

# Does Summary Evaluation Survive Translation to Other Languages?

Anonymous ACL submission

## Abstract

The creation of a quality summarization dataset is an expensive, time-consuming effort, requiring the production and evaluation of summaries by both trained humans and machines. The returns to such an effort would increase significantly if the dataset could be used in additional languages without repeating human annotations. To investigate how much we can trust machine translation of summarization datasets, we translate the English *SummEval* dataset to seven languages and compare performances across automatic evaluation measures. We explore equivalence testing as the appropriate statistical paradigm for evaluating correlations between human and automated scoring of summaries. We also consider the effect of translation on the relative performance between measures. We find some potential for dataset reuse in languages similar to the source and along particular dimensions of summary quality.

## 1 Introduction

A large summarization dataset includes thousands of texts and human-written summaries (for example, CNN/Daily Mail (Hermann et al., 2015)). In order to make it applicable for wider research, it may also contain machine-generated summaries by many models, accompanied by human and machine evaluations of the quality of the generated summaries (Fabbri et al., 2021). The human annotation alone is a complicated effort, requiring careful planning and setup (Kryscinski et al., 2020; Tang et al., 2021; Iskender et al., 2021).

What purpose do the human annotations serve? Their main utility is serving as a benchmark for automated evaluation measures. Researchers design measures to closely approximate human judgment in order to increase the pace of summarization model improvement. As summarization resources grow for English-language models, it becomes increasingly important to consider whether we can

repurpose these datasets for use in other languages as well.

Given a method that could produce flawless translations, the original human annotations quite clearly remain useful, as the relative rankings of the summaries would be invariant. In this scenario, comparing automated measures in another language with the English human scores produces valid conclusions.

In reality, translation will introduce some distortions - both mild and extreme - that can spoil the utility of the original annotations. While a "uniform" distortion over all texts would preserve the relations among evaluations measures, this too is an unrealistic assumption as translation will correct and simplify some texts, introduce errors into others, and push components of text quality like relevance, coherence, and fluency in different directions (Fomicheva et al., 2021; Freitag et al., 2021). We are left to ask how to determine whether it is still practical to rely on the original human annotations for at least some quality measures and alternate languages?

In this paper, we seek to address this question through two quantitative explorations of automated evaluation measures under translation. First, we determine how often the correlation between a given measure and the original human annotations remains equivalent under translation. Second, we consider if one automated measure aligns more closely with human judgment than another in English, how often their relative positions are maintained after the translation. We conduct this investigation using the SummEval dataset (Fabbri et al., 2021), the largest corpus of English-language human annotated text summaries widely available. We translate this dataset from English to seven languages and evaluate the correlations between automated summary evaluation measures and human annotations. Using equivalence tests, we show that some aspects of summary quality ranking are pre-

served under translation for languages with similar alphabets and grammars to English. While we find some reasons for optimism about the potential for dataset reuse, our work clearly demonstrates that more research is needed to make translated datasets useful for a diverse set of languages.

## 2 Data and Models

We focus our analysis on the portion of SummEval<sup>1</sup> that includes human annotations. It consists of 100 texts, each accompanied by 11 human-written reference summaries and 17 machine-generated summaries produced by different models. Each machine-generated summary is annotated by three experts and five crowd workers using a 5-point scale for four quality measures: coherence, consistency, fluency, and relevance. For simplicity, we create a composite rating by averaging the expert scores for each quality of a given text-summary pair.

We translate all 100 source texts, 1100 human reference summaries, and 1700 machine-generated summaries into seven languages, French, German, Italian, Spanish, Afrikaans, Hindi, and Russian, using translation models trained and uploaded to the Hugging Face Model Hub by Helsinki-NLP<sup>2</sup> and accessed via the transformers library (Wolf et al., 2020). The specific models used for translation are named ‘opus-mt-L1-L2’, where one of L1 or L2 is ‘en’ (English), and the other is one of the languages ‘af’, ‘de’, ‘es’, ‘fr’, ‘hi’, ‘it’, or ‘ru’.

We used ‘bert-base-multilingual-cased’ as the underlying model for BLANC and ESTIME. While other choices of underlying model could produce higher correlations with human annotations in English, this multilingual model was selected to provide a more uniform performance across languages. BERTScore relies on ‘bert-base-multilingual-cased’ for all languages except English, for which it uses the model ‘roberta-large’<sup>3</sup>. ESTIME embeddings were taken from the 10th transformer block layer instead of the final 12th layer. We followed Vasilyev and Bohannon (2021), where it was shown that for the larger model ‘bert-large-uncased-whole-word-masking’ the 21st layer delivers the better performance than the 24th and final layer.

In each language version of the dataset, we

<sup>1</sup><https://github.com/Yale-LILY/SummEval>

<sup>2</sup><https://huggingface.co/Helsinki-NLP>

<sup>3</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

score machine-generated summaries with a few common or promising automated evaluation measures that could be applied to all eight languages. We calculate the following truly automated (not needing human written reference summaries) measures: Jensen-Shannon (Louis and Nenkova, 2009), ESTIME (Vasilyev and Bohannon, 2021)<sup>4</sup> and BLANC (Vasilyev et al., 2020)<sup>5</sup>. We also calculate the following reference-based automatic evaluation measures: BLEU (Papineni et al., 2002), BERTScore-F1<sup>6</sup> (Zhang et al., 2020), and ROUGE (Lin, 2004) as ROUGE-1,2,L<sup>7</sup>. We use the same original human annotations provided by the SummEval dataset as annotations in each of the seven translated languages.

We calculate correlations between automated evaluation measures in each language and the human annotations on the original English dataset. We seek to answer whether these correlations are reasonably independent of the language. In other words, can we rely on such correlations to provide consistent judgement of evaluation measures in other languages?

## 3 Comparisons within Measures

### 3.1 Simple Correlations

It has become standard in the summarization literature to judge the performance of an automated measure by the correlation of its scores with human evaluation of summaries (e.g. Zhang et al. (2020), Deutsch et al. (2021)). Figure 1 shows Spearman’s  $\rho$  and Kendall’s  $\tau$  correlation coefficients between the expert human evaluations and the automated measures run on the English summaries found in the SummEval dataset.

The correlations are consistently weak, indicating that the measures rely on different features than human evaluations of a summary. ESTIME, BERTScore, and Jensen-Shannon all demonstrate somewhat higher correlations in at least some measures of quality, perhaps reflecting a more nuanced approach to summary scoring.

Automated evaluation of summarization models is still an evolving field. While most measures disagree with human judgment often, they are still widely used as points of comparison across model

<sup>4</sup><https://github.com/PrimerAI/blanc/tree/master/estime>

<sup>5</sup><https://github.com/PrimerAI/blanc>

<sup>6</sup>[https://github.com/Tiiiger/bert\\_score](https://github.com/Tiiiger/bert_score)

<sup>7</sup><https://github.com/google-research/google-research/tree/master/rouge>

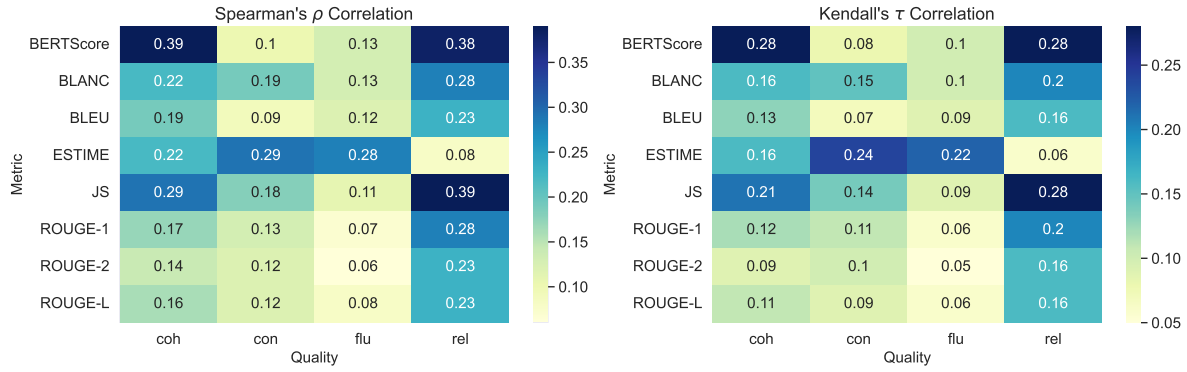


Figure 1: Spearman's  $\rho$  and Kendall's  $\tau$  correlations of expert human scores (coherence, consistency, fluency, relevance) with automated evaluation measures for the original English summaries. Note: JS (Jensen-Shannon) and ESTIME correlations are negated.

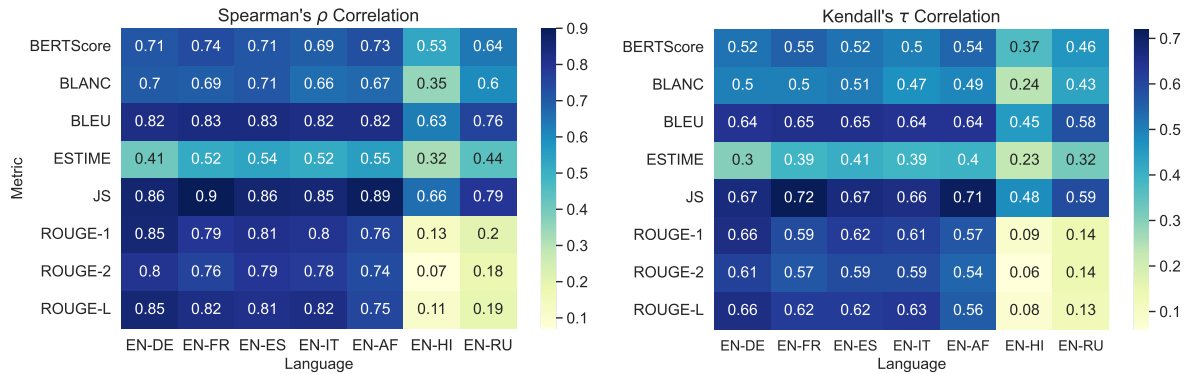


Figure 2: Spearman's  $\rho$  and Kendall's  $\tau$  correlations between automated evaluation measures in English and in translated languages German (DE), French (FR), Spanish (ES), Italian (IT), Afrikaans (AF), Hindi (HI), and Russian (RU).

174 outputs. Therefore, it remains highly relevant to  
 175 determine whether translation preserves the judgments  
 176 rendered by the automated measures.

177 We may consider an evaluation measure to be  
 178 useful under translation if the scores it assigns to  
 179 summaries are consistent across languages, perhaps  
 180 in absolute value but at least in the rank ordering  
 181 of summaries. Therefore such a measure would exhibit  
 182 high correlation between its values on English  
 183 summaries and those for the summaries translated  
 184 to other languages. Figure 2 shows Spearman's  $\rho$   
 185 and Kendall's  $\tau$  correlation coefficients between  
 186 the automated measures run on the English corpus  
 187 and each translated corpus.

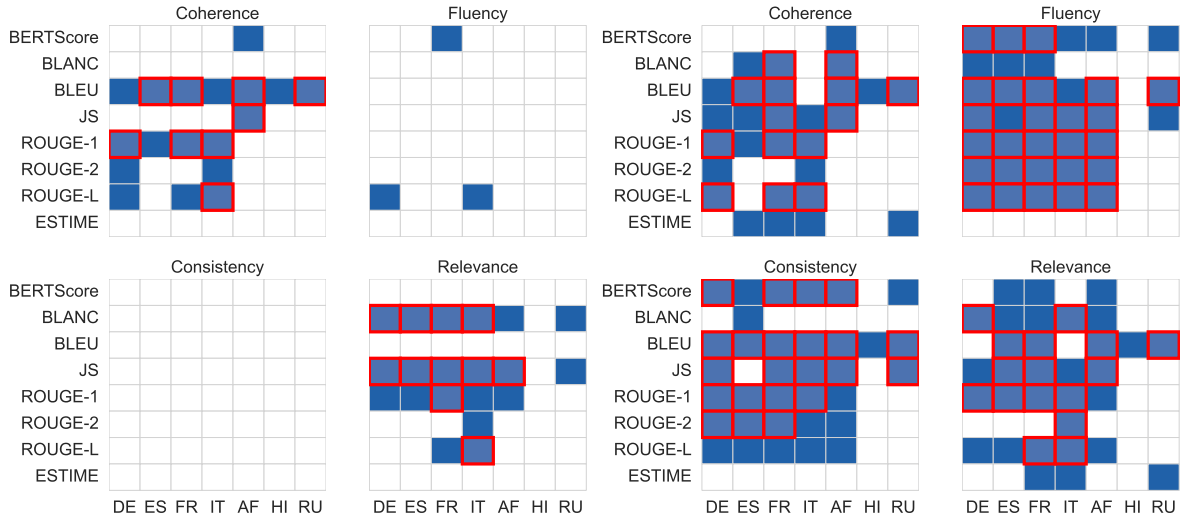
188 For a given measure, the correlations across lan-  
 189 guages are generally much stronger than those be-  
 190 tween automated measures and human evaluations  
 191 in English seen in Figure 1. For languages with  
 192 the strongest correlations to the English measures,  
 193 this result provides some promise that translation  
 194 might introduce minimal additional noise, meaning

195 the evaluation measure provides consistent signal  
 196 across languages.

197 The reference-based measures generally show  
 198 stronger correlations ( $\rho > 0.6$ ,  $\tau > 0.5$ ) between  
 199 English and German, French, Spanish, Italian, and  
 200 Afrikaans translations. For Russian and Hindi,  
 201 they show weaker correlations, drastically so for  
 202 ROUGE measures. Among the reference-free mea-  
 203 sures, Jensen-Shannon and BLANC demonstrate  
 204 similar patterns of performance. These results at  
 205 least suggest that measures may prove useful when  
 206 translating datasets to languages with similar ori-  
 207 gins (here Italic or Germanic languages). However,  
 208 ESTIME shows weak correlations across languages  
 209 with a smaller drop in correlation between Western  
 210 European derived languages and Hindi and Rus-  
 211 sian.

### 3.2 Significance Tests 212

213 Given the promising results in Section 3.1, we  
 214 seek to test whether correlations between an au-



(a) TOST with standard deviation margin of equivalence

(b) TOST with constant 0.05 margin of equivalence

Figure 3: Results of tests of equivalence for each automated measure (y-axis), language (x-axis), and quality measure (coherence, consistency, fluency, relevance). Blue squares indicate p-value  $\leq 0.05$  while red highlights indicate the result remained significant after applying Benjamini-Yekutieli correction for FDR control. *Left*: Results for TOST with standard deviation margin of equivalence. *Right*: Results for TOST with constant 0.05 margin of equivalence.

tomated measure and the original expert scores are statistically invariant when run on the English and translated summaries. Since human evaluations are split into four qualities - coherence, consistency, fluency, relevance - we consider correlations separately along each measure. For example, we look to answer whether the correlation between English BLANC scores and English expert scores for relevance is equivalent to the correlation between German BLANC scores and English expert scores for relevance. We consider this a natural test of an automated measure’s utility after translation, as we hope measures will reflect human judgment in a consistent and predictable manner across languages.

Since we are interested in demonstrating a lack of statistical difference between two correlations,  $\rho_1$  and  $\rho_2$ , we cannot use a typical hypothesis test with null hypothesis  $H_0 : \rho_1 = \rho_2$ . Such a test would only suggest equivalence by failing to reject the null hypothesis, which could simply occur due to a lack of statistical power.

Instead, we turn to equivalence tests, a paradigm which effectively reverses null and alternative hypotheses, ie.  $H_0 : \rho_1 \neq \rho_2$ . We explore two such tests, Two One-Sided Tests (TOST) and Anderson-Hauck tests, and call for additional research to standardize their use for summarization evaluation.

### 3.3 Two One-Sided Tests (TOST)

In the TOST procedure (Schuirmann, 1987), we must set a margin of equivalence,  $\Delta_E$ , within which we consider two test statistics to be equivalent. Then for two correlations,  $\rho_1$  and  $\rho_2$ , we have null and alternative hypotheses:

$$H_0 : \rho_1 - \rho_2 < -\Delta_E \text{ or } \rho_1 - \rho_2 > \Delta_E$$

$$H_1 : -\Delta_E < \rho_1 - \rho_2 < \Delta_E$$

While in a field like medicine, the margin might be well defined by a chemical process, we lack a strong prior for choosing a relevant margin. We explore several options and consider the sensitivity of p-values to our choices when evaluating the validity of the tests’ conclusions.

The Kendall rank correlation differences considered do not follow a normal distribution, and we use bootstrap resampling (Efron and Tibshirani, 1993) to generate an empirical distribution. For a given translation language, automated evaluation measure, and quality measure, we sample across (text, summary, and reference summary) tuples. (Note for reference-based summaries - BERTScore, BLEU, and ROUGE - a more complete bootstrap procedure would account for the stochasticity present in the choice of reference summaries themselves. We provide an illustrative example in Appendix B.)

While permutation-based tests have been shown to have higher power in summarization evaluation than bootstrap resampling (Deutsch et al., 2021), permutation tests assume null hypothesis  $H_0 : \rho_1 = \rho_2$  and are not simply adapted to our case. We apply a multiple testing correction to the p-values calculated due to the large number of tests considered. We use the Benjamini-Yekutieli procedure (Benjamini and Yekutieli, 2001) to account for dependence among correlation measures and control the false discovery rate (FDR) at level  $\alpha = 0.05$ .

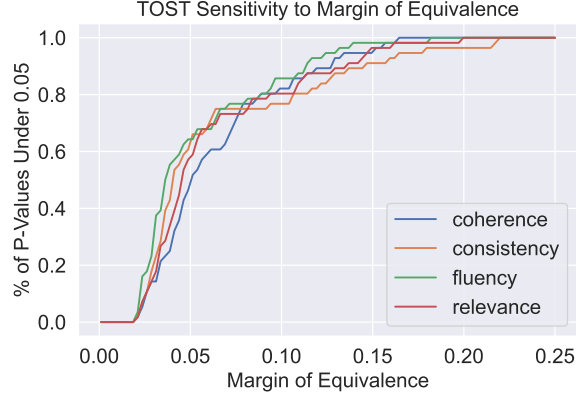
We consider several relevant equivalence margins with different trade-offs. We try a "constant margin" of 0.05 across all measures and qualities; a "standard deviation margin" using the standard deviations for correlations between individual experts and an automated measure; and a "maximum difference margin" calculated as the largest absolute difference in correlations between individual experts and an automated measure. Under the constant margin, 56% of correlations are equivalent before FDR correction and 31% after. Under the max difference margin, 42% of correlations are equivalent before correction and 28% after. Finally, under the standard deviation margin, 17% of tests are equivalent before and 8% after correction.

We present the full results of the TOST procedure with a standard deviation margin in the left panel and a constant margin in the right panel of Figure 3. While both panels demonstrate interesting patterns of equivalence, we focus on the standard deviation margin as it is tailored to each language-measure pair, relies on a less arbitrary value of expected variation under equivalence, and is more conservative than the other margins considered. The max difference and constant margins found much higher rates of equivalence under translation.

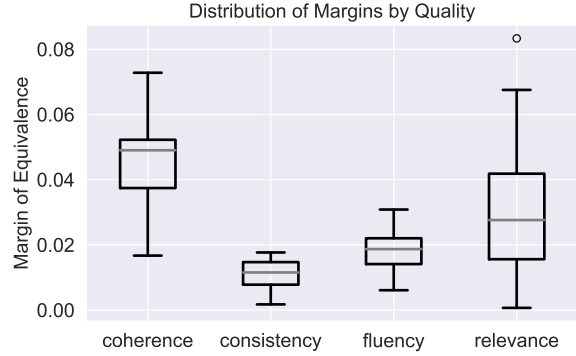
Examining the results, we can note a few clear patterns. First, as seen under the simple correlation analysis, the Italic and Germanic languages have a higher number of significant results than Hindi or Russian. We may still consider using translated summarization datasets from English to languages considered "close." However, there are few significant results in the fluency or consistency qualities. Therefore the automated measures may only be useful under translation along specific dimensions of quality. Looking at the correlations in English between automated measures and expert judgments

in Figure 1, fluency and consistency also tend to have much lower correlations than coherence and relevance.

Additionally, the choice of equivalence margin has a consequential impact on results. Figure 4a shows how the number of significant p-values changes in response to an increasing margin of equivalence. Given the apparent sensitivity to changes in the margin, further research is warranted into how the performance of translation and summarization systems relates to the correlations measured here.



(a) Sensitivity of TOST to the margin of equivalence. Small changes in the margin can result in a large change in the percent of tests with significant p-values.



(b) Distribution of margins by quality under the standard deviation margin. Using the variation observed among individual human annotations produces more strict thresholds of equivalence for consistency and fluency and more lenient ones for coherence and relevance.

Figure 4: Measuring the impact of margins of equivalence on the TOST results.

Therefore, the lack of significance for the fluency and consistency qualities can be attributed to both the capabilities of the automated measures and how the standard deviation margin varies across qualities. We already expect from Figure 1 that measures may be capturing a large amount of noise for

fluency and consistency and would fare poorly under translation, resulting in fewer equivalent results. However, the amount of inter-rater disagreement also plays a significant role in determining equivalence by expanding or contracting the margins. Figure 4b highlights the differences in standard deviation margins for each quality across automated measures. Consistency and fluency had smaller margins with tighter distributions, with median margins of 0.012 and 0.019 and inter-quartile ranges (IQRs) of 0.007 and 0.008 respectively. By contrast, coherence and relevance had median margins 0.049 and 0.028 with IQRs 0.015 and 0.026 respectively. Thus human annotators showed stronger agreement on consistency and fluency, presenting a higher threshold for equivalence after translation.

### 3.4 Anderson-Hauck Tests

While TOST provides a non-parametric route towards equivalence testing, we consider an additional parametric test that may improve statistical power. The Anderson-Hauck test is an equivalence testing procedure for dependent correlation coefficients which uses an approximate non-central t-distribution to calculate p-values (Anderson and Hauck, 1983). Prior comparisons with TOST demonstrated that Anderson-Hauck can trade some additional Type-I error for higher power (Counsell and Cribbie, 2015).

We consider the same margins of equivalence and apply Benjamini-Yekutieli for FDR control at level  $\alpha = 0.05$ . A similar pattern emerges when considering results under different margins, and under the standard deviation margin we reject the null hypothesis in under 1% of tests.

The pattern of equivalence is largely the same as that found under TOST but with greater sparsity of significant results. Ultimately while the tests hint towards the ability to reuse summarization datasets in similar languages to English, we are only able to detect equivalence in a minority of cases. Our analysis relies predominantly on the TOST results since it does not rely on distributional assumptions for the differences in correlations and has a more robust literature to follow.

## 4 Comparisons between Measures

While our statistical tests focus on the absolute correlation between automated and human scores, we can instead consider the automated measures relative to one another. If one measure correlates better

than another with human scores in the original English dataset, would it still be better in a translated (non-English) dataset? Additionally, we can return the dataset back to English to get a sense of the distortion introduced by the translation process.

To estimate the consistency with which one measure dominates another, we turn to bootstrap resampling of the summary evaluations. We select 10,000 bootstrap samples from the 1700 text-summary-references tuples. Let  $P$  represent the fraction of samples in which one measure is better than another for a given measure-measure pair; we consider a pair "resolved" if one measure outperforms another in at least 97.5% of all the resamplings, ie.  $P \geq 0.975$  in the original English dataset. Using Kendall rank correlations, the number of resolved measure-measure pairs is 64% for relevance, 61% for coherence, 56% for consistency, and 42% for fluency. With a baseline reading of how stable the measure rankings are in English, we can ask what happens with these resolved pairs when the dataset is translated.

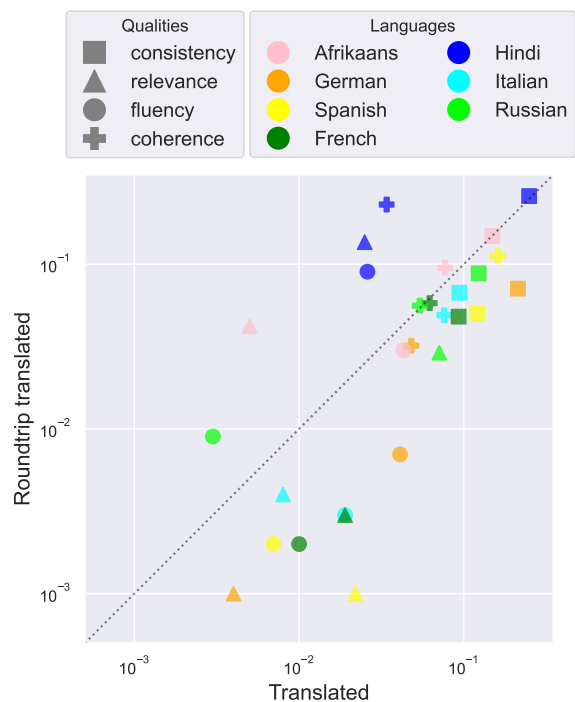


Figure 5: Result of bootstrapping: average shift in probability  $P$  of one measure being better than another, when the evaluation data are translated to another language (x-axis) and then translated back to English (y-axis). The average is taken over all measure-measure pairs that had  $P \geq 0.975$  in English.

For most languages and qualities the shift of  $P$  is less than 0.1, the largest is 0.25 (consistency,

Hindi). Many resolved measure-measure pairs become unresolved after translation, though no shift is drastic enough to reverse which measure ranks higher in a majority of samples (i.e. crossing  $P = 0.5$ ). Figure 5 suggests that in most cases our conclusion about comparing two measures will not change with translation.

Along its x-axis, Figure 5 shows how much on average the fraction  $P$  changes (increases or decreases) after translation for resolved measure-measure pairs, where the average is over a given language and quality measure.

A round-trip translation returns each summary to its source language, allowing us to more effectively isolate the effect of translation quality on the consistency of automated measures. The dashed line  $y = x$  seen in Figure 5 represents points where the round-trip translation causes an equally-sized shift as the forward translation. We note that the observed shifts are mostly under the diagonal - the shifts caused by translation are to some degree reversed when we return to English.

While the shifts for round-trip translations are on average smaller than for one-way, they demonstrate that translation is far from perfect and introduces enough noise to be detected by the summarization evaluation measures. Notably, the points above the diagonal come from Hindi, Russian and Afrikaans round-trip translation. This confirms our intuition that a translation to languages more distant from English is more risky for the survival of the summary evaluation. We hope further research may reveal additional ways to use the round-trip translation for the criteria of survival.

## 5 Discussion

The results presented significant differences among automated summarization measures and their relationships to the four quality measures. We seek to build an intuition for these findings and make use of qualitative exploration to ground our understanding

We can review the scores for the 1700 summaries in reduced dimensions using principal components analysis (PCA). Figure 6 shows each 1700-dimensional vector projected onto the first two principal components, which collectively explain 38.5% of the variance. There are four vectors of human expert scores, corresponding to the quality measures coherence, consistency, fluency, and relevance, averaged over the three individual experts. Each automated measure (for example, ROUGE-2)

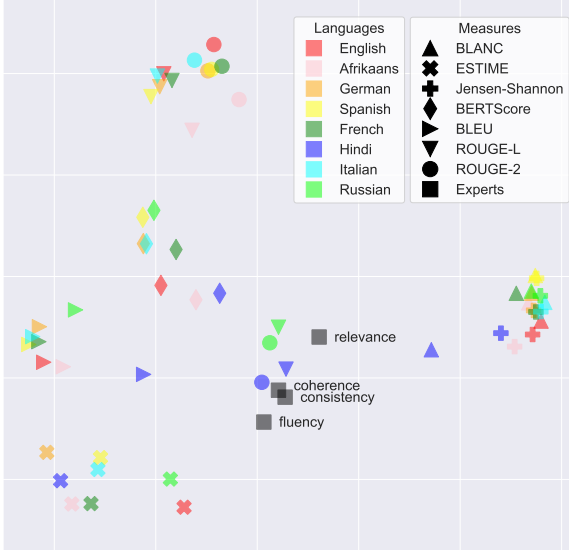


Figure 6: PCA plot of summary quality scores. All scores were transformed to ranks before PCA, to reduce subjectivity of the respective scales. Note the human expert scores in black squares exist for the English dataset only.

produced eight 1700-dimensional vectors, one for each language.

PCA can be used to disentangle the sources of divergence among evaluation measures under translation. The plot helps highlight the relative strength of translation over the summarization evaluation methods themselves. If machine translation added significant noise to the summaries, we would expect the relative position of language-specific scores in Figure 6 to be inconsistent across evaluation measures. Instead, we generally observe tight clusters for each evaluation measure with shared relative positions among the languages (at least when ignoring Hindi and Russian).

This pattern reflects the "stability" of evaluation measures undergoing translation found in Section 4. The PCA recasts translation as a shift in geometric space; across measures, the location occupied by each language is a similar vector shift from its corresponding English point. The exercise in round-trip translation is an indicator of reversibility for this geometric shift. The qualities and languages that occupy the bottom of Figure 5 are most unchanged by the translation process. On the other hand, measures like ESTIME that break this pattern highlight the non-uniformity of the distortion introduced by translation and indicate that it may be more prudent to rely on measures where the distortion is consistent and predictable.

This closer look at the effects of translation also helps disentangle the sources of noise that degraded the correlations studied in Section 3. A measure like ESTIME shows strong correlation with the human evaluations of consistency and fluency in English, but its unusual response to translation is a strong explanatory factor for why its relationships to human annotations were not found to be equivalent in other languages. Consistency also tends to show larger shifts in measure-measure pair rankings in Figure 5, adding another reason that translation would cause greater degradation to ESTIME’s performance. Similarly, among the Germanic and Italic languages, relevance and fluency appear to be least affected by translation. Any lack of equivalence found for these qualities is then more likely to be caused by the abilities of the automated measures rather than the caliber of translation. Comparisons within and between measures can serve as a guide for how much to trust an automated measure under translation and where sources of noise may arise.

We note a few curious observations from Figure 6 in Appendix A.

## 6 Conclusion

In this paper, we probed how well automated evaluations of summaries remain consistent on texts translated to other languages. We focused on the SummEval dataset and considered its translation to French, German, Italian, Spanish, Afrikaans, Hindi, and Russian.

To answer whether English human annotations can be trusted in other languages, at least for specific qualities, we explored tests of equivalence as a gauge of consistency after translation. We found that translation can preserve correlations of evaluation metrics with the English human scores for coherence or relevance but could not conclude the same for fluency or consistency.

A complete answer to our query is a challenging task, since moving to another language affects not only the dataset, but also the measures themselves. While definitely proving that the original human annotations cannot be reused is likely impossible, our results suggest that there are clear differences in performance based on the choice of target language, automated measure, and notion of quality.

We call for additional research into summary evaluation metrics that can survive translation, as it offers a relatively simple path towards extending

NLP capabilities for lower resource languages. Future work could identify how changes in the margin of equivalence equate to deterioration of model performance. Additionally, this line of research could be extended to a larger selection of languages and automated evaluation measures.

## References

- Sharon Anderson and Walter W. Hauck. 1983. [A new procedure for testing equivalence in comparative bioavailability and other clinical trials](#). *Communications in Statistics - Theory and Methods*, 12(23):2663–2692.
- Yoav Benjamini and Daniel Yekutieli. 2001. [The control of the false discovery rate in multiple testing under dependency](#). *The Annals of Statistics*, 29(4):1165 – 1188.
- Alyssa Counsell and Robert A. Cribbie. 2015. [Equivalence tests for comparing correlation and regression coefficients](#). *British Journal of Mathematical and Statistical Psychology*, 68(2):292–309.
- 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.
- Bradley Efron and Robert J. Tibshirani. 1993. *An Introduction to the Bootstrap*. Number 57 in Monographs on Statistics and Applied Probability. Chapman & Hall/CRC, Boca Raton, Florida, USA.
- Alexander R. Fabbri, Wojciech Kryściński, Bryan McCann, Caiming Xiong, Richard Socher, and Dragomir Radev. 2021. [SummEval: Re-evaluating summarization evaluation](#). *Transactions of the Association for Computational Linguistics*, 9:391–409.
- Marina Fomicheva, Lucia Specia, and Nikolaos Aletras. 2021. [Translation error detection as rationale extraction](#).
- Markus Freitag, George Foster, David Grangier, Viresh Ratnakar, Qijun Tan, and Wolfgang Macherey. 2021. [Experts, errors, and context: A large-scale study of human evaluation for machine translation](#).
- Karl Moritz Hermann, Tomáš Kočiský, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. 2015. [Teaching machines to read and comprehend](#). In *Advances in Neural Information Processing Systems 28*, pages 1693–1701. Curran Associates, Inc.
- Neslihan Iskender, Tim Polzehl, and Sebastian Möller. 2021. [Reliability of human evaluation for text summarization: Lessons learned and challenges ahead](#). In *Proceedings of the Workshop on Human Evaluation of NLP Systems (HumEval)*, pages 86–96. Association for Computational Linguistics (2021).



594	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. <a href="#">Evaluating the factual consistency of abstractive text summarization</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing</i> , pages 9332–9346. Association for Computational Linguistics.	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. <a href="#">BERTScore: Evaluating text generation with bert</a> . <i>arXiv</i> , arXiv:1904.09675v3.	650 651 652 653
600	Chin-Yew Lin. 2004. <a href="#">ROUGE: A package for automatic evaluation of summaries</a> . In <i>Proceedings of Workshop on Text Summarization Branches Out</i> , pages 74–81. Association for Computational Linguistics.	<b>A Observations from PCA</b>	654
604	Annie Louis and Ani Nenkova. 2009. <a href="#">Automatically evaluating content selection in summarization without human models</a> . In <i>Proceedings of the 2009 Conference on Empirical Methods in Natural Language Processing</i> , pages 306–314. Association for Computational Linguistics.	The locations of the measures in Figure 6 after translation largely remain close to the original English version, except Hindi and Russian points. The locations show interesting patterns. The reference-based measures, based on hard (ROUGE, BLEU) or soft (BERTScore) overlap of tokens between the summary and the human-written reference summaries, are in the same top left quadrant with respect to the human scores. The reference-free measures BLANC and Jensen-Shannon are on the opposite side. It is natural for both BLANC and Jensen-Shannon to be on the relevance side of the human scores: BLANC estimates how well a text can be reconstructed from its summary, and Jensen-Shannon considers the Kullback–Leibler divergence between the summary and the text. For ESTIME, however, as for a consistency-oriented measure, it makes sense to be on the consistency and fluency side of the human scores, rather than on the relevance side.	655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674
610	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. <a href="#">BLEU: a method for automatic evaluation of machine translation</a> . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics (ACL)</i> , pages 311–318, Philadelphia. Association for Computational Linguistics.	For most measures, the translated scores are often closer to the expert evaluations than the English scores. Strangely, it is especially true for Hindi and, in the case of ROUGE, for Russian. One possible explanation is that the translation simplifies the phrases and the choice of words, thus making it easier for some evaluation measures, at least along some dimensions. The pattern associated with ESTIME is distinct from other measures: the non-English scores for ESTIME are almost always further away from the human scores. This suggests that maybe ESTIME is sensitive enough to require a higher quality translation. We cannot blame the underlying multilingual model, because both BLANC and BERTScore use the same model.	675 676 677 678 679 680 681 682 683 684 685 686 687 688 689
616	Donald J. Schuirmann. 1987. A comparison of the two one-sided tests procedure and the power approach for assessing the equivalence of average bioavailability. <i>Journal of Pharmacokinetics and Biopharmaceutics</i> , 15:657–680.	<b>B Bootstrap with Reference-Summaries</b>	690
621	Xiangru Tang, Alexander R. Fabbri, Ziming Mao, Griffin Adams, Borui Wang, Haoran Li, Yashar Mehdad, and Dragomir Radev. 2021. <a href="#">Investigating crowdsourcing protocols for evaluating the factual consistency of summaries</a> . <i>arXiv</i> , arXiv:2109.09195.	Throughout the paper we used bootstrapping with resampling of the (text, summary, references) tuples, where the references are the reference summaries needed by some measures (BERTScore, BLEU, ROUGE). For each text in SummEval (Fabbri et al., 2021), there are 11 reference summaries, and a full bootstrap for the reference-based measures should also include a resampling of the reference summaries themselves.	691 692 693 694 695 696 697 698 699
626	Oleg Vasilyev and John Bohannon. 2021. <a href="#">Estime: Estimation of summary-to-text inconsistency by mismatched embeddings</a> . In <i>Proceedings of the 2nd Workshop on Evaluation and Comparison of NLP Systems</i> , pages 94–103. Association for Computational Linguistics.		
632	Oleg Vasilyev, Vedant Dharnidharka, and John Bohannon. 2020. <a href="#">Fill in the BLANC: Human-free quality estimation of document summaries</a> . In <i>Proceedings of the First Workshop on Evaluation and Comparison of NLP Systems</i> , pages 11–20. Association for Computational Linguistics.		
638	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. <a href="#">Transformers: State-of-the-art natural language processing</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45. Association for Computational Linguistics (2020).		

700 The impact of this added source of randomness  
 701 can be seen by constructing confidence intervals  
 702 for the estimated correlation between an evaluation  
 703 measure and human scores. When we add resam-  
 704 pling over reference summaries, confidence inter-  
 705 vals widen and require more time and resources  
 706 to compute. In Table 1 we illustrate the widen-  
 707 ing of the confidence interval on an example using  
 708 BERTScore correlations with SummEval human  
 709 expert scores (in the original English SummEval  
 710 dataset). We ran 500K reference summaries resam-  
 711 plings, recomputing scores and correlations. The  
 712 BERTScore is a peculiar and convenient case for  
 713 bootstrap resampling of reference summaries, be-  
 714 cause the score is defined as a max score over the  
 715 scores taken individually for each reference sum-  
 716 mary (Zhang et al., 2020).

	Kendall's $\tau$			Spearman's $\rho$		
	low	high	widen	low	high	widen
<b>coherence</b>	0.245	0.307	0.011	0.345	0.428	0.015
<b>consistency</b>	0.041	0.117	0.002	0.052	0.148	0.003
<b>fluency</b>	0.062	0.135	0.004	0.080	0.175	0.006
<b>relevance</b>	0.246	0.310	0.026	0.338	0.424	0.035

Table 1: The columns 'low' and 'high' are the confidence boundaries from bootstrap without resampling reference summaries, for BERTScore correlations with expert human scores (coherence, consistency, fluency, relevance). The column 'widen' is the widening of the confidence interval as a result of adding the resampling of the reference summaries to the bootstrap resampling. Kendall's Tau correlation is Tau-c. The confidence boundaries are for 0.025 and 0.975 percentiles. The bootstrapping used 500K resamplings.

717 The low and high correlation values are given in  
 718 the table for bootstrap without resampling of ref-  
 719 erence summaries, as corresponding to 0.025 and  
 720 0.975 percentiles of the distribution. The 'widen'  
 721 column in the table shows how much the confi-  
 722 dence interval ('high minus low') changed after in-  
 723 cluding resampling of the 11 reference summaries  
 724 into the bootstrapping. Some quality measures  
 725 are especially affected by the change, with confi-  
 726 dence intervals for Kendall correlation widening  
 727 by 40% for relevance and by 17% for coherence  
 728 (for Spearman's correlations, correspondingly, 42%  
 729 and 18%). Notice that the relevance and coherence  
 730 are exactly the qualities in which BERTScore is  
 731 reported as a strong measure (Vasilyev and Bohan-  
 732 non, 2021).