shmatova, Steinthór Steingrímsson, and Vilém Zouhar. 2024a. Findings of the WMT24 general machine translation shared task: The LLM era is here but MT is not solved yet. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1–46, 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, 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.
- 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. 2022a. 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.
- Tom Kocmi and Christian Federmann. 2023a. GEMBA-MQM: Detecting translation quality error spans with GPT-4. In *Proceedings of the Eighth Conference on Machine Translation*, pages 768–775, Singapore. Association for Computational Linguistics.
- Tom Kocmi and Christian Federmann. 2023b. Large language models are state-of-the-art evaluators of translation quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*, pages 193–203, Tampere, Finland. European Association for Machine Translation.
- Tom Kocmi, Christian Federmann, Roman Grundkiewicz, Marcin Junczys-Dowmunt, Hitokazu Matsushita, and Arul Menezes. 2021. To ship or not to ship: An extensive evaluation of automatic metrics for machine translation. In *Proceedings of the Sixth Conference on Machine Translation*, pages 478–494, Online. Association for Computational Linguistics.
- Tom Kocmi, Hitokazu Matsushita, and Christian Federmann. 2022b. MS-COMET: More and better human judgements improve metric performance. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 541–548, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

654

655

656

600

Tom Kocmi, Vilém Zouhar, Eleftherios Avramidis, Roman Grundkiewicz, Marzena Karpinska, Maja Popović, Mrinmaya Sachan, and Mariya Shmatova. 2024b. Error span annotation: A balanced approach for human evaluation of machine translation. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1440–1453, Miami, Florida, USA. Association for Computational Linguistics.

542

543

555

557

559

560

561

562

566

568

569

570

572

573

577

578

579

580

581

583

584

588

591

594

595

- Samuel Läubli, Rico Sennrich, and Martin Volk. 2018. Has machine translation achieved human parity? a case for document-level evaluation. In *Proceedings* of the 2018 Conference on Empirical Methods in Natural Language Processing, pages 4791–4796, Brussels, Belgium. Association for Computational Linguistics.
- Yilun Liu, Xiaosong Qiao, Zhanglin Wu, Su Chang, Min Zhang, Yanqing Zhao, Song Peng, Shimin Tao, Hao Yang, Ying Qin, Jiaxin Guo, Minghan Wang, Yinglu Li, Peng Li, and Xiaofeng Zhao. 2022. Partial could be better than whole. HW-TSC 2022 submission for the metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 549–557, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Chi-kiu Lo. 2019. YiSi a unified semantic MT quality evaluation and estimation metric for languages with different levels of available resources. In Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1), pages 507–513, Florence, Italy. Association for Computational Linguistics.
- Arle. Language Technology Lab) Lommel, Hans. Language Technology Lab) Uszkoreit, and Aljoscha.
 Language Technology Lab) Burchardt. 2014a. Multidimensional quality metrics (mqm) : a framework for declaring and describing translation quality metrics. *Translation*, (12):455–463.
- Arle Richard Lommel, Maja Popovic, and Aljoscha Burchardt. 2014b. Assessing inter-annotator agreement for translation error annotation. In *MTE: Workshop* on Automatic and Manual Metrics for Operational Translation Evaluation. International Conference on Language Resources and Evaluation (LREC-14), located at LREC 14, May 26-31, Reykjavik, Iceland. LREC.
- Federico Martelli, Stefano Perrella, Niccolò Campolungo, Tina Munda, Svetla Koeva, Carole Tiberius, and Roberto Navigli. 2024. Dibimt: A gold evaluation benchmark for studying lexical ambiguity in machine translation. *Computational Linguistics*, pages 1–72.
- Nitika Mathur, Timothy Baldwin, and Trevor Cohn.
 2019. Putting evaluation in context: Contextual embeddings improve machine translation evaluation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2799–2808, Florence, Italy. Association for Computational Linguistics.

- Nitika Mathur, Johnny Wei, Markus Freitag, Qingsong Ma, and Ondřej Bojar. 2020. Results of the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 688–725, Online. Association for Computational Linguistics.
- Ananya Mukherjee, Hema Ala, Manish Shrivastava, and Dipti Misra Sharma. 2020. Mee : An automatic metric for evaluation using embeddings for machine translation. In 2020 IEEE 7th International Conference on Data Science and Advanced Analytics (DSAA), pages 292–299.
- Ananya Mukherjee and Manish Shrivastava. 2022a. REUSE: REference-free UnSupervised quality estimation metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 564–568, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Ananya Mukherjee and Manish Shrivastava. 2022b. Unsupervised embedding-based metric for MT evaluation with improved human correlation. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 558–563, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Ananya Mukherjee and Manish Shrivastava. 2023. MEE4 and XLsim : IIIT HYD's submissions' for WMT23 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 800–805, Singapore. Association for Computational Linguistics.
- Ananya Mukherjee and Manish Shrivastava. 2024. chrF-S: Semantics is all you need. In *Proceedings of the Ninth Conference on Machine Translation*, pages 470–474, Miami, Florida, USA. Association for Computational Linguistics.
- Subhajit Naskar, Daniel Deutsch, and Markus Freitag. 2023. Quality estimation using minimum Bayes risk. In *Proceedings of the Eighth Conference on Machine Translation*, pages 806–811, Singapore. 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.
- Stefano Perrella, Lorenzo Proietti, Pere-Lluís Huguet Cabot, Edoardo Barba, and Roberto Navigli. 2024a. Beyond correlation: Interpretable evaluation of machine translation metrics. In *Proceedings of the* 2024 Conference on Empirical Methods in Natural Language Processing, pages 20689–20714, Miami, Florida, USA. Association for Computational Linguistics.
- Stefano Perrella, Lorenzo Proietti, Alessandro Scirè, Edoardo Barba, and Roberto Navigli. 2024b.

Guardians of the machine translation metaevaluation: Sentinel metrics fall in! In *Proceedings* of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 16216–16244, Bangkok, Thailand. Association for Computational Linguistics.

657

672

675

676

677

678

679

687

696

697

704

707

709

710

711

712

713

714

- Stefano Perrella, Lorenzo Proietti, Alessandro Scirè, Niccolò Campolungo, and Roberto Navigli. 2022.
 MaTESe: Machine translation evaluation as a sequence tagging problem. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 569–577, Abu Dhabi, United Arab Emirates (Hybrid). 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.
- Maja Popović. 2017. chrF++: words helping character n-grams. In *Proceedings of the Second Conference on Machine Translation*, pages 612–618, Copenhagen, Denmark. 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.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020a. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference* on Empirical Methods in Natural Language Processing (EMNLP), pages 2685–2702, Online. Association for Computational Linguistics.
- Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020b. Unbabel's participation in the WMT20 metrics shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 911–920, Online. Association for Computational Linguistics.

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. 715

716

717

719

724

725

726

727

729

732

734

735

736

738

739

740

741

742

743

744

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

- Parker Riley, Daniel Deutsch, George Foster, Viresh Ratnakar, Ali Dabirmoghaddam, and Markus Freitag. 2024. Finding replicable human evaluations via stable ranking probability. In Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pages 4908–4919, Mexico City, Mexico. Association for Computational Linguistics.
- Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020a. 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.
- Thibault Sellam, Amy Pu, Hyung Won Chung, Sebastian Gehrmann, Qijun Tan, Markus Freitag, Dipanjan Das, and Ankur Parikh. 2020b. Learning to evaluate translation beyond English: BLEURT submissions to the WMT metrics 2020 shared task. In *Proceedings of the Fifth Conference on Machine Translation*, pages 921–927, Online. Association for Computational Linguistics.
- Matthew Snover, Bonnie Dorr, Rich Schwartz, Linnea Micciulla, and John Makhoul. 2006. A study of translation edit rate with targeted human annotation. In *Proceedings of the 7th Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 223–231, Cambridge, Massachusetts, USA. Association for Machine Translation in the Americas.
- Peter Stanchev, Weiyue Wang, and Hermann Ney. 2019. EED: Extended edit distance measure for machine translation. In *Proceedings of the Fourth Conference on Machine Translation (Volume 2: Shared Task Papers, Day 1)*, pages 514–520, Florence, Italy. Association for Computational Linguistics.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp

- 773 774 775
- 7
- 778 779
- 78
- 78
- 78
- 785
- 786 787

791

7

7

796

778

8 8

802 803

8

805 806 807

809 810

811

812

813 814 815

816 817

- 818
- 819 820

821 822 823

824 825

826 827

827 828 829 Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. 2022. No language left behind: Scaling human-centered machine translation. *Preprint*, arXiv:2207.04672.

- Brian Thompson, Nitika Mathur, Daniel Deutsch, and Huda Khayrallah. 2024. Improving statistical significance in human evaluation of automatic metrics via soft pairwise accuracy. In *Proceedings of the Ninth Conference on Machine Translation*, pages 1222–1234, Miami, Florida, USA. Association for Computational Linguistics.
- Brian Thompson and Matt Post. 2020a. Automatic machine translation evaluation in many languages via zero-shot paraphrasing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 90–121, Online. Association for Computational Linguistics.
- Brian Thompson and Matt Post. 2020b. Paraphrase generation as zero-shot multilingual translation: Disentangling semantic similarity from lexical and syntactic diversity. In *Proceedings of the Fifth Conference on Machine Translation*, pages 561–570, Online. Association for Computational Linguistics.
- Giorgos Vernikos, Brian Thompson, Prashant Mathur, and Marcello Federico. 2022. Embarrassingly easy document-level MT metrics: How to convert any pretrained metric into a document-level metric. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 118–128, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Vasiliy Viskov, George Kokush, Daniil Larionov, Steffen Eger, and Alexander Panchenko. 2023. Semantically-informed regressive encoder score. In *Proceedings of the Eighth Conference on Machine Translation*, pages 815–821, Singapore. Association for Computational Linguistics.
- Yu Wan, Keqin Bao, Dayiheng Liu, Baosong Yang, Derek F. Wong, Lidia S. Chao, Wenqiang Lei, and Jun Xie. 2022a. Alibaba-translate China's submission for WMT2022 metrics shared task. In *Proceedings of the Seventh Conference on Machine Translation (WMT)*, pages 586–592, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.
- Yu Wan, Dayiheng Liu, Baosong Yang, Haibo Zhang, Boxing Chen, Derek Wong, and Lidia Chao. 2022b. UniTE: Unified translation evaluation. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 8117–8127, Dublin, Ireland. Association for Computational Linguistics.
- Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. CharacTer: Translation edit rate on character level. In *Proceedings of the First Conference on Machine Translation: Volume*

2, Shared Task Papers, pages 505–510, Berlin, Germany. Association for Computational Linguistics.

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

- Zhanglin Wu, Yilun Liu, Min Zhang, Xiaofeng Zhao, Junhao Zhu, Ming Zhu, Xiaosong Qiao, Jingfei Zhang, Ma Miaomiao, Zhao Yanqing, Song Peng, Shimin Tao, Hao Yang, and Yanfei Jiang. 2023. Empowering a metric with LLM-assisted named entity annotation: HW-TSC's submission to the WMT23 metrics shared task. In *Proceedings of the Eighth Conference on Machine Translation*, pages 822–828, Singapore. Association for Computational Linguistics.
- Jin Xu, Yinuo Guo, and Junfeng Hu. 2020. Incorporate semantic structures into machine translation evaluation via UCCA. In *Proceedings of the Fifth Conference on Machine Translation*, pages 934–939, Online. Association for Computational Linguistics.
- Wenda Xu, Xian Qian, Mingxuan Wang, Lei Li, and William Yang Wang. 2023. SESCORE2: Learning text generation evaluation via synthesizing realistic mistakes. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics* (Volume 1: Long Papers), pages 5166–5183, Toronto, Canada. Association for Computational Linguistics.
- Wenda Xu, Yi-Lin Tuan, Yujie Lu, Michael Saxon, Lei Li, and William Yang Wang. 2022. Not all errors are equal: Learning text generation metrics using stratified error synthesis. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 6559–6574, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. HellaSwag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4791–4800, Florence, Italy. Association for Computational Linguistics.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert. In *Proceedings of the International Conference on Learning Representations*.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Preprint*, arXiv:2306.05685.
- Vilém Zouhar, Shuoyang Ding, Anna Currey, Tatyana Badeka, Jenyuan Wang, and Brian Thompson. 2024. Fine-tuned machine translation metrics struggle in unseen domains. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 488–500, Bangkok, Thailand. Association for Computational Linguistics.

- 893 894
- 897
- 898

- 901
- 902 903
- 904 905

906

907

908

909

910 911

913 914

916

915

917

919

- 921
- 922 923

925

926

927

A **Human Annotations**

We illustrate, at a high level, how each considered annotation protocol works:

- Multidimensional Quality Metrics (MQM) requires annotators to identify error spans in the translated text, specifying error category and severity, to be selected among Neutral, Minor, Major, and Critical. A translation quality score is derived by assigning a penalty to each error span depending on severity (Lommel et al., 2014a).
- Error Span Annotation (ESA) requires annotators to identify error spans in the translated text, specify error severity, and later assign a scalar quality score from 0 to 100 to the translation (Kocmi et al., 2024b).
 - Scalar Quality Metrics (SQM) requires annotators to assign a scalar quality score from 0 to 6 to the translated text. Following (Freitag et al., 2021a), we use 'pSQM' to refer to SQM conducted by professional annotators.²
 - Direct Assessments + Scalar Quality Metrics (Kocmi et al., 2022a, DA+SQM) requires raters to assign a scalar quality score from 0 to 100 to the translated text. Raters are presented with seven labeled tick marks describing translation quality levels at various score thresholds, similar to the SQM protocol.

Here, for each set of annotations employed in this work (i.e., those reported in Table 1), we indicate the reference paper that released them:

- The MQM-based and pSQM-based annotations for the test sets 2020 EN \rightarrow DE and 2020 $ZH \rightarrow EN$ have been released by Freitag et al. (2021a).
- The MQM-based annotations for the test sets 2022 EN \rightarrow DE and 2022 EN \rightarrow ZH have been released by Freitag et al. (2022) and Riley et al. (2024).
- The DA+SOM-based annotations for the test sets 2022 EN \rightarrow DE and 2022 EN \rightarrow ZH have been released by Kocmi et al. (2022a).

 One set of MQM-based annotations for the test set 2023 EN \rightarrow DE has been released by Freitag et al. (2023).

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

- The ESA-based annotations and the other set of MQM-based annotations for the test set 2023 EN \rightarrow DE have been released by Kocmi et al. (2024b).
- The ESA-based annotations for the test set 2024 EN \rightarrow ES have been released by Kocmi et al. (2024a).
- The MQM-based annotations for the test set 2024 EN \rightarrow ES have been released by Freitag et al. (2024).

B Metrics

Table 3 lists the complete set of automatic evaluators considered in this work.

Full Rankings С

Tables 4, 5, 6, and 7 present the ranking containing all tested evaluators.

D **Meta-Evaluation Measures**

In this section, we describe the two Meta-Evaluation measures used in our work, as listed in Section 2.3.

D.1 Soft Pairwise Accuracy (SPA)

Thompson et al. (2024) introduced Soft Pairwise Accuracy (SPA) as an extension of Pairwise Accuracy (Kocmi et al., 2021, PA).

Given a test set t, which consists of N_t source segments and M_t translations generated by the respective M_t MT systems (as described in Section 2.1), PA counts how often an evaluator e ranks system pairs in the same order as the ground truth g. Let a_{ij} be 1 if evaluator e ranks systems i and j in the same order as the ground truth and 0 otherwise, where $i, j \in \{0, ..., Mt\}$. Then, PA is defined as:

$$PA = \binom{N}{2}^{-1} \sum_{i=0}^{M_t} \sum_{j=i+1}^{M_t} a_{ij}$$
(1)

SPA extends PA by incorporating the confidence with which an evaluator and the ground truth rank two MT systems. Confidence is represented using statistical p-values. Specifically, p_{ij}^e denotes the *p*-value associated with the statistical hypothesis that system i is better than system j according to

²In this work, we use only annotations produced by professional annotators or translators. Therefore, we exclude cSQM and Direct Assessments (DA) - which were crowdsourced from the 2020 test sets.

Metric	Reference paper	Metric	Reference paper
all-rembert-20	(Mathur et al., 2020)		(A
BAQ_dyn	(Mathur et al., 2020)	metametrics_mt_mqm	(Anugrana et al., 2024)
BAQ_static	(Mathur et al., 2020)	metametrics_mt_mqm_qe*	(Anugraha et al., 2024)
BERT-base-L2	(Mathur et al., 2020)	MetricX-23-QE-XXL*	(Juraska et al., 2023)
BERT-large-L2	(Mathur et al., 2020)	MetricX-23-XXL	(Juraska et al., 2023)
BERTScore	(Zhang et al., 2020)	MetricX-24-Hybrid	(Juraska et al., 2024)
BLCOM_1	(Freitag et al., 2024)	MetricX-24-Hybrid-QE*	(Juraska et al., 2024)
BLEU	(Papineni et al., 2002)	metricx xxl MQM 2020	(Freitag et al., 2022)
BLEURI	(Sellam et al., 2020a)	mre-score-labse-regular	(Viskov et al., 2023)
BLEURI-0.1-all	(Mathur et al., 2020)	MS-COMET-22	(Kocmi et al. 2022b)
BLEURI-0.1-en	(Mathur et al., 2020)	MS COMET OF 22*	(Kocmi et al. 2022b)
BLEURI-0.2 DI EUDT 20	(Mathur et al., 2020)	OpenVirvi Bert*	(Koelin et al., 20220) (Koelin et al., 2010)
BLEURI-20	(Sellam et al., 2020a)	OpenKiwi-Bert*	(Kepler et al., 2019)
DIEURI-combi	(Mathur et al., 2020)	OpenKiwi-XLMR*	(Kepler et al., 2019)
bright go*	(Senam et al., 2020b) (Emittag et al., 2024)	parbleu	(Bawden et al., 2020)
Calibri COMET22	(Freitag et al., 2024) (Freitag et al., 2023)	parchrf++	(Bawden et al., 2020)
Calibri-COMET22-OF*	(Freitag et al. 2023)	paresim-1	(Bawden et al., 2020)
CharacTER	(Wang et al. 2016)	prism	(Thompson and Post, 2020a)
chrF	(Popović 2015)	prismRef	(Thompson and Post, 2020a,b)
chrF++	(Popović 2017)	PrismRefMedium	(Thompson and Post, 2020a,b)
chrfS	(Mukheriee and Shriyastava, 2024)	PrismRefSmall	(Thompson and Post 2020a b)
COMET	(Rei et al., 2020b)	prismSrc*	(Thompson and Post, 2020a,b)
COMET-20	(Rei et al., 2020a)	Pandam avanama*	(Fraitag at al. 2022)
COMET-22	(Rei et al., 2022a)	Random-systiame"	(Fleftag et al., 2023)
COMET-2R	(Rei et al., 2020b)	REUSE*	(Mukherjee and Shrivastava, 2022a)
COMET-HTER	(Rei et al., 2020b)	sentBLEU	(Papineni et al., 2002)
COMET-MQM	(Rei et al., 2020b)	sentinel-cand-mqm*	(Perrella et al., 2024b)
COMET-QE*	(Rei et al., 2021)	sentinel-ref-mqm	(Perrella et al., 2024b)
COMET-Rank	(Rei et al., 2020b)	sentinel-src-mqm*	(Perrella et al., 2024b)
COMETKiwi*	(Rei et al., 2022b)	SEScore	(Xu et al., 2022)
CometKiwi-XL*	(Rei et al., 2023)	sescoreX	(Xu et al., 2023)
CometKiwi-XXL*	(Rei et al., 2023)	spBLEU	(Team et al., 2022)
cometoid22-wmt22*	(Gowda et al., 2023)	SWSS+METEOR	(Xu et al. 2020)
damonmonli	(Freitag et al., 2024)	TFR	(Snover et al 2006)
docWMT22CometDA	(Vernikos et al., 2022)	tokongram E	(Draano et al. 2023b)
docWMT22CometKiwiDA*	(Vernikos et al., 2022)	UniTE	$(W_{ex} \rightarrow z_{e}^{1} - 2022h_{e})$
eBLEU	(ElNokrashy and Kocmi, 2023)	UMIE	(Wan et al., 2022b,a)
EED	(Stanchev et al., 2019)	UnilE-src*	(wan et al., 2022b)
embed_nama	(Mothur et al., 2025a)	XCOMET	(Guerreiro et al., 2024)
f200spBI EU	(Mathur et al., 2019) (Team et al., 2022)	XCOMET-Ensemble	(Guerreiro et al., 2024)
GEMBA-MOM*	(Kocmi and Federmann, 2023a)	XCOMET-QE*	(Guerreiro et al., 2024)
gemba esa*	(Freitag et al. 2024)	XCOMET-QE-Ensemble*	(Guerreiro et al., 2024)
HWTSC-Teacher-Sim*	(1 in et al 2022)	XLsim	(Mukherjee and Shrivastava, 2023)
KG-BERTScore*	(Wu et al., 2023)	XLsimMqm*	(Mukherjee and Shrivastava, 2023)
MaTESe	(Perrella et al., 2022)	YiSi-0	(Lo, 2019)
MaTESe-QE*	(Perrella et al., 2022)	YiSi-1	$(L_0, 2019)$
mBERT-L2	(Mathur et al., 2020)	YiSi-2*	$(L_0, 2019)$
mbr-metricx-qe*	(Naskar et al., 2023)	Visi combi	(Mothur et al. 2020)
MEE	(Mukherjee et al., 2020)		(Mathur et al., 2020)
MEE4	(Mukherjee and Shrivastava, 2022b)	yisi1-translate	(wathur et al., 2020)

Table 3: List of all automatic evaluators considered, i.e., MT metrics, associated with their reference paper. If some metrics do not have a dedicated reference paper, we provide the Metrics Shared Task results paper in which they were submitted.

evaluator e, while p_{ij}^g represents the corresponding p-value for the ground truth g. SPA is then defined as follows:

970

971

972

973

974

975

976

977

978

979

980

981

$$SPA = \binom{N}{2}^{-1} \sum_{i=0}^{M_t} \sum_{j=i+1}^{M_t} 1 - |p_{ij}^g - p_{ij}^e| \quad (2)$$

Thus, SPA rewards an evaluator for expressing confidence levels similar to those of the ground truth and penalizes deviations.

D.2 Pairwise Accuracy with Tie Calibration (acc^{*}_{ea})

Deutsch et al. (2023) introduced $\operatorname{acc}_{eq}^*$ to account for tied scores in Meta-Evaluation. Unlike PA and SPA, $\operatorname{acc}_{eq}^*$ is a segment-level measure, meaning it evaluates a metric's ability to estimate the quality of individual translations rather than MT systems. Specifically, acc_{eq}^* counts how often an evaluator *e* ranks pairs of translations of the same source in the same order as the ground truth *g*, accounting for tied scores.

982

983

984

985

986

987

988

989

990

991

992

993

994

995

Let C be the number of translation pairs ranked in the same order by both the evaluator e and the ground truth g. Similarly, let D denote the pairs ranked in the opposite order. The terms T_e and T_g represent pairs tied only in the evaluator's scores and only in the ground truth, respectively. Lastly, T_{eg} refers to pairs tied in both the evaluator's scores and the ground truth. $\operatorname{acc}_{eg}^*$ is then defined as:

$$\operatorname{acc}_{eq}^{*} = \frac{C + T_{eg}}{C + D + T_{e} + T_{q} + T_{eq}}$$
(3)

Tie Calibration Many automatic metrics produce assessments on a continuous scale, such as the real numbers in the interval [0, 1]. As a consequence, these metrics rarely, if ever, produce tied scores, resulting in $T_e \approx 0$ and $T_{eg} \approx 0$. The Tie Calibration algorithm addresses this issue by estimating a threshold value ϵ_e for each evaluator e, such that two assessments e_i and e_j are considered tied if $|e_i - e_j| \leq \epsilon_e$.

996

997

998

1002

1005

1006

1007

1008

1009

1010

1011

1014

1015

1016

1017

1018

1019

1020 1021

1022

1023

1025 1026

1027 1028

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1043

E Fair Extraction of Evaluators from Human Annotations

The human evaluation campaigns conducted by Freitag et al. (2021a) and Riley et al. (2024) produced multiple annotations for each translation (MQM and pSQM produced three annotations per translation, as reported in the columns '2020' and '2022' in Table 1). As described in Section 2.1, these annotation campaigns distributed the annotation workload among multiple raters. Since some of these annotations were used as ground truth, and since we are interested in measuring the performance of independent evaluators, we prevent the same rater from contributing to both a human evaluator and the ground truth, or to two distinct evaluators, simultaneously. For example, in the '2020' $EN \rightarrow DE$ test set, six raters provided three annotations per translation. We would like to extract three human evaluators from these annotations, using one as the ground truth and the other two as human evaluators (MQM-2020-2 and MQM-2020-3 in Table 2). To achieve this, we need to partition the six raters, for instance, into three groups of two raters each. However, not all raters annotated the entire set of source segments, and the workload distribution did not allow for a rater assignment that covers all annotated segments. Therefore, to maximize the number of segments in our test sets, we need to solve the following optimization problem: Find the largest subset of segments for which the set of raters can be partitioned into three disjoint groups i.e., the three human evaluators.

More generally, let us define a test set $t = \{s_1, ..., s_{N_t}\}$ as a set of N_t segments. Each segment was annotated by k out of R raters, with $S = \{r_1, ..., r_R\}$ representing the set of raters. Our goal is to determine a partition $\Pi = \{S_1, ..., S_k\}$ of S and a subset $u \subseteq t$ such that u is the largest

subset in which every segment has been annotated1044by exactly one rater from each of the k sets in the1045partition Π .1046

To solve this optimization problem, we formulate it as an Integer Linear Programming (ILP)1047problem and solve it using the PuLP³ Python library. We applied this procedure to the '2020' and '2022' test sets.1050

³https://coin-or.github.io/pulp/.

	EN→DE				ZH→EN			
	SI	PA	ac	c_{eq}^{*}	SI	PA	ac	c_{eq}^{*}
Metric	Rank	Acc.	Rank	Acc.	Rank	Acc.	Rank	Acc.
MQM-2020-2	1	96.45	1	58.86	1	88.10	1	55.70
pSQM-1	1	95.59	6	49.41	1	79.16	13	43.89
MQM-2020-3	2	90.39	2	56.84	1	92.06	2	52.80
BLEURT-0.2	2	86.81	4	50.81	2	72.59	3	50.57
pSQM-2	2	85.87	9	46.97	1	89.33	9	46.77
BLEURT-20	2	85.52	3	51.68	3	67.46	4	50.12
pSQM-3	2	84.61	6	49.38	1	87.94	7	47.88
all-rembert-20	3	79.19	4	51.04	3	66.41	3	50.61
BLEURT-extended	3	75.55	5	50.21	3	64.00	3	50.74
COMET-MQM	4	71.39	7	48.21	4	55.43	6	48.49
BLEURT-0.1-all	4	71.38	7	48.63	2	71.04	5	49.54
COMET	4	71.09	8	47.36	4	56.01	5	49.28
COMET-QE*	4	70.59	8	47.82	3	58.37	8	47.09
COMET-HTER	5	65.71	8	47.62	4	54.79	5	49.30
mBERT-L2	5	65.03	10	45.48	4	56.49	6	48.97
COMET-2R	6	58.12	9	46.43	4	55.99	4	50.20
COMET-Rank	6	54.78	14	41.31	3	58.16	14	43.57
OpenKiwi-XLMR*	6	53.25	11	44.11	4	53.29	8	47.23
OpenKiwi-Bert*	6	52.01	16	39.98	3	59.55	11	45.13
prism	6	51.92	11	43.59	4	57.88	8	47.56
Yisi-combi	7	51.10	12	42.63	-	_	_	_
bleurt-combi	7	51.10	12	42.63	-	_	-	_
esim	7	50.72	14	41.35	4	52.90	10	46.19
chrF	7	49.86	13	42.05	5	47.70	13	44.09
EED	7	49.81	15	40.94	5	45.41	14	43.64
paresim-1	7	49.54	14	41.37	4	53.34	10	46.15
chrF++	7	48.87	13	41.99	5	48.96	12	44.27
YiSi-1	7	48.79	12	42.70	4	52.74	7	48.01
CharacTER	7	47.71	16	40.45	5	48.84	13	44.01
BLEURT-0.1-en	7	47.43	15	40.96	4	57.29	7	48.26
YiSi-0	7	46.23	17	39.78	5	46.47	14	43.60
TER	7	45.98	16	40.15	6	39.68	15	43.34
parchrf++	7	45.57	13	42.25	5	48.68	12	44.25
MEE	7	45.31	14	41.61	4	52.91	13	43.94
sentBLEU	7	44.41	15	41.07	4	50.45	15	43.37
parbleu	8	41.38	15	41.01	4	50.28	15	43.43
yisi1-translate	8	39.76	12	42.60	4	52.28	11	44.70
YiSi-2*	8	38.44	18	34.36	5	43.35	12	44.60

Table 4: 2020

	EN→DE				EN→ZH			
	SI	PA	ace	$\operatorname{acc}_{eq}^*$		SPA		c_{eq}^*
Metric	Rank	Acc.	Rank	Acc.	Rank	Acc.	Rank	Acc.
MetricX-23-QE-XXL*	1	94.89	3	57.64	2	83.92	2	47.43
MQM-2022-2	1	94.49	6	55.55	2	80.82	3	47.05
MQM-2022-3	1	92.59	1	61.06	1	87.22	2	47.56
MetricX-23-XXL	2	92.34	2	59.27	1	87.69	1	48.43
COMET-22	2	91.63	5	56.51	2	84.08	3	46.74
COMET-20	2	91.28	9	52.42	2	80.56	7	43.81
CometKiwi*	2	89.51	7	53.77	3	75.36	8	43.21
BLEURT-20	3	88.20	7	53.33	3	77.80	7	43.84
metricx_xxl_MQM_2020	3	88.10	3	57.43	1	87.04	3	46.89
COMET-QE*	3	85.51	10	51.69	3	78.33	7	43.61
MS-COMET-22	3	85.37	8	53.13	1	85.18	6	44.92
CometKiwi-XXL*	3	84.43	7	53.27	2	81.25	2	47.28
UniTE	4	82.77	4	57.03	2	83.88	5	45.86
UniTE-src*	4	81.55	6	55.00	4	65.74	7	43.53
CometKiwi-XL*	4	81.13	8	52.73	2	81.56	4	46.33
YiSi-1	4	78.91	13	48.26	4	70.72	8	43.23
MATESE	5	78.03	7	53.48	_	_	_	_
BERTScore	5	75.61	14	47.57	4	70.69	8	43.28
SEScore	5	75.16	12	50.45	_	_	_	_
MS-COMET-QE-22*	5	74.44	12	50.37	2	78.84	9	42.51
MEE4	5	74.19	15	46.81	_	_	_	_
chrF	5	73.05	16	46.38	3	72.67	10	41.87
f200spBLEU	5	71.04	15	46.84	4	71.76	10	41.85
HWTSC-Teacher-Sim*	5	69.68	13	48.10	4	68.43	11	40.53
DA+SQM	6	66.61	16	46.03	2	82.95	12	36.26
MATESE-QE*	6	65.42	11	51.06	_	_	_	_
BLEU	6	65.00	15	46.51	4	67.31	13	34.28
REUSE*	7	37.95	17	43.58	5	33.46	12	35.89

Table 5: 2022

	EN→DE						
	SI	PA	$\operatorname{acc}_{ea}^{*}$				
Metric	Rank	Acc.	Rank	Acc.			
GEMBA-MQM*	1	96.07	3	57.50			
CometKiwi-XL*	1	95.65	4	57.20			
CometKiwi-XXL*	1	95.51	3	57.56			
MetricX-23-XXL	1	94.98	2	60.13			
MetricX-23-QE-XXL*	1	94.82	1	61.44			
COMET	1	94.59	5	56.12			
BLEURT-20	1	94.35	5	55.63			
docWMT22CometDA	1	94.26	6	54.59			
Calibri-COMET22-QE*	1	94.18	12	48.48			
XCOMET-QE-Ensemble*	1	93.99	3	58.30			
sescoreX	1	93.94	6	55.22			
XCOMET-Ensemble	2	93.75	2	59.91			
cometoid22-wmt22*	2	93.66	3	57.60			
DA+SQM	2	93.00	13	46.41			
CometKiwi*	2	92.55	4	57.36			
docWMT22CometKiwiDA*	2	92.44	6	54.36			
MS-COMET-QE-22*	2	92.33	6	54.93			
ESA-1	2	92.21	13	46.49			
KG-BERTScore*	2	92.05	5	56.50			
mbr-metricx-qe*	2	91.89	3	58.03			
Calibri-COMET22	2	91.19	9	51.53			
MaTESe	2	89.89	7	52.79			
mre-score-labse-regular	2	89.78	8	52.51			
mqm-2023-2	3	89.45	14	43.08			
prismRef	3	89.32	11	50.18			
f200spBLEU	3	89.20	8	51.86			
ESA-2	3	89.14	11	49.77			
YiSi-1	3	88.88	7	53.33			
XLsim	3	88.56	9	51.32			
BLEU	3	87.84	9	51.17			
BERTscore	3	87.03	9	51.54			
MEE4	3	86.55	8	51.85			
eBLEU	3	85.37	10	50.50			
tokengram_F	4	84.83	9	50.99			
chrF	4	83.88	10	50.91			
embed_llama	4	81.46	12	48.10			
Random-sysname*	5	59.50	15	40.37			
prismSrc*	6	27.58	14	42.33			

Table 6: 2023

	EN→ES					
	SI	PA	ac	c_{eq}^{*}		
Metric	Rank	Acc.	Rank	Acc.		
CometKiwi-XXL*	1	86.12	4	67.24		
gemba_esa*	1	85.72	3	67.68		
COMET-22	1	82.37	5	66.60		
bright-qe*	1	81.77	4	67.39		
ESA	2	80.12	8	63.84		
XCOMET-QE*	2	80.10	3	67.99		
metametrics_mt_mqm_hybrid_kendall	2	80.10	1	68.95		
XCOMET	2	79.96	2	68.67		
MetricX-24-Hybrid	2	79.75	1	69.20		
BLCOM_1	2	79.17	6	65.02		
MetricX-24-Hybrid-QE*	2	79.09	2	68.92		
sentinel-cand-mqm*	2	78.54	5	66.39		
BLEURT-20	2	75.96	7	64.48		
metametrics_mt_mqm_qe_kendall.seg.s*	3	73.29	4	67.49		
CometKiwi*	3	71.74	5	66.51		
PrismRefMedium	3	70.93	11	61.39		
PrismRefSmall	3	70.52	10	61.51		
YiSi-1	3	70.51	11	61.44		
BERTScore	3	67.75	11	61.41		
chrF	3	66.73	13	61.05		
damonmonli	3	66.37	9	62.10		
chrfS	4	64.31	11	61.37		
spBLEU	4	63.19	12	61.08		
BLEU	5	60.67	13	61.04		
MEE4	5	60.36	10	61.57		
sentinel-ref-mqm	6	44.19	13	61.04		
sentinel-src-mqm*	6	44.19	13	61.04		
XLsimMqm*	6	39.25	12	61.11		

Table 7: 2024