

Simpson’s Paradox and the Accuracy-Fluency Tradeoff in Translation

Anonymous ACL submission

Abstract

A good translation should be faithful to the source and should respect the norms of the target language. We address a theoretical puzzle about the relationship between these objectives. On one hand, intuition and some prior work suggest that accuracy and fluency should trade off against each other, and that capturing every detail of the source can only be achieved at the cost of fluency. On the other hand, quality assessment researchers often suggest that accuracy and fluency are highly correlated and difficult for human raters to distinguish (Callison-Burch et al., 2007). We show that the tension between these views is an instance of Simpson’s paradox, and that accuracy and fluency are positively correlated at the level of the corpus but trade off at the level of individual source segments. We further suggest that the relationship between accuracy and fluency is best evaluated at the segment (or sentence) level, and that the trade off between these dimensions has implications both for assessing translation quality and developing improved MT systems.

1 Introduction

No translation can simultaneously satisfy all possible goals, and translation is therefore an art of navigating competing objectives (Darwish, 2008). Many objectives are discussed in the literature, but two in particular seem especially fundamental. The first is accuracy (also known as fidelity or adequacy), or the goal of preserving the information in the source text (ST). The second is fluency, or the goal of producing target text (TT) that respects the norms of the target language (TL) and is easy for the recipient to process (Kunilovskaya, 2023).

Here we study the relationship between accuracy and fluency and work with two operationalizations of these notions. The first relies on human judgments of accuracy and fluency collected in prior work on translation quality estimation (Castilho et al., 2018). The second relies on probabilities

estimated using neural machine translation (NMT) models. Given a source-translation pair (x, y) , $p(x|y)$ corresponds to accuracy, and $p(y)$ corresponds to fluency (Teich et al., 2020). $p(x|y)$ will be low if y fails to preserve all of the information in x , and $p(y)$ will be low if y violates the norms of the target language. To highlight that model estimates $p(x|y)$ and $p(y)$ are related to but distinct from human ratings of accuracy and fluency, we refer to $p(x|y)$ as accuracy_M and $p(y)$ as fluency_M.

Some parts of the literature argue that accuracy trades off with fluency. In Figure 1a, the blue dots are translations of the same source segment, and Table 1 shows three translations that illustrate the same kind of tradeoff. A translator choosing between these alternatives cannot simultaneously maximize accuracy and fluency, because the most accurate translations are not the most fluent, and vice versa. Teich et al. (2020) argues that accuracy_M and fluency_M should trade off in this way, and the same view is implicitly captured by noisy-channel models of translation (Brown et al., 1993), which aim to generate translations y that maximize $p(y|x) \propto p(x|y)p(y)$. Typically these models include weights for the two components $p(x|y)$ and $p(y)$ that can be interpreted as the extent to which accuracy_M is prioritized over fluency_M, or vice versa (Yu et al., 2016; Yee et al., 2019; Yu et al., 2020; Müller et al., 2020).

An opposing view of the relationship between accuracy and fluency emerges from the literature on quality estimation. Here the common wisdom is that accuracy and fluency are highly correlated and practically indistinguishable to human annotators (Callison-Burch et al., 2007; Banchs et al., 2015; Mathur, 2021, but see Djiko 2019; Sulem et al. 2020). As a result, accuracy and fluency are conflated as a single assessment score in recent WMT General Machine Translation Tasks, with more emphasis given to accuracy than fluency (Farhad et al., 2021; Kocmi et al., 2022, 2023).

Translation		accuracy	fluency	accuracy _M	fluency _M	$\log p(\mathbf{y} \mathbf{x})$
(i) Ich gab Ihnen eine Rückerstattung des Buches.		23.0	25.0	-10.81	-56.0	-10.31
(ii) Ich habe Ihnen eine Rückerstattung des Buches ausgestellt.		24.3	24.7	-6.13	-64.0	-12.13
(iii) Ich stellte Ihnen eine Rückerstattung des Buches aus.		25.0	23.0	-6.44	-70.0	-14.75

Table 1: Translations of “*I issued you a refund of the book.*” from English to German, which correspond to three of the orange dots in Figure 1. Option (i) is acceptable but *gab* (past tense of give) is less accurate than the conjugations of *ausstellen* (issue) used in (ii) and (iii). Option (iii) is the least natural because *stellte ... aus* (Präteritum tense) is typically used only in formal writing.

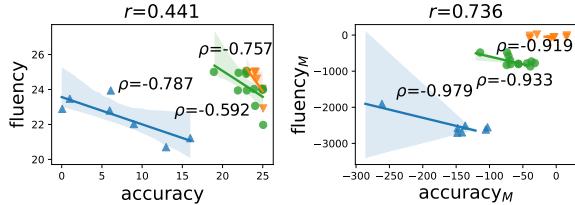


Figure 1: Simpson’s paradox. Each panel shows translations of three source segments indexed by color and marker shape. At the source segment level, accuracy and fluency (left) and $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$ (right) both show negative correlations. At the corpus level, both pairs of dimensions show positive correlations ($p < .05$; see panel labels). Source segments and translations are drawn from past WMT General Task submissions and data points have been jittered for clarity.

We argue that the conflict between these views is an instance of Simpson’s paradox (Yuan et al., 2021), which occurs when a relationship at one level of analysis (e.g. the corpus level) is reversed at a different level (e.g. the segment or sentence level). Figure 1 shows how the correlation between accuracy and fluency can be positive over a miniature corpus including translations of three source segments even though the correlation for each individual source segment is negative. Of the two levels of analysis, the segment level is the appropriate level for understanding how humans and machine translation systems should choose among possible translations of a source segment. The central goal of our work is therefore to establish that the correlation between accuracy and fluency is negative at the level of individual source segments.¹

2 Tradeoff between $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$

Because accuracy_M and fluency_M have formal definitions, we start with these dimensions.

2.1 Theoretical formulation and simulation

Let \mathbf{Y} be a finite set of translations of source segment \mathbf{x} , and let $\vec{p}_{\mathbf{x}|\mathbf{y}}$ and $\vec{p}_{\mathbf{y}}$ denote log probabil-

ity vectors that include accuracy_M and fluency_M scores for all $\mathbf{y} \in \mathbf{Y}$.² We use the Pearson correlation between the two vectors:

$$\rho = \text{corr}(\vec{p}_{\mathbf{x}|\mathbf{y}}, \vec{p}_{\mathbf{y}}) \quad (1)$$

to quantify the tradeoff between accuracy_M and fluency_M across translations of \mathbf{x} . If $\rho > 0$ there is no tradeoff, and the translations with higher accuracy_M also tend to have higher fluency_M. If $\rho < 0$ the dimensions trade off, and improving a translation along one dimension tends to leave it worse along the other.

Suppose that a translator is choosing among several good translations \mathbf{y} that all have near-maximal values of $p(\mathbf{y}|\mathbf{x})$. Because $p(\mathbf{y}|\mathbf{x}) \propto p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$, it follows that accuracy_M and fluency_M trade off within the set of candidates. To validate this informal argument, we ran simulations in which \mathbf{x} and \mathbf{y} are both numeric vectors drawn from a Gaussian joint distribution $P(\mathbf{x}, \mathbf{y})$.³ For each source “segment” \mathbf{x} there are infinitely many translations \mathbf{y} , and we focus on the best, or those with highest $p(\mathbf{y}|\mathbf{x})$.

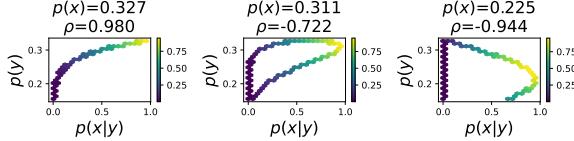
We initially assume that both \mathbf{x} and \mathbf{y} are one-dimensional vectors. Figure 2a shows the relationship between $p(\mathbf{y})$ and $p(\mathbf{x}|\mathbf{y})$ for 3 samples \mathbf{x} . Each point in each panel corresponds to a candidate translation \mathbf{y} , and candidates with highest $p(\mathbf{y}|\mathbf{x})$ are shown in yellow. The correlation above each panel results from applying Equation 1 to all translations with $p(\mathbf{y}|\mathbf{x})$ above the 90th percentile (i.e. all points in the brightest part of each plot). The first source “segment” \mathbf{x} (leftmost panel) has relatively high probability $p(\mathbf{x})$, and no tradeoff is observed in this case. The tradeoff emerges, however, and becomes increasingly strong as \mathbf{x} moves away from the mode of the distribution $p(\mathbf{x})$.

Figure 2b shows that the tradeoff persists when the dimensionality of \mathbf{x} and \mathbf{y} is increased. The

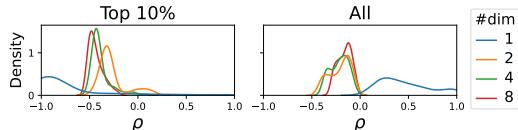
²There are infinitely many possible translations, but here we only consider a pool of human and machine generations.

³All simulation details are provided in Appendix.

¹Code will be made available upon publication.



(a) Simulation with one-dimensional x and y . The three panels correspond to three different source “segments” x of decreasing probability $p(x)$, and the points in each panel are candidate translations y . Brighter colors indicate translations with larger $p(y|x)$. Pearson correlations across translations ranked in the top 10% based on $p(y|x)$ are shown at the top of each panel.



(b) Kernel density plots of tradeoffs across the top 10% (left) and across all translation choices (right). The tradeoff persists in higher dimensional space, and is stronger when selecting only y with the highest values of $p(y|x)$.

Figure 2: Tradeoffs between $p(x|y)$ and $p(y)$ in synthetic data.

density plot for each dimensionality is based on a sample of 100 source “segments” (rather than the 3 in Figure 2a), and at all dimensionalities the majority of source “segments” induce tradeoffs. The tradeoffs are stronger (i.e. correlations more negative) when the candidate translations consist of the y with highest $p(y|x)$, but for all dimensions except $n = 1$ most source “segments” still induce a tradeoff even if all candidate translations are considered. These results provide theoretical grounds to expect that tradeoffs will also occur in real translations generated by humans and machines.

2.2 Human and machine translation

We now show that human and machine translations show the same tradeoff between accuracy _{M} and fluency _{M} , which correspond to $p(x|y)$ and $p(y)$ estimated by an NMT model.

Data. We analyze 15 translation studies from CRITT TPR-DB (CRITT) that include 13 language pairs (Carl et al., 2016b). We also use a subset of the Russian Learner Translator Corpus (RLTC) that has been aligned at the sentence level by Kunilovskaya (2023). For machine translation, we use WMT test sets that are annotated with Multidimensional Quality Metrics labels (MTMQM) (Freitag et al., 2021a,b; Zerva et al., 2022; Freitag et al., 2023). To reduce spurious correlations, we remove duplicate translations and source segments with fewer than four unique translations. Additional details are provided in the Appendix.

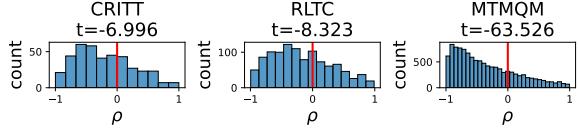


Figure 3: Tradeoffs between estimated $p(x|y)$ and $p(y)$ across source segments from three corpora. Paired-sample t-tests against randomly permuted $p(y)$ and $p(x|y)$ are shown at the top of each panel.

Models. We use NLLB-200’s 3.3B variant model (Costa-jussà et al., 2022) to estimate $p(y|x)$ and $p(x|y)$.⁴ For consistency, we also extract $p(y)$ based on the same model, skipping all inputs except for special tokens (e.g., <eos> tags).⁵ All probabilities are log scaled.

Results. Figure 3 is analogous to the densities in Figure 2b, and shows distributions of tradeoff scores for source segments in CRITT, RLTC and MTMQM. In all three cases most source segments induce tradeoffs (i.e. produce negative correlations). To test for statistical significance we compared the actual distributions against randomly permuted data. The results of all paired-sample t-tests are significant ($p < .001$), and are included in the figure.⁶ When samples are aggregated at the corpus level, $p(y)$ and $p(x|y)$ show significant positive correlations ($p < .001$) for CRITT ($r = .625$), RLTC ($r = .685$) and MTMQM ($r = .675$), showing that Simpson’s paradox applies in all three cases.

3 Tradeoff between accuracy and fluency

We now turn to human ratings of accuracy and fluency, and demonstrate that the two are again negatively correlated at the segment level.

Data. Only RLTC and MTMQM are rated by human annotators. The subset of RLTC released by Kunilovskaya (2023) includes accuracy and fluency scores derived from error annotations. For MTMQM, we aggregate accuracy scores based on “Accuracy” errors, and aggregate errors in “Fluency”, “Terminology”, “Style” and “Locale convention” for fluency rating. Targets that are labelled “Non-translation” receive scores of zero for both accuracy and fluency scores.⁷

⁴NLLB model card

⁵To ensure reproducibility across models, we repeat our analysis in Appendix using M2M100 (Fan et al., 2021).

⁶Each permuted data set is created by randomly shuffling the pairings of $p(x|y)$ and $p(y)$ within the set of possible translations of each source segment.

⁷Further details of accuracy and fluency score derivation

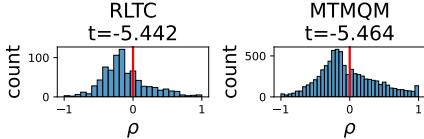


Figure 4: Tradeoffs between human ratings of accuracy and fluency across segments from two corpora. Paired-sample t-tests against randomly permuted scores are shown at the top of each panel.

Results. Figure 4 shows correlations at the level of individual source segments. The majority of correlations are negative, and paired-sample t-tests reveal that both distributions are significantly ($p < .001$) different from distributions obtained from random permutations. The results therefore suggest that accuracy and fluency (as rated by humans) trade off at the level of individual segments. At the corpus level, accuracy and fluency are positively correlated for MTMQM ($r = .396, p < .001$), but not RLTC ($r = -.085, p < .001$), suggesting that Simpson’s paradox applies only to MTMQM.

Figure 4 is directly analogous to Figure 3, and we expected that source segments which showed strong tradeoffs (i.e. extreme negative correlations) in Figure 3 would also show strong tradeoffs in Figure 4. The two tradeoff measures, however, were uncorrelated,⁸ which suggests that accuracy_M and fluency_M overlap only partially with human ratings of accuracy and fluency.

A similar conclusion is suggested by Figure 5, which shows Pearson correlations of translation probability ($p(\mathbf{y}|\mathbf{x})$), accuracy_M ($p(\mathbf{x}|\mathbf{y})$) and fluency_M ($p(\mathbf{y})$) with human ratings of accuracy and fluency for RLTC and MTMQM.⁹ As expected, $p(\mathbf{x}|\mathbf{y})$ shows a higher correlation with accuracy than fluency, and $p(\mathbf{y})$ shows the opposite pattern. Figure 5 however, suggests that $p(\mathbf{x}|\mathbf{y})$ is not superior to $p(\mathbf{y}|\mathbf{x})$ as a predictor of accuracy, and that $p(\mathbf{y})$ is not superior to $p(\mathbf{y}|\mathbf{x})$ as a predictor of fluency. One reason why our model estimates of accuracy and fluency depart from human ratings is that accuracy_M and fluency_M are sensitive to segment length. For example, a longer segment will have lower fluency_M than a shorter segment even if the two are both perfectly idiomatic.

are available in Appendix.

⁸The Pearson correlations between the two tradeoff measures for RLTC and MTMQM are $r = .005, p = .885$ and $r = .024, p = .04$.

⁹Values are ranked by percentile.

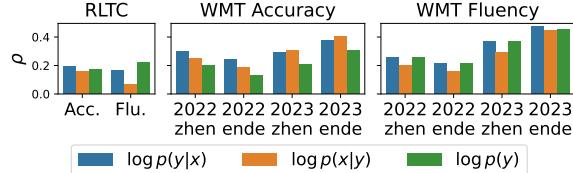


Figure 5: accuracy_M and fluency_M predict human accuracy and fluency ratings for RLTC and WMT submissions to the general translation task in 2022 and 2023. zhen and ende refer to Chinese-English and English-German language pairs. All correlations reported are significant ($p < .001$).

4 Conclusion

We showed that accuracy and fluency and $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$ both trade off when translating individual source segments. This finding suggests that current protocols for assessing translation quality may need to be adjusted. Human assessments for recent WMT General Task are performed with Direct Assessment and Scalar Quality Metrics (DA+SQM) (Koci et al., 2022, 2023). This approach conflates meaning preservation and grammar into a single score indicative of overall quality of a translation. In contrast, MQM is much more costly, but produces highly detailed scores that use multiple sub-categories for both accuracy and fluency. Future approaches could therefore consider a middle ground that extends DA+SQM to include accuracy and fluency as independent aspects as in WMT16 (Bojar et al., 2016). This direction would allow automatic MT evaluation metrics such as BLEURT (Sellam et al., 2020) and COMET (Rei et al., 2022) (both fine-tuned to DA scores) to be adapted to provide independent scores for accuracy and fluency.

Our results also suggest the value of developing MT models that navigate the accuracy-fluency tradeoff in human-like ways. In some settings (e.g. translating legal texts) accuracy is more important (Popović, 2020; Martindale and Carpuat, 2018; Vela and Tan, 2015; Specia et al., 2011; Martindale et al., 2019), but in others (e.g. translating informal conversation) fluency may take priority (Poibeau, 2022; Frankenberg-Garcia, 2022). One natural approach builds on noisy channel models (Yu et al., 2016; Yee et al., 2019; Müller et al., 2020), which incorporate both $p(\mathbf{y})$ and $p(\mathbf{x}|\mathbf{y})$ along with trade-off parameters that specify the relative weights of the two. Tuning these parameters for specific registers may allow a model to find the right balance between accuracy and fluency in each case.

282 5 Limitations

283 Although we provided evidence for both accuracy-
284 fluency and accuracy_M-fluency_M tradeoffs in
285 translation, we did not explore factors which pre-
286 dict which source segments are likely to produce
287 the greatest tradeoffs. Outside of our simulation
288 we do not have access to ground-truth values of
289 $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$, and are only able to approximate
290 these values using specific NMT models. Our work
291 is also limited by the fact that MTMQM only in-
292 cludes translations generated by certain kinds of
293 NMT models, and it is possible that our results do
294 not generalize to translations generated by other
295 types of models, such as statistical or rule-based
296 MT systems. Finally, both RLTC and MTMQM
297 have accuracy and fluency ratings derived from
298 error annotations that are very similar in range. This
299 constraint makes quality assessment and compari-
300 son at the segment level challenging.

301 Ethics Statement

302 We do not foresee any potential risks and harmful
303 use of our work. Our analyses are based on licensed
304 data which are freely available for academic use.

305 References

- 306 Fabio Alves and José Luiz Gonçalves. 2013. Investi-
307 gating the conceptual-procedural distinction in the
308 translation process: A relevance-theoretic analysis of
309 micro and macro translation units. *Target. Interna-*
310 *tional Journal of Translation Studies*, 25(1):107–124.
- 311 Rafael E Banchs, Luis F D’Haro, and Haizhou Li. 2015.
312 Adequacy–fluency metrics: Evaluating mt in the con-
313 tinuous space model framework. *IEEE/ACM Trans-*
314 *actions on Audio, Speech, and Language Processing*,
315 23(3):472–482.
- 316 Ondrej Bojar, Rajen Chatterjee, Christian Federmann,
317 Yvette Graham, Barry Haddow, Matthias Huck, Anton-
318 io Jimeno Yepes, Philipp Koehn, Varvara Loga-
319 cheva, Christof Monz, et al. 2016. Findings of
320 the 2016 conference on machine translation (wmt16).
321 In *First conference on machine translation*, pages
322 131–198. Association for Computational Linguistics.
- 323 Peter F Brown, Stephen A Della Pietra, Vincent J
324 Della Pietra, Robert L Mercer, et al. 1993. The math-
325 ematics of statistical machine translation: Parameter
326 estimation.
- 327 Chris Callison-Burch, Cameron Shaw Fordyce, Philipp
328 Koehn, Christof Monz, and Josh Schroeder. 2007.
329 (meta-) evaluation of machine translation. In *Pro-*
330 *ceedings of the Second Workshop on Statistical Ma-*
331 *chine Translation*, pages 136–158.

- 332 Michael Carl, Akiko Aizawa, and Masaru Yamada. 332
333 2016a. English-to-Japanese translation vs. dictation 333
334 vs. post-editing: Comparing translation modes in a 334
335 multilingual setting. In *Proceedings of the Tenth In- 335
336 ternational Conference on Language Resources and 336
337 Evaluation (LREC’16)*, pages 4024–4031. 337
- 338 Michael Carl and M Cristina Toledo Báez. 2019. Ma- 338
339 chine translation errors and the translation process: 339
340 A study across different languages. *Journal of Spe- 340
341 cialised Translation*, 31:107–132. 341
- 342 Michael Carl, Moritz Schaeffer, and Srinivas Banga- 342
343 loore. 2016b. The CRITT translation process research 343
344 database. In *New directions in empirical translation 344
345 process research*, pages 13–54. Springer. 345
- 346 Sheila Castilho, Stephen Doherty, Federico Gaspari, 346
347 and Joss Moorkens. 2018. Approaches to human 347
348 and machine translation quality assessment. *Transla- 348
349 tion quality assessment: From principles to practice*, 349
350 pages 9–38. 350
- 351 Marta R Costa-jussà, James Cross, Onur Çelebi, Maha 351
352 Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe 352
353 Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, 353
354 et al. 2022. No language left behind: Scaling 354
355 human-centered machine translation. *arXiv preprint 355
356 arXiv:2207.04672*. 356
- 357 Ali Darwish. 2008. *Optimality in translation*. 357
358 Writoscope Publishers. 358
- 359 Gabriel Armand Djako. 2019. *Lexical ambiguity in 359
360 machine translation and its impact on the evalua- 360
361 tion of output by users*. Ph.D. thesis, Saarländische 361
362 Universitäts-und Landesbibliothek. 362
- 363 Barbara Dragsted. 2010. Coordination of reading and 363
364 writing processes in translation: An eye on uncharted 364
365 territory. In *Translation and Cognition*, pages 41–62. 365
366 John Benjamins Publishing Company. 366
- 367 Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi 367
368 Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep 368
369 Baines, Onur Çelebi, Guillaume Wenzek, Vishrav 369
370 Chaudhary, et al. 2021. Beyond english-centric mul- 370
371 tilingual machine translation. *Journal of Machine 371
372 Learning Research*, 22(107):1–48. 372
- 373 Akhbardeh Farhad, Arkhangorodsky Arkady, Biesialska 373
374 Magdalena, Bojar Ondřej, Chatterjee Rajen, Chaud- 374
375 hary Vishrav, Marta R Costa-jussà, España-Bonet 375
376 Cristina, Fan Angela, Federmann Christian, et al. 376
377 2021. Findings of the 2021 conference on machine 377
378 translation (wmt21). In *Proceedings of the Sixth 378
379 Conference on Machine Translation*, pages 1–88. As- 379
380 sociation for Computational Linguistics. 380
- 381 Ana Frankenberg-Garcia. 2022. Can a corpus-driven 381
382 lexical analysis of human and machine translation 382
383 unveil discourse features that set them apart? *Target*, 383
384 34(2):278–308. 384

385	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. <i>Transactions of the Association for Computational Linguistics</i> , 9:1460–1474.	442
386		443
387		444
388		445
389		446
390		
391	Markus Freitag, Nitika Mathur, Chi-ku Lo, Eleftherios Avramidis, Ricardo Rei, Brian Thompson, Tom Kocmi, Frédéric Blain, Daniel Deutsch, Craig Stewart, et al. 2023. Results of wmt23 metrics shared task: Metrics might be guilty but references are not innocent. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 578–628.	447
392		448
393		449
394		450
395		451
396		
397		
398	Markus Freitag, Ricardo Rei, Nitika Mathur, Chi-ku 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 <i>Proceedings of the Sixth Conference on Machine Translation</i> , pages 733–774.	452
399		453
400		454
401		
402		
403		
404		
405	Kristian Tangsgaard Hvelplund Jensen, Annette C Sjørup, and Laura Winther Balling. 2009. Effects of 11 syntax on 12 translation. <i>Copenhagen Studies in Language</i> , 38:319–336.	460
406		461
407		462
408		463
409		464
410		
411		
412		
413		
414		
415		
416		
417	Tom Kocmi, Eleftherios Avramidis, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Markus Freitag, Thamme Gowda, Roman Grundkiewicz, et al. 2023. Findings of the 2023 conference on machine translation (wmt23): Llms are here but not quite there yet. In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 1–42.	465
418		466
419		467
420		468
421		469
422		
423		
424	Tom Kocmi, Rachel Bawden, Ondřej Bojar, Anton Dvorkovich, Christian Federmann, Mark Fishel, Thamme Gowda, Yvette Graham, Roman Grundkiewicz, Barry Haddow, et al. 2022. Findings of the 2022 conference on machine translation (wmt22). In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 1–45.	470
425		471
426		472
427		473
428	Maria Kunilovskaya. 2023. <i>Translationese indicators for human translation quality estimation (based on English-to-Russian translation of mass-media texts)</i> . Ph.D. thesis, University of Wolverhampton.	474
429		475
430		476
431		
432		
433		
434	Marianna Martindale and Marine Carpuat. 2018. Fluency over adequacy: A pilot study in measuring user trust in imperfect mt. In <i>Proceedings of the 13th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)</i> , pages 13–25.	477
435		478
436		479
437		480
438		481
439	Marianna Martindale, Marine Carpuat, Kevin Duh, and Paul McNamee. 2019. Identifying fluently inadequate output in neural and statistical machine translation. In <i>Proceedings of Machine Translation Summit XVII: Research Track</i> , pages 233–243.	482
440		483
441		484
442	Bartolomé Mesa-Lao. 2014. Gaze behaviour on source texts: An exploratory study comparing translation and post-editing. In <i>Post-editing of machine translation: Processes and applications</i> , pages 219–245. Cambridge Scholars Publishing.	485
443		486
444		487
445		488
446		489
447	Mathias Müller, Annette Rios Gonzales, and Rico Sennrich. 2020. Domain robustness in neural machine translation. In <i>Proceedings of the 14th Conference of the Association for Machine Translation in the Americas (Volume 1: Research Track)</i> , pages 151–164.	490
448		491
449		492
450		493
451		494
452	Jean Nitzke. 2019. <i>Problem solving activities in post-editing and translation from scratch: A multi-method study</i> . Language Science Press.	495
453		496
454		497
455	Dagmara Płońska. 2016. Problems of literality in french-polish translations of a newspaper article. <i>New directions in empirical translation process research: exploring the CRITT TPR-DB</i> , pages 279–291.	498
456		499
457		500
458		501
459		502
460	Thierry Poibeau. 2022. On “human parity” and “super human performance” in machine translation evaluation. In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 6018–6023.	503
461		504
462		505
463		506
464		507
465	Maja Popović. 2020. Relations between comprehensibility and adequacy errors in machine translation output. In <i>Proceedings of the 24th Conference on Computational Natural Language Learning</i> , pages 256–264.	508
466		509
467		510
468		511
469		512
470	Ricardo Rei, José GC De Souza, Duarte Alves, Chrysoula Zerva, Ana C Farinha, Taisiya Glushkova, Alon Lavie, Luisa Coheur, and André FT Martins. 2022. Comet-22: Unbabel-ist 2022 submission for the metrics shared task. In <i>Proceedings of the Seventh Conference on Machine Translation (WMT)</i> , pages 578–585.	513
471		514
472		515
473		516
474		517
475		518
476		519
477	Márcia Schmaltz, Igor AL da Silva, Adriana Pagano, Fabio Alves, Ana Luísa V Leal, Derek F Wong, Lidia S Chao, and Paulo Quaresma. 2016. Cohesive relations in text comprehension and production: An exploratory study comparing translation and post-editing. <i>New Directions in Empirical Translation Process Research: Exploring the CRITT TPR-DB</i> , pages 239–263.	520
478		521
479		522
480		523
481		524
482		525
483		526
484		527
485	Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020. Bleurt: Learning robust metrics for text generation. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7881–7892.	528
486		529
487		530
488		531
489		532
490	Annette Camilla Sjørup. 2013. <i>Cognitive effort in metaphor translation: An eye-tracking and key-logging study</i> . Frederiksberg: Copenhagen Business School (CBS).	533
491		534
492		535
493		536
494		537
495	Lucia Specia, Najeh Hajlaoui, Catalina Hallett, and Wilker Aziz. 2011. Predicting machine translation adequacy. In <i>Proceedings of Machine Translation Summit XIII: Papers</i> .	538
496		539
497		540

498 Elior Sulem, Omri Abend, and Ari Rappoport. 2020. Semantic structural decomposition for neural machine
499 translation. In *Proceedings of the ninth joint conference*
500 on lexical and computational semantics, pages
501 50–57.

503 Elke Teich, José Martínez Martínez, and Alina
504 Karakanta. 2020. Translation, information theory
505 and cognition. *The Routledge Handbook of Translation*
506 and Cognition, pages 9781315178127–24.

507 Bram Vanroy. 2021. *Syntactic difficulties in translation*.
508 Ph.D. thesis, Ghent University.

509 Mihaela Vela and Liling Tan. 2015. Predicting machine
510 translation adequacy with document embeddings. In
511 *Proceedings of the Tenth Workshop on Statistical*
512 *Machine Translation*, pages 402–410.

513 Lucas Nunes Vieira, Natalie Zelenka, Roy Youdale, Xi-
514 aochun Zhang, and Michael Carl. 2023. Translating
515 science fiction in a CAT tool: Machine translation
516 and segmentation settings. *Translation & Interpreting*,
517 15(1):216–235.

518 Kyra Yee, Yann Dauphin, and Michael Auli. 2019.
519 Simple and effective noisy channel modeling for
520 neural machine translation. In *Proceedings of the*
521 *2019 Conference on Empirical Methods in Natural*
522 *Language Processing and the 9th International Joint*
523 *Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5696–5701.

525 Lei Yu, Phil Blunsom, Chris Dyer, Edward Grefenstette,
526 and Tomas Kociský. 2016. The neural noisy channel.
527 In *International Conference on Learning Representations*.

529 Lei Yu, Laurent Sartran, Wojciech Stokowiec, Wang
530 Ling, Lingpeng Kong, Phil Blunsom, and Chris Dyer.
531 2020. Better document-level machine translation
532 with bayes’ rule. *Transactions of the Association for*
533 *Computational Linguistics*, 8:346–360.

534 Fei Yuan, Longtu Zhang, Huang Bojun, and Yaobo
535 Liang. 2021. Simpson’s bias in nlp training. In
536 *Proceedings of the AAAI Conference on Artificial*
537 *Intelligence*, volume 35, pages 14276–14283.

538 Chrysoula Zerva, Frédéric Blain, Ricardo Rei, Piyawat
539 Lertvittayakumjorn, José GC De Souza, Steffen Eger,
540 Diptesh Kanodia, Duarte Alves, Constantin Orăsan,
541 Marina Fomicheva, et al. 2022. Findings of the wmt
542 2022 shared task on quality estimation. In *Proceedings*
543 of the Seventh Conference on Machine Translation
544 (WMT), pages 69–99.

545 A Appendix

546 A.1 Simulation

547 We assume \mathbf{x} and \mathbf{y} are drawn from a Gaussian
548 joint distribution $P(\mathbf{x}, \mathbf{y})$ where both \mathbf{x} and \mathbf{y} are
549 centered at zero. We set an initial square matrix
550 A with dimensionality equal to the total number

of dimensions in \mathbf{x} and \mathbf{y} combined. Assuming
all elements in \mathbf{x} and \mathbf{y} have $\sigma^2 = 1$ and pairwise
positive covariance, all diagonal elements of A
are set to 1 and other elements 0.7. To ensure
the covariance matrix is positive semi-definite, we
replace the initial matrix A with a final covariance
matrix defined as $A^\top A$.

Based on this setup we first draw 100 random
samples of \mathbf{x} from $P(\mathbf{x})$. For each \mathbf{x} we then
sample 100,000 values of \mathbf{y} from a distribution
 $q(\mathbf{y}) = \prod_i q(y_i)$, where each element y_i of \mathbf{y} is
sampled uniformly within two standard deviations
of its mean. We then compute $p(\mathbf{x}|\mathbf{y})$, $p(\mathbf{y})$ and
 $p(\mathbf{y}|\mathbf{x})$ for each (\mathbf{x}, \mathbf{y}) pair using the known joint
 $P(\mathbf{x}, \mathbf{y})$.

566 A.2 Data specification

567 A.2.1 Corpora

CRITT Translation Process Research Database
(Carl et al., 2016b) is a collection of translation
behavioural data in the area of Translation Pro-
cess Research. From the public CRITT database
we obtain 15 studies across 13 pairs of languages:
RUC17 (enzh, Carl and Báez, 2019), ENJA15 (enja,
Carl et al., 2016a), NJ12 (enhi, Carl et al., 2016b),
STC17 (enzh, Carl and Báez, 2019), SG12 (ende,
Nitzke, 2019), ENDU20 (ennl, Vanroy, 2021),
BML12 (enes, Mesa-Lao, 2014), ACS08 (daen,
Sjørup, 2013), MS13 (ptzh, Schmaltz et al., 2016),
JLG10 (pten, Alves and Gonçalves, 2013), BD13
(daen, Dragsted, 2010), LWB09 (daen, Jensen et al.,
2009), DG01 (plfr, Płońska, 2016), BD08 (daen,
Dragsted, 2010) and CREATIVE (enzh, Vieira
et al., 2023).¹⁰ After deduplication and removing
source segments with fewer than 4 unique transla-
tions, the total number of source segments included
is 399, each with an average of 10.9 unique transla-
tions.

RLTC is a subset of the Russian Learner Trans-
lator Corpus that has been segment-aligned by Ku-
nilovskaya (2023). We include a total of 1079
source segments from 5 genres: ‘Essay’, ‘Infor-
mational’, ‘Speech’, ‘Interview’ and ‘Educational’.
The average number of unique translations for each
source segment is 10.5.

MTMQM is obtained from (Freitag et al.,
2021a), which contains translations of TED talks
and news data from the test sets of WMT general

¹⁰<https://sites.google.com/site/centrerevolutioninnovation/tpr-db/public-studies>

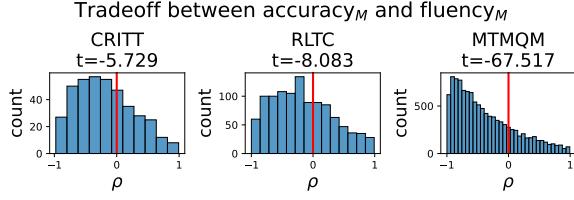


Figure 6: Histogram of tradeoffs between estimated $p(\mathbf{x}|\mathbf{y})$ and $p(\mathbf{y})$ estimated by M2M100, which is analogous to Figure 3 in the main text. When analyzed at the corpus level, the correlations ($p < .001$) for CRITT, RLTC and MTMQM are .689, .703 and .801 respectively.

task between 2020 and 2023.¹¹ The translations are annotated with MQM labels. After preprocessing we are left with 11219 source segments and an average of 9.9 unique translations per source segment.

A.2.2 Accuracy and fluency scores by MQM

For MTMQM, we aggregate accuracy and fluency penalty score by error annotations, and based on weights recommended by Freitag et al. (2021a). Major and minor errors receive penalties of 5 and 1 respectively. Fluency/Punctuation is assigned a penalty of 0.1. We calculate the final rating as $s_c = \max(0, 25 - e_c)$, where e_c denotes the total penalty in error category c .¹² Because some systems submit the same translation but receive different ratings, we average these scores and remove the duplicate entries.

A.3 Alternative result with M2M100 translation model

In Figure 6 and 7, we replicate our findings of accuracy _{M} and fluency _{M} in Section 2 and 3 with estimates of M2M100 (1.2B variant) (Fan et al., 2021).¹³

A.4 Tradeoff examples

Tables 2, 3 and 4 include the full set of translations plotted in Figure 1. The tables specify accuracy, fluency, accuracy _{M} , fluency _{M} and translation probability $p(\mathbf{y}|\mathbf{x})$ for each segment. All translations listed are submissions to WMT General Task between 2020 to 2022.

¹¹<https://github.com/google/wmt-mqm-human-evaluation>

¹²The maximum score is set at 25 because a translation can only get a maximum of 25 penalty score.

¹³https://huggingface.co/facebook/m2m100_1.2B

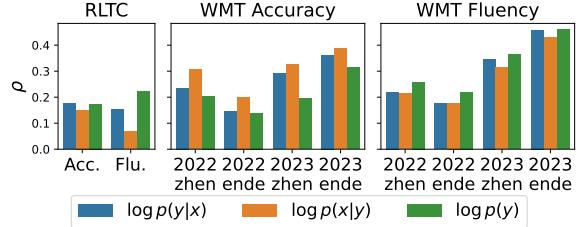


Figure 7: accuracy _{M} and fluency _{M} estimated by M2M100 are also predictive of human accuracy and fluency ratings ($p < .05$). The figure is analogous to Figure 5.

Ich gab Ihnen eine Rückerstattung des Buches. {accuracy: 23.0, fluency: 25.0, accuracy _{M} : -10.81, fluency _{M} : -56.0, log $p(\mathbf{y} \mathbf{x})$: -10.31}
Ich habe dir eine Rückerstattung des Buches ausgestellt. {accuracy: 23.0, fluency: 25.0, accuracy _{M} : -5.84, fluency _{M} : -62.5, log $p(\mathbf{y} \mathbf{x})$: -12.44}
Ich habe dir das Buch zurückerstattet. {accuracy: 23.0, fluency: 25.0, accuracy _{M} : -17.5, fluency _{M} : -44.25, log $p(\mathbf{y} \mathbf{x})$: -7.63}
Ich habe Ihnen das Buch erstattet. {accuracy: 24.0, fluency: 25.0, accuracy _{M} : -15.19, fluency _{M} : -43.25, log $p(\mathbf{y} \mathbf{x})$: -9.06}
Ich habe Ihnen das Buch zurückerstattet. {accuracy: 24.2, fluency: 25.0, accuracy _{M} : -17.25, fluency _{M} : -43.5, log $p(\mathbf{y} \mathbf{x})$: -7.28}
Ich habe Ihnen eine Rückerstattung des Buches ausgestellt. {accuracy: 24.3, fluency: 24.67, accuracy _{M} : -6.13, fluency _{M} : -64.0, log $p(\mathbf{y} \mathbf{x})$: -12.13}
Ich stellte Ihnen eine Rückerstattung des Buches aus. {accuracy: 25.0, fluency: 23.0, accuracy _{M} : -6.44, fluency _{M} : -70.0, log $p(\mathbf{y} \mathbf{x})$: -14.75}
Ich habe Ihnen eine Rückerstattung für das Buch erteilt. {accuracy: 25.0, fluency: 24.0, accuracy _{M} : -11.56, fluency _{M} : -63.0, log $p(\mathbf{y} \mathbf{x})$: -14.19}

Table 2: Translations of *I issued you a refund of the book*. (plotted in orange in Figure 1).

Ashanti Development arbeitet seit fast 20 Jahren mit einer wachsenden Anzahl von Gemeinden in der Region Ashanti in Ghana zusammen und unterstützt sie in den Bereichen Wasser und sanitäre Einrichtungen, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft.

{accuracy: 19.0, fluency: 25.0, accuracy_M: -120.5, fluency_M: -498.0, log p($\mathbf{y}|\mathbf{x}$): -27.0}

Ashanti Development arbeitet seit fast zwanzig Jahren mit einer ständig wachsenden Anzahl von Gemeinden in der Region Ashanti in Ghana zusammen, engagiert sich mit Gemeinden und unterstützt Wasser und Sanitärversorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft. Gemeinschaften erlangen das Wissen, um ihre eigene Entwicklung einzubetten und zu unterstützen.

{accuracy: 22.0, fluency: 24.0, accuracy_M: -47.5, fluency_M: -748.0, log p($\mathbf{y}|\mathbf{x}$): -47.25}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Zahl von Gemeinden in der Ashanti-Region in Ghana zusammen, engagiert sich für Gemeinden und bietet Unterstützung in den Bereichen Wasser und sanitäre Einrichtungen, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft. Communities erwerben das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 22.0, fluency: 25.0, accuracy_M: -46.5, fluency_M: -832.0, log p($\mathbf{y}|\mathbf{x}$): -49.0}

Ashanti Development arbeitet seit 20 Jahren mit einer immer größeren Zahl von Gemeinden in der Region Ashanti in Ghana zusammen, engagiert sich mit Gemeinden und unterstützt Wasser und Sanitärversorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft.

{accuracy: 23.0, fluency: 24.9, accuracy_M: -101.0, fluency_M: -516.0, log p($\mathbf{y}|\mathbf{x}$): -39.25}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Anzahl von Gemeinden in der Ashanti-Region Ghanas zusammen, indem es sich mit Gemeinden beschäftigt und ihnen Unterstützung in den Bereichen Wasser und Sanitärversorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft bietet.

{accuracy: 23.0, fluency: 25.0, accuracy_M: -98.5, fluency_M: -652.0, log p($\mathbf{y}|\mathbf{x}$): -29.625}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Zahl von Gemeinden in der Ashanti-Region in Ghana zusammen, arbeitet mit Gemeinden zusammen und unterstützt sie in den Bereichen Wasser und Abwasserentsorgung, Bildung, Gesundheitswesen, Baumpflanzung und Landwirtschaft. Gemeinschaften erwerben das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 23.0, fluency: 24.0, accuracy_M: -53.0, fluency_M: -828.0, log p($\mathbf{y}|\mathbf{x}$): -42.5}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Anzahl von Gemeinden in der Ashanti-Region in Ghana zusammen, engagiert sich für Gemeinden und unterstützt sie bei Wasser- und Sanitärversorgung, Bildung, Gesundheitswesen, Baumpflanzung und Landwirtschaft. Gemeinschaften gewinnen das Wissen, um ihre eigene Entwicklung einzubetten und zu unterstützen.

{accuracy: 24.0, fluency: 23.0, accuracy_M: -47.5, fluency_M: -784.0, log p($\mathbf{y}|\mathbf{x}$): -45.25}

Ashanti Development arbeitet seit fast 20 Jahren mit einer stetig wachsenden Anzahl von Gemeinden in der Ashanti-Region in Ghana zusammen, engagiert sich mit Gemeinden und bietet Unterstützung in den Bereichen Wasserversorgung und Abwasserentsorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft. Gemeinden erwerben das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 24.0, fluency: 24.0, accuracy_M: -49.5, fluency_M: -848.0, log p($\mathbf{y}|\mathbf{x}$): -42.25}

Ashanti Development arbeitet seit fast 20 Jahren mit einer stetig wachsenden Anzahl von Gemeinschaften in der Ashanti-Region von Ghana zusammen, engagiert sich in den Gemeinschaften und bietet Unterstützung in den Bereichen Wasser und Sanitär, Bildung, Gesundheitswesen, Baumpflanzung und Landwirtschaft. Die Gemeinschaften erwerben das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 25.0, fluency: 22.0, accuracy_M: -50.25, fluency_M: -828.0, log p($\mathbf{y}|\mathbf{x}$): -43.0}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Zahl von Gemeinden in der Ashanti-Region in Ghana zusammen, engagiert sich für Gemeinden und leistet Unterstützung bei Wasser- und Sanitärversorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft. Gemeinschaften erlangen das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 25.0, fluency: 24.0, accuracy_M: -45.75, fluency_M: -816.0, log p($\mathbf{y}|\mathbf{x}$): -45.0}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Zahl von Gemeinden in der Ashanti-Region in Ghana zusammen und unterstützt sie in den Bereichen Wasserversorgung und Abwasserentsorgung, Bildung, Gesundheitsversorgung, Baumpflanzung und Landwirtschaft. Die Gemeinden erlangen das Wissen, um ihre eigene Entwicklung zu fördern und zu unterstützen.

{accuracy: 25.0, fluency: 24.0, accuracy_M: -74.0, fluency_M: -768.0, log p($\mathbf{y}|\mathbf{x}$): -42.0}

Ashanti Development arbeitet seit fast 20 Jahren mit einer ständig wachsenden Zahl von Gemeinden in der Ashanti-Region in Ghana zusammen, engagiert sich für Gemeinden und leistet Unterstützung bei Wasser- und Sanitärversorgung, Bildung, Gesundheitswesen, Baumpflanzung und Landwirtschaft. Gemeinschaften erwerben das Wissen, um ihre eigene Entwicklung zu verankern und zu unterstützen.

{accuracy: 25.0, fluency: 24.0, accuracy_M: -46.25, fluency_M: -812.0, log p($\mathbf{y}|\mathbf{x}$): -46.25}

Table 3: Translations of *Ashanti Development has been working with an ever-expanding number of communities in the Ashanti region of Ghana for approaching 20 years, engaging with communities and providing support with water and sanitation, education, healthcare, tree planting and farming. Communities gain the knowledge to embed and support their own development.* These translations are plotted in green in Figure 1.

Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor 2009 vom Subaru-Teleskop gefunden wurde. „Es ist vernünftig, einen Protohaufen in der Nähe eines massereichen Objekts wie Himiko zu finden. Wir sind jedoch überrascht zu sehen, dass Himiko nicht im Zentrum des Protohaufens lag, sondern am Rande 500 Millionen Lichtjahre vom Zentrum entfernt.“ Sagte Masami Ouchi, ein Teammitglied am Nationalen Astronomischen Observatorium von Japan und der Universität von Tokio, die Himiko im Jahr 2009 entdeckte, dass die Beziehung zwischen den Himiko und den Himiko-Klöstern noch immer nicht verstanden wird.

{accuracy: 0.0, fluency: 22.9, accuracy_M: -286.0, fluency_M: -1904.0, log p($\mathbf{y}|\mathbf{x}$): -139.0}

""""""""Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor vom Subaru-Teleskop im Jahr 2009 gefunden wurde. """"""""Es ist vernünftig, einen Protocluster in der Nähe eines massiven Objekts wie Himiko zu finden. Wir sind jedoch überrascht zu sehen, dass Himiko nicht im Zentrum des Protoclusters, sondern am Rand 500 Millionen Lichtjahre vom Zentrum entfernt war"""""""""", sagte Masami Ouchi, ein Teammitglied am Nationalen Astronomischen Observatorium von Japan und der Universität von Tokio, der Himiko im Jahr 2009 entdeckte. Ironischerweise soll die mythologische Königin Himiko auch abgeschieden von ihrem Volk gelebt haben. Ouchi fährt fort: """"""""Es ist immer noch nicht verstanden, warum Himiko nicht im Zentrum liegt. Diese Ergebnisse werden ein Schlüssel für das Verständnis der Beziehung zwischen Haufen und massiven Galaxien sein"""""""""

{accuracy: 1.0, fluency: 23.4, accuracy_M: -125.0, fluency_M: -2624.0, log p($\mathbf{y}|\mathbf{x}$): -103.5}

""""""""Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor vom Subaru-Teleskop im Jahr 2009 gefunden wurde. „Es ist vernünftig, einen Protocluster in der Nähe eines massiven Objekts, wie Himiko, zu finden. Allerdings sind wir überrascht zu sehen, dass Himiko nicht im Zentrum des Protoclusters, sondern am Rande 500 Millionen Lichtjahre vom Zentrum entfernt war.“, sagte Masami Ouchi, ein Teammitglied am Nationalen Astronomischen Observatorium von Japan und der Universität von Tokio, der Himiko im Jahr 2009 entdeckte. Ironischerweise soll die mythologische Königin Himiko auch abgeschieden von ihrem Volk gelebt haben. Ouchi fährt fort: """"""""Es ist immer noch nicht verstanden, warum Himiko nicht im Zentrum liegt. Diese Ergebnisse werden ein Schlüssel für das Verständnis der Beziehung zwischen Haufen und und massiven galaxien sein."""""""""

{accuracy: 6.0, fluency: 24.0, accuracy_M: -121.0, fluency_M: -2688.0, log p($\mathbf{y}|\mathbf{x}$): -143.0}

""""""""Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor 2009 vom Subaru-Teleskop gefunden wurde. „Es ist vernünftig, einen Protocluster in der Nähe eines massiven Objekts wie Himiko zu finden. Wir sind jedoch überrascht zu sehen, dass sich Himiko nicht im Zentrum des Protoclusters befand, sondern am Rande 500 Millionen Lichtjahre vom Zentrum entfernt“, sagte Masami Ouchi, Teammitglied am National Astronomical Observatory of Japan und der Universität Tokio, der Himiko 2009 entdeckte. Ironischerweise soll die mythologische Königin Himiko auch abgeschieden von ihrem Volk gelebt haben. Ouchi fährt fort: """"""""Es ist immer noch nicht verstanden, warum Himiko nicht im Zentrum liegt. Diese Ergebnisse werden ein Schlüssel für das Verständnis der Beziehung zwischen Haufen und massiven Galaxien sein.""""""""

{accuracy: 6.0, fluency: 22.7, accuracy_M: -126.0, fluency_M: -2592.0, log p($\mathbf{y}|\mathbf{x}$): -123.0}

Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor 2009 vom Subaru-Teleskop gefunden wurde. „Es ist vernünftig, einen Protocluster in der Nähe eines massiven Objekts wie Himiko zu finden. Wir sind jedoch überrascht zu sehen, dass sich Himiko nicht im Zentrum des Protoclusters befand, sondern am Rand 500 Millionen Lichtjahre vom Zentrum entfernt“, sagte Masami Ouchi, Teammitglied am National Astronomical Observatory of Japan und der Universität Tokio, der Himiko 2009 entdeckte. Ironischerweise soll die mythologische Königin Himiko auch abseits ihres Volkes im Kloster gelebt haben. Ouchi fährt fort: „Es ist immer noch nicht verstanden, warum Himiko sich nicht im Zentrum befindet. Diese Ergebnisse werden ein Schlüssel zum Verständnis der Beziehung zwischen Haufen und massiven Galaxien sein.“

{accuracy: 9.0, fluency: 22.0, accuracy_M: -131.0, fluency_M: -2512.0, log p($\mathbf{y}|\mathbf{x}$): -108.0}

""""""""Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das 2009 vom Subaru-Teleskop gefunden wurde. """"""""Es ist vernünftig, einen Protokluster in der Nähe eines massiven Objekts zu finden, wie z Himiko. Wir sind jedoch überrascht zu sehen, dass sich Himiko nicht in der Mitte des Protoklusters befand, sondern am Rand von 500 Millionen Lichtjahren vom Zentrum entfernt. """""""" sagte Masami Ouchi, ein Teammitglied des Nationalen Astronomischen Observatoriums Japans und der Universität Tokio, das Himiko 2009 entdeckte. Ironischerweise soll die mythologische Königin Himiko auch im Kloster von ihrem Volk gelebt haben. Ouchi fährt fort: """"""""Es ist immer noch nicht klar, warum Himiko nicht im Zentrum liegt. Diese Ergebnisse werden ein Schlüssel zum Verständnis der Beziehung zwischen Clustern und massiven Galaxien sein."""""""""

{accuracy: 13.0, fluency: 20.7, accuracy_M: -132.0, fluency_M: -2688.0, log p($\mathbf{y}|\mathbf{x}$): -127.0}

""""""""Interessanterweise war eine der 12 Galaxien in z66OD ein riesiges Objekt mit einem riesigen Gaskörper, bekannt als Himiko, das zuvor vom Subaru-Teleskop im Jahr 2009 gefunden wurde. """"""""Es ist vernünftig, einen Protocluster in der Nähe eines massiven Objekts wie Himiko zu finden. Wir sind jedoch überrascht zu sehen, dass Himiko nicht im Zentrum des Protoclusters lag, sondern am Rande 500 Millionen Lichtjahre vom Zentrum entfernt"""""""""", sagte Masami Ouchi, Teammitglied am Nationalen Astronomischen Observatorium Japans und der Universität Tokio, der Himiko 2009 entdeckte. Ironischerweise soll auch die mythologische Königin Himiko von ihrem Volk abgeschottet gelebt haben. Ouchi fährt fort: """"""""Es ist immer noch nicht klar, warum Himiko nicht in der Mitte liegt. Diese Ergebnisse werden ein Schlüssel zum Verständnis der Beziehung zwischen Clustern und massiven Galaxien sein."""""""""

{accuracy: 16.0, fluency: 21.3, accuracy_M: -122.5, fluency_M: -2624.0, log p($\mathbf{y}|\mathbf{x}$): -111.0}

Table 4: Translations of """"Interestingly, one of the 12 galaxies in z66OD was a giant object with a huge body of gas, known as Himiko, which was found previously by the Subaru Telescope in 2009. """"""""It is reasonable to find a protocluster near a massive object, such as Himiko. However, we're surprised to see that Himiko was located not in the center of the protocluster, but on the edge 500 million light-years away from the center."""" """" said Masami Ouchi, a team member at the National Astronomical Observatory of Japan and the University of Tokyo, who discovered Himiko in 2009. Ironically, the mythological queen Himiko is also said to have lived cloistered away from her people. Ouchi continues, """"""""It is still not understood why Himiko is not located in the center. These results will be a key for understanding the relationship between clusters and massive galaxies."""" """"""""

These translations are plotted in blue in Figure 1.