

CrossSum: Beyond English-Centric Cross-Lingual Abstractive Text Summarization for 1500+ Language Pairs

Abhik Bhattacharjee¹, Tahmid Hasan¹, Wasi Uddin Ahmad²,
Yuan-Fang Li³, Yong-Bin Kang⁴, Rifat Shahriyar¹

Bangladesh University of Engineering and Technology (BUET)¹, University of California,
Los Angeles², Monash University³, Swinburne University of Technology⁴

{tahmidhasan, rifat}@cse.buet.ac.bd, abhik@ra.cse.buet.ac.bd

Abstract

We present CrossSum, a large-scale cross-lingual abstractive summarization dataset comprising 1.7 million article-summary samples in 1500+ language pairs. We create CrossSum by aligning identical articles written in different languages via cross-lingual retrieval from a multilingual summarization dataset. We propose a multi-stage data sampling algorithm to effectively train a cross-lingual summarization model capable of summarizing an article in any target language. We also propose LaSE, a new metric for automatically evaluating model-generated summaries and showing a strong correlation with ROUGE. Performance on ROUGE and LaSE indicate that pretrained models fine-tuned on CrossSum consistently outperform baseline models, even when the source and target language pairs are linguistically distant. To the best of our knowledge, CrossSum is the largest cross-lingual summarization dataset and the first-ever that does not rely solely on English as the pivot language. We are releasing the dataset, alignment and training scripts, and the models to spur future research on cross-lingual abstractive summarization. The resources can be found at <https://github.com/csebuetnlp/CrossSum>.

1 Introduction

Cross-lingual summarization is the task of generating a summary in a target language given a source text in another language. The task is challenging as it combines summarization and translation in one task, both challenging tasks in their own right. Earlier approaches to cross-lingual summarization thus employed pipeline methods like translate-then-summarize (Leuski et al., 2003) or summarize-then-translate (Wan et al., 2010). Not only computationally expensive, having to use multiple models, these approaches also suffer from error-propagation (Zhu et al., 2019) from one model to another, degrading the overall performance.

Input Article: [...] 新型コロナウイルスに対し、様々な既存の治療法の効果を試す世界的規模の臨床試験の一貫として、**デキサメタゾン**が試された。(Dexamethasone was tested as part of a global clinical trial to test the effectiveness of various existing therapies against the new **coronavirus**.) [...] その結果、人工呼吸器を必要とする**重症患者**の致死率が3割下がり。(As a result, the case fatality rate of **critically ill patients** who require a ventilator is reduced by 30%.) [...] ボリス・ジョンソン英首相は「イギリス**科学界**の素晴らしい成果」を歓迎し。(British Prime Minister Boris Johnson welcomed "the great achievements of the British **scientific community**".) [...] 「しかもこれは、**世界中で手に入る薬だ**」("And this is a **medicine available all over the world**".) [...] きわめて**安い**ステロイド剤だった (but a very **cheap** steroid that has been used for a long time.)

Summary: **বিজ্ঞানীরা** বলছেন **ডেক্সামেথাসোন** নামে **সস্তা** ও **সহজলভ** একটি **ঔষধ** **করোনাভাইরাসে** **গুরুতর** **অসুস্থ** **রোগীদের** **জীবন** **রক্ষা** করতে সাহায্য করবে। (**Scientists** say a **cheap** and **readily available drug** called **dexamethasone** will help save the lives of **critically ill patients** with **coronavirus**.)

Figure 1: A sample article-summary pair from CrossSum, the article is written in Japanese and the summary is in Bengali. We translate the texts in English for better understanding. Word and phrases of the article relevant to the summary are color-coded.

The success of sequence-to-sequence (seq2seq) models (Cho et al., 2014; Sutskever et al., 2014) and the advances in Transformer-based models (Vaswani et al., 2017; Rothe et al., 2020) have aided in the emergence of end-to-end methods that can produce cross-lingual summaries with one single model (Zhu et al., 2019). The availability of cross-lingual summarization datasets (Ladhak et al., 2020; Perez-Beltrachini and Lapata, 2021) has also helped this task gain popularity in recent times. However, these datasets cover only a few languages, have few samples for training and evaluation, or use English as the pivot language (i.e., the target language always remains English), thereby limiting the applicability to a great extent.

To democratize cross-lingual summarization beyond high-resource languages, in this work, we introduce **CrossSum**, a large-scale cross-lingual ab-

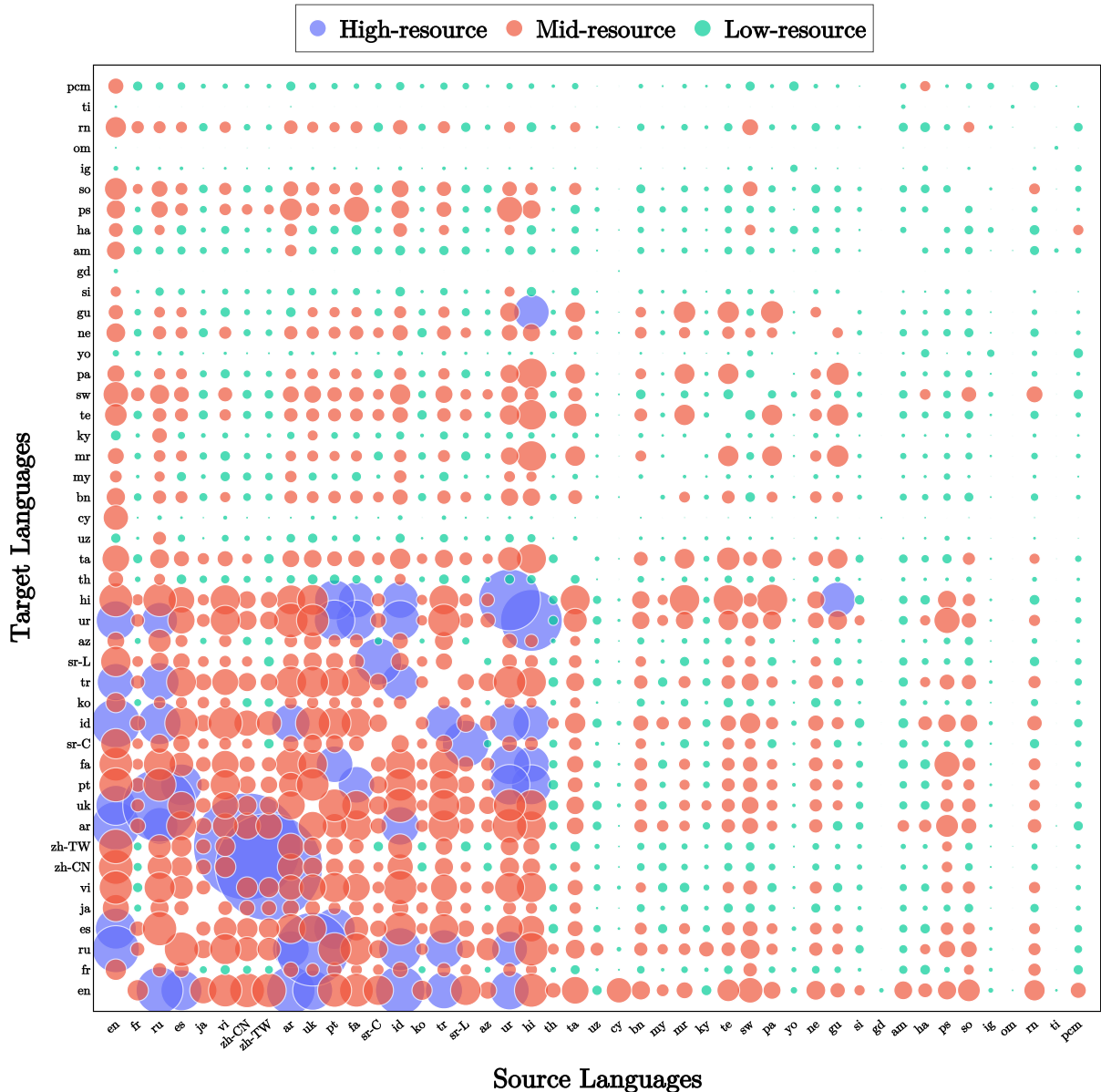


Figure 2: A bubble plot depicting the article-summary statistics of the CrossSum dataset. The radii of the bubbles are proportional to the number of article-summary pairs for the corresponding language pair. Languages in the axes are sorted by the number of their Wikipedia entries to show a sequential contrast from high- to low-resource languages. We consider a language pair as low-resource in CrossSum if the number of samples is below 500, mid-resource for 500 to less than 5000, and high-resource for pairs exceeding 5000.

stractive summarization dataset containing 1.7 million article-summary samples in 1500+ language pairs by aligning identical articles written in different languages via cross-lingual retrieval from the multilingual XL-Sum (Hasan et al., 2021) dataset covering 45 languages. We design a multistage sampling algorithm for successful training of multilingual models that can generate a summary in any target language for a source article in any language (i.e., a many-to-many summarization model). We also propose LaSE, an automatic metric for evaluating cross-lingual summaries when reference

summaries in the target language may not be available (but available in another language), potentially opening new doors to evaluate low-resource language pairs. We also show a strong correlation between ROUGE and LaSE, validating the reliability of LaSE. For the very first time, we perform cross-lingual summarization on a broad and diverse set of languages without relying on English as the standalone pivot language, consistently outperforming several many-to-one and one-to-many models, as well as summarize-then-translate baselines.

To the best of our knowledge, CrossSum is the

first publicly available cross-lingual summarization dataset for a large number of language pairs. We are releasing the dataset, alignment and training scripts, and models hoping that these resources will encourage the community to push the boundaries of cross-lingual abstractive summarization beyond the English and other high-resource languages.

2 The CrossSum Dataset

The idea of curating a cross-lingual summarization dataset is to pair the source text of an article A with the summary of another identical article B written in a different language and vice-versa, with the availability of a multilingual dataset where different languages have identical contents. Language-agnostic sentence representations (Artetxe and Schwenk, 2019a; Feng et al., 2022) have achieved state-of-the-art results in cross-lingual text mining (Zweigenbaum et al., 2017; Artetxe and Schwenk, 2019b), and therefore, provide a way to search identical contents across languages.

Two contemporary works have compiled large-scale multilingual summarization datasets, namely XL-Sum (Hasan et al., 2021) (1.35M samples in 45 languages) and MassiveSumm (Varab and Schluter, 2021) (28.8M samples in 92 languages). Though substantially larger than the other, MassiveSumm is not publicly available. Since public availability is crucial for promoting open research, we opted for the other alternative, XL-Sum, which is distributed under a non-commercial research license. XL-Sum has another benefit: all articles are crawled from a single source, BBC News. We observed that BBC publishes similar news contents in different languages and follow similar summarization strategies; hence it would increase the quality and quantity of the mined article-summary pairs.

For simplicity, we perform the similarity search over summaries only. To ensure maximum quality, we set two strong prerequisites for a summary S_A in language A to be paired with another summary S_B of language B:

1. S_B must be the nearest neighbor of S_A among all summaries in B, and vice-versa.
2. The similarity between S_A and S_B must be above the threshold, τ .

To measure similarity, we use the inner products of Language-agnostic BERT Sentence Representation (LaBSE) (Feng et al., 2022) (a unit vector for an input text sequence). We set the minimum similarity threshold as the average threshold

($\tau = 0.7437$) of all languages that maximized respective F_1 score for LaBSE in the BUCC mining tasks (Zweigenbaum et al., 2017).¹

Induced Pairs We noticed that many summaries, despite being nearest neighbors, were filtered out because of the threshold, although interestingly, both were matched with the exact same summary in a different language. To accommodate these pairs into CrossSum, we introduce ‘*induced pairs*.’ Formally, two summaries S_A, S_B in languages A, B are induced pairs if they are nearest neighbors of one another in A, B, their similarity score is below τ , and both are matched with S_C in language C as valid pairs (S_A, S_B), (S_B, S_C) (or through a chain of matched pairs in other languages).

We observed that induced pairs are prevalent if their languages are distant or low-resource. LaBSE uses contrastive learning (Guo et al., 2018; Yang et al., 2019) to rank parallel sentences over non-parallels. Since parallel pairs are mostly found for high-resource and linguistically close languages, we hypothesize that LaBSE fails to assign high similarity to sentences from languages that are not. We thus try to incorporate the induced pairs into CrossSum through a simple graph-based algorithm:

We represent all summaries as vertices in a graph and draw edges between two vertices if the summaries are matched as valid pairs. Then we find the connected components in the graph and draw edges (i.e., induced pairs) between all vertices in a component. Again to ensure quality, before computing the induced pairs, we use the max-flow min-cut theorem (Dantzig and Fulkerson, 1955) considering the similarity scores as edge weights to limit the size of each component to 50 vertices (since ideally a component should have at most 45 vertices, one summary from each language) and set the minimum threshold to $\tau' = (\tau - 0.10)$.

We finally assembled the original matched pairs and induced pairs to create the CrossSum dataset. Figure 2 shows the article-summary statistics for all language pairs in CrossSum.

Implicit Leakage We initially made the train-dev-test splits respecting the original XL-Sum split: and performed an initial assessment of the dataset using the splits by training a many-to-one model

¹Around 90% F_1 score is achieved using LaBSE in the BUCC tasks, hence it is expected that not all alignment will be correct in CrossSum. Since Hasan et al. (2021) reported summaries around this percentage to be good-quality in XL-Sum, we went ahead with this threshold.

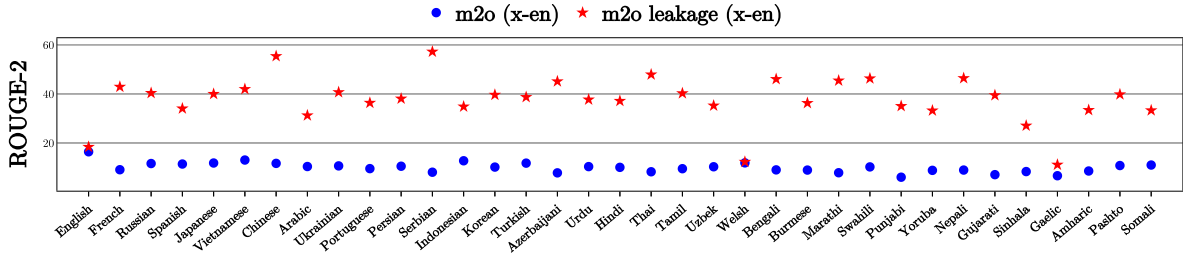


Figure 3: Training on the dataset respecting the original XL-Sum splits causes absurdly high ROUGE scores (marked red) in many-to-one models due to implicit data leakage. Therefore, we split taking the issue into account and consequently, models trained on the new set (marked blue) does not exhibit any unusual spike in ROUGE-2.

(articles written in any source language being summarized into one target language) in a supervised fashion. Upon evaluating the model, we found very high ROUGE-2 scores (up to 60) for many language pairs, even reaching as high as 80 for some (Figure 3). For contrast, Hasan et al. (2021) reported ROUGE-2 in the 10-20 range in the multilingual summarization task.

We inspected the model outputs and found that many summaries were exactly the same as the references. Through closer inspection, we found that all the articles, the summaries of which are exact copies of references, had their identical counterparts in some other language occurring in the training set. During training, the models were successfully able to align the representations of identical articles (albeit written in different languages) and were able to generate the exact same output by memorizing from the training sample. While models should undoubtedly be credited for being able to make these cross-lingual mappings, this is not ideal for benchmarking purposes as this creates unusually high ROUGE scores. We denote this phenomenon as ‘*implicit leakage*’ and make a new dataset split to avoid this. Before proceeding, we deduplicate the XL-Sum² dataset using semantic similarity, considering two summaries S_A, S'_A in language A to be duplicates if their LaBSE representations have similarity above 0.95. We take advantage of the component graph mentioned previously to handle the leakage and assign all article-summary pairs originating from a single component in the training (dev/test) set of CrossSum, creating an even 80%-10%-10% split for all language pairs. Since identical articles no longer appear in the train set of one language and dev/test set of another, the leakage is not observed anymore (Fig-

ure 3). We further validated this by inspecting the model outputs and found no exact copies.

3 Training & Evaluation Methodologies

In this section, we discuss the multistage sampling strategy for training cross-lingual text generation models and our proposed metric for evaluating model-generated summaries.

3.1 Multistage Language Sampling

From Figure 2, we can see that CrossSum is heavily imbalanced in terms of samples for different language pairs, and thus training directly without upsampling low-resource languages may result in their degraded performance. Conneau et al. (2020) used a probability smoothing technique for upsampling in multilingual pretraining and sampled all data points of a batch from one language. However, if we did the same for the language pairs in CrossSum, many batches would have duplicate samples since many pairs do not have enough examples, and at the same time, many would not be sampled during training for lack of enough training steps (due to a limitation of computational resources from our side). To address this, we adapt their algorithm to introduce a multistage upsampling method and ensure either the source or the target texts of a batch are sampled from the same language.

Let L_1, L_2, \dots, L_n be the languages of a cross-lingual source-target dataset. Let c_{ij} be the number of training samples where the source is from L_i and target from L_j . We compute the probabilities of the source languages by

$$p_i = \frac{\sum_{k=1}^n c_{ik}}{\sum_{j=1}^n \sum_{k=1}^n c_{jk}} \forall i \in \{1, 2, \dots, n\}$$

We then use an exponent smoothing factor α and normalize the probabilities

$$q_i = \frac{p_i^\alpha}{\sum_{j=1}^n p_j^\alpha} \forall i \in \{1, 2, \dots, n\}$$

²XL-Sum is deduplicated using lexical overlap methods only. But due to the risk of implicit leakage, which is not lexical, we further perform semantic deduplication.

Given the source language L_i , we now compute the probabilities of its target languages.

$$p_{j|i} = \frac{c_{ij}}{\sum_{k=1}^n c_{ik}} \forall j \in \{1, 2, \dots, n\}$$

We again smooth $p_{j|i}$ by a factor of β and obtain the normalized probabilities

$$q_{j|i} = \frac{p_{j|i}^\beta}{\sum_{k=1}^n p_{k|i}^\beta} \forall j \in \{1, 2, \dots, n\}$$

We analogously compute p_j and $p_{i|j}$ and, using them, describe the training algorithm with multi-stage sampling in Algorithm 1.

Note that the proposed algorithm can be applied to any cross-lingual seq2seq task where both the source and target languages are imbalanced.

3.2 Evaluating Summaries Across Languages

A sufficient number of reference samples are essential for the reliable evaluation of model-generated summaries. However, for many CrossSum language pairs, especially low-resource ones, even the training sets are very small, let alone their test sets. Being able to evaluate using reference summaries written in a different language would allow evaluation in a broad range of languages, especially for which there are inadequate references in the target language. Embedding-based similarity metrics (Zhang et al., 2020; Zhao et al., 2019) have gained popularity in the last few years. We draw inspiration from them and design a similarity metric that does not rely on the lexical overlap between the generated and reference texts. As a result, this new metric can effectively measure similarity across languages in a language-independent manner. We consider three essential factors for our metric:

1. Meaning Similarity: The generated summary and the reference summary should convey the same meaning irrespective of their language. Just like our alignment procedure from Section 2, we use LaBSE to compute the meaning similarity between the generated (s_{gen}) and reference summary (s_{ref}):

$$MS(s_{gen}, s_{ref}) = \text{emb}(s_{gen}) \cdot \text{emb}(s_{ref})^T,$$

where, $\text{emb}(s)$ denotes the embedding vector output of LaBSE for input text s .

2. Language Confidence: The metric should identify, with high confidence, that the summary is indeed being generated in the target language. As such, we use the *fastText* language-ID classifier (Joulin et al., 2017) to obtain the language probability distribution of the generated summary and

Algorithm 1: A pseudocode of the multi-stage sampling algorithm.

Input: $D_{ij} \forall i, j \in \{1, 2, \dots, n\}$: training data with source/target languages
 L_i/L_j ;
 $c_{ij} \leftarrow |D_{ij}| \forall i, j \in \{1, 2, \dots, n\}$;
 m : number of mini-batches.

```

1 Compute  $p_i, p_j, p_{j|i}, p_{i|j}$  using  $c_{ij}$ 
2 while (Model Not Covered) do
3    $batch \leftarrow \phi$ 
4   Sample  $r \sim Unif(0, 1)$ 
5   if  $r > 0.5$  then
6     Sample  $L_i \sim p_i$ 
7     for  $i \leftarrow 1$  to  $m$  do
8       Sample  $L_j \sim p_{j|i}$ 
9       Create mini-batch  $mb$  from  $D_{ij}$ 
10       $batch \leftarrow batch \cup \{mb\}$ 
11    end
12  end
13  else
14    Sample  $L_j \sim p_j$ 
15    for  $j \leftarrow 1$  to  $m$  do
16      Sample  $L_i \sim p_{i|j}$ 
17      Create mini-batch  $mb$  from  $D_{ij}$ 
18       $batch \leftarrow batch \cup \{mb\}$ 
19    end
20  end
21  Optimize model using  $batch$ 
22 end

```

define the Language Confidence (LC) as:

$$LC(s_{gen}, s_{ref}) = \begin{cases} 1, & \text{if } L_{ref} = \text{argmax } P(L_{gen}) \\ P(L_{gen} = L_{ref}), & \text{otherwise} \end{cases}$$

3. Length Penalty: Generated summaries should not be unnecessarily long, and the metric should penalize long summaries. While model-based metrics may indicate how similar a generated summary is to its reference and its language, it is not clear how they can be used to determine its brevity. As such, we adapt the BLEU (Papineni et al., 2002) brevity penalty to measure the length penalty of generated summaries:

$$LP(s_{gen}, s_{ref}) = \begin{cases} 1, & \text{if } |s_{gen}| \leq |s_{ref}| + c \\ \exp(1 - \frac{|s_{gen}|}{|s_{ref}| + c}), & \text{otherwise} \end{cases}$$

The languages of s_{gen} and s_{ref} may not be the same, and identical texts may vary in length across languages. Hence, we used a length offset c

to avoid penalizing generated summaries slightly longer than the references. By examining the standard deviation of mean summary lengths of the languages, we set $c = 6$.

We finally define our metric, **Language-agnostic Summary Evaluation (LaSE)** score as follows.

$$\text{LaSE}(s_{gen}, s_{ref}) = \text{MS}(s_{gen}, s_{ref}) \times \text{LC}(s_{gen}, s_{ref}) \times \text{LP}(s_{gen}, s_{ref}) \quad (1)$$

4 Experiments & Benchmarks

We aim to train one model to generate summaries in any target language for an input article from another language by providing explicit cross-lingual supervision. Fine-tuning pretrained language models (Devlin et al., 2019; Xue et al., 2021) have shown state-of-the-art results on monolingual and multilingual abstractive text summarization (Rothe et al., 2020; Hasan et al., 2021). Many pretrained multilingual generative models are currently available, some prominent ones being mBART (Liu et al., 2020), CRISS (Tran et al., 2020), mT5 (Xue et al., 2021). Though CRISS is pretrained with a cross-lingual objective, which better suits our use case, in contrast to the multilingual objective of mBART and mT5, we choose mT5 for fine-tuning because of its broad coverage of 101 languages with support for 41 languages from CrossSum.

We compare our proposed multistage many-to-many (m2m) model with the standard unistage m2m model as well as many-to-one (m2o) and one-to-many (o2m) models, standards for cross-lingual summarization. We train four different m2o and o2m models using four highly spoken and typologically diverse pivot languages: English, Hindi, Arabic, and Russian. As another baseline, we use a summarize-then-translate pipeline. First, we fine-tune mT5 on our proposed split of the in-language data to obtain a multilingual summarization model. Then we use the M2M-100 model (Fan et al., 2021) (418M parameters variant) to translate the summaries into the target language.

We fine-tune mT5-base with the multistage algorithm with batch size 256, mini-batch size 32 on CrossSum (together with the in-language samples) with $\alpha = 0.5$, $\beta = 0.75$. The unistage m2m model is sampled with $\alpha = 0.25$, and each batch is packed with 8 mini-batches, each sample of which being taken from one language pair. m2o and o2m models are also trained in the same manner. All models are trained for 25k steps on 8 Nvidia Tesla

Target Lang.	ROUGE-2 vs. LaSE-in-lang. Pearson/Spearman	LaSE-in-lang vs. LaSE-out-lang. Pearson/Spearman
English	0.923/0.821	0.931/0.929
Hindi	0.967/1.000	0.940/0.600
Arabic	0.963/1.000	0.924/1.000
Russian	0.477/0.489	0.024/0.257

Table 1: Correlation analysis of ROUGE-2 and LaSE for different target languages.

P100 GPUs for 3 days. We discard a language pair from training if it has fewer than 30 training samples to prevent too many duplicates in a mini-batch. We limit the input to 512 and output to 84 tokens and use language-specific BOS (beginning of sequence) tokens (Wu et al., 2016) for guiding the decoder to generate summaries in the intended target language during inference and use a length penalty of 0.6. We show the evaluation results using ROUGE-2 and LaSE in Figure 4 and 5. Results indicate that the m2m model trained with our proposed algorithm consistently outperforms the unistage sampling model, the m2o and o2m models, and the summarize-then-translate pipeline.

5 Analysis & Discussions

Zero-shot/few-shot cross-lingual transfer: Experiments are done in Section 4 in fully supervised fashion. However, for many low-resource language pairs, samples are not available. Hence, it is attractive to be able to perform zero-shot cross-lingual transfer without relying on any labeled examples. To this end, we fine-tune mT5 with the in-language (both source and target are in the same language) samples only in a multilingual fashion and, during inference, vary the target language. Unfortunately, the model fails at generating cross-lingual summaries and performs in-language summarization instead. We also fine-tune m2o models in a zero-shot setting (with only the in-language samples of the target language). Here, the model can generate non-trivial summaries but still lags behind fully supervised models (results in the Appendix). We do not perform any few-shot experiments and leave them as potential future directions.

How reliable is LaSE? To validate the reliability of LaSE, we show its correlation with ROUGE-2. To further posit that LaSE is language-agnostic and can be effectively evaluated with references in a different language from the target, we swap the

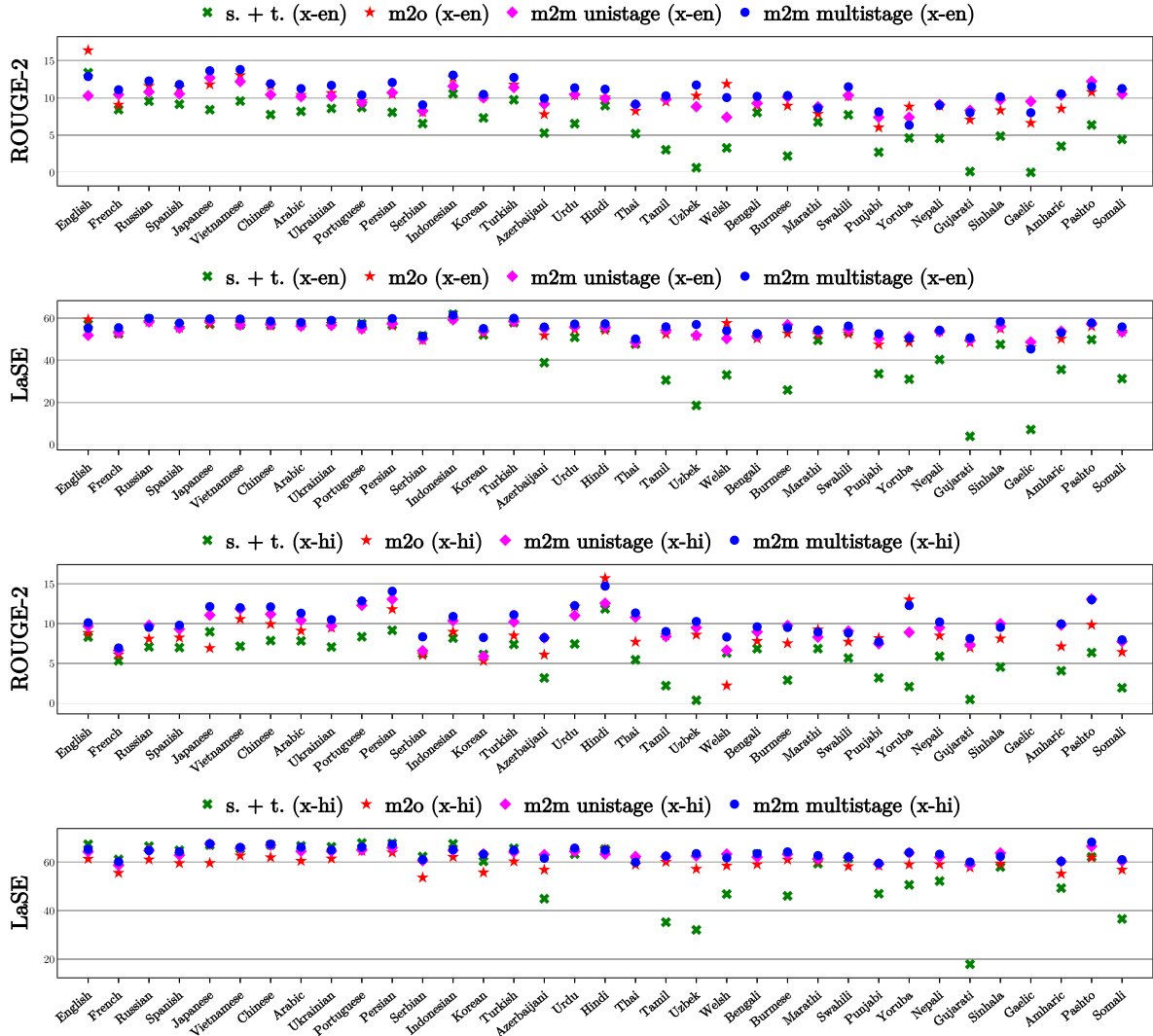


Figure 4: ROUGE-2 and LaSE scores for English and Hindi as target pivots as the sources languages vary. Scores indicate that our many-to-many (m2m) model with multistage sampling significantly outperforms the one-to-many models, summarize-then-translate and unistage m2m baselines models on most languages. The comparisons with other pivots are shown in the appendix due to space restrictions.

reference texts with the references in the language of the source text and show the correlation between the two variants of LaSE. We present the Pearson and Spearman correlation coefficients in Table 1. Since we were concerned that data scarcity would question the reliability of evaluation, we only take those language pairs into account that have at least 500 test samples. Results show that there is a high correlation between ROUGE-2 and LaSE for English, Hindi, and Arabic, and moderate for Russian. On the other hand, we find a strong correlation even when the references are swapped for the three above-mentioned languages. However, for Russian, we observed little to no correlation. We wish to investigate this discrepancy in the future and find ways to mitigate this.

6 Related Works

Pipeline-based methods were popular at the beginning stages of cross-lingual summarization research (Leuski et al., 2003; Orasan and Chiorean, 2008; Wan et al., 2010), breaking it into two sequential summarization and translation tasks. The lack of large datasets was a major obstruction towards end-to-end methods. End-to-end methods that performed cross-lingual summarization with a single model gained popularity with the emergence of neural models. Using a synthetic dataset, Zhu et al. (2019); Cao et al. (2020) performed cross-lingual summarization with a dual Transformer architecture in a multitask framework, while Bai et al. (2021) propose a single encoder-decoder for better transfer across tasks. Until recently, cross-lingual

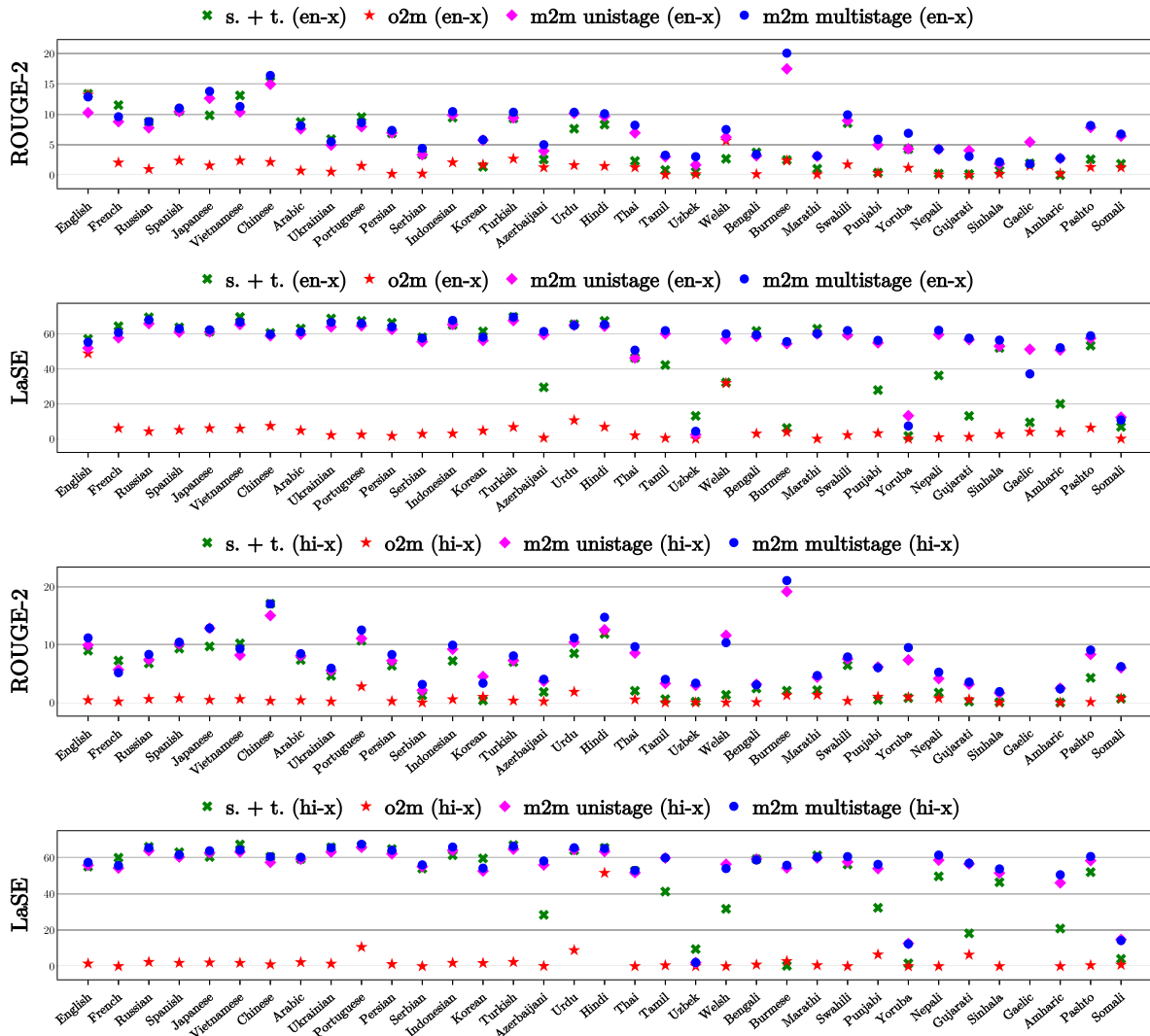


Figure 5: ROUGE-2 and LaSE scores for English and Hindi as source pivots as the target languages vary. Scores indicate that our many-to-many (m2m) model with multistage sampling significantly outperforms the one-to-many models, summarize-then-translate and unistage m2m baselines models on most languages. The comparisons with other pivots are shown in the appendix due to space restrictions.

summarization was limited to English-Chinese pair only due to the lack of benchmark datasets. To promote the task beyond them, [Ladhak et al. \(2020\)](#) introduced Wikilingua, a large-scale many-to-one dataset with English as the pivot language, while [Perez-Beltrachini and Lapata \(2021\)](#) introduced XWikis, containing 4 European languages in 12 many-to-many directions.

7 Conclusion & Future Works

In this paper, we present CrossSum, a large-scale, non-English-centric cross-lingual abstractive summarization dataset containing 1.7 million samples across 1500+ language pairs. CrossSum provides the first publicly available cross-lingual summarization dataset and benchmarks for many of these

pairs. We also make the alignment scripts available for the researchers, which will help produce better alignments. Furthermore, we introduced a new multistage sampling algorithm that can be generalized to any cross-lingual generation task and a new language-agnostic metric for evaluating cross-lingual summaries when references in the target languages may not be available. Additionally, we demonstrate that training one multilingual model can help better cross-lingual summarization than baselines. Moreover, CrossSum can also be helpful in zero-shot cross-lingual settings.

In the future, we will investigate the use of our dataset for other summarization tasks, e.g., multi-document ([Fabbri et al., 2019](#)) and multi-modal summarization ([Zhu et al., 2018](#)).

Acknowledgements

This work was performed using the OzSTAR national facility at the Swinburne University of Technology. The OzSTAR program receives funding in part from the Astronomy National Collaborative Research Infrastructure Strategy (NCRIS) allocation provided by the Australian Government.

References

- Mikel Artetxe and Holger Schwenk. 2019a. [Margin-based parallel corpus mining with multilingual sentence embeddings](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3197–3203, Florence, Italy. Association for Computational Linguistics.
- Mikel Artetxe and Holger Schwenk. 2019b. [Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond](#). *Transactions of the Association for Computational Linguistics*, 7:597–610.
- Yu Bai, Yang Gao, and Heyan Huang. 2021. [Cross-lingual abstractive summarization with limited parallel resources](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6910–6924, Online. Association for Computational Linguistics.
- Yue Cao, Hui Liu, and Xiaojun Wan. 2020. [Jointly learning to align and summarize for neural cross-lingual summarization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6220–6231, Online. Association for Computational Linguistics.
- Kyunghyun Cho, Bart van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. [Learning phrase representations using RNN encoder–decoder for statistical machine translation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- George Bernard Dantzig and Delbert Ray Fulkerson. 1955. On the max flow min cut theorem of networks. Technical report, RAND CORP SANTA MONICA CA.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota, USA. Association for Computational Linguistics.
- Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and Dragomir Radev. 2019. [Multi-news: A large-scale multi-document summarization dataset and abstractive hierarchical model](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 1074–1084, Florence, Italy. Association for Computational Linguistics.
- Angela Fan, Shruti Bhosale, Holger Schwenk, Zhiyi Ma, Ahmed El-Kishky, Siddharth Goyal, Mandeep Baines, Onur Celebi, Guillaume Wenzek, Vishrav Chaudhary, et al. 2021. [Beyond english-centric multilingual machine translation](#). *Journal of Machine Learning Research*, 22(107):1–48.
- Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. [Language-agnostic BERT sentence embedding](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland. Association for Computational Linguistics.
- Mandy Guo, Qinlan Shen, Yinfei Yang, Heming Ge, Daniel Cer, Gustavo Hernandez Abrego, Keith Stevens, Noah Constant, Yun-Hsuan Sung, Brian Strope, and Ray Kurzweil. 2018. [Effective parallel corpus mining using bilingual sentence embeddings](#). In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 165–176, Belgium, Brussels. Association for Computational Linguistics.
- Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. [XLsum: Large-scale multilingual abstractive summarization for 44 languages](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4693–4703, Online. Association for Computational Linguistics.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. [Bag of tricks for efficient text classification](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.
- Faisal Ladhak, Esin Durmus, Claire Cardie, and Kathleen McKeown. 2020. [WikiLingua: A new benchmark dataset for cross-lingual abstractive summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4034–

- 4048, Online. Association for Computational Linguistics.
- Anton Leuski, Chin-Yew Lin, Liang Zhou, Ulrich Germann, Franz Josef Och, and Eduard Hovy. 2003. Cross-lingual c* st* rd: English access to hindi information. *ACM Transactions on Asian Language Information Processing (TALIP)*, 2(3):245–269.
- Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation. *Transactions of the Association for Computational Linguistics*, 8:726–742.
- Constantin Orasan and Oana Andreea Chiorean. 2008. Evaluation of a cross-lingual romanian-english multi-document summariser. In *Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC'08)*, Marrakech, Morocco. European Language Resources Association (ELRA). [Http://www.lrec-conf.org/proceedings/lrec2008/](http://www.lrec-conf.org/proceedings/lrec2008/).
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. **Bleu: a method for automatic evaluation of machine translation**. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Laura Perez-Beltrachini and Mirella Lapata. 2021. **Models and datasets for cross-lingual summarisation**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 9408–9423, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sascha Rothe, Shashi Narayan, and Aliaksei Severyn. 2020. **Leveraging pre-trained checkpoints for sequence generation tasks**. *Transactions of the Association for Computational Linguistics*, 8:264–280.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. **Sequence to sequence learning with neural networks**. In *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS 2014)*, pages 3104–3112, Montreal, Canada.
- Chau Tran, Yuqing Tang, Xian Li, and Jiatao Gu. 2020. **Cross-lingual retrieval for iterative self-supervised training**. In *Advances in Neural Information Processing Systems*, volume 33, pages 2207–2219. Curran Associates, Inc.
- Daniel Varab and Natalie Schluter. 2021. **MassiveSumm: a very large-scale, very multilingual, news summarisation dataset**. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 10150–10161, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. **Attention is all you need**. In *Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017)*, page 6000–6010, Long Beach, California, USA.
- Xiaojun Wan, Huiying Li, and Jianguo Xiao. 2010. **Cross-language document summarization based on machine translation quality prediction**. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 917–926, Uppsala, Sweden. Association for Computational Linguistics.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. **Google’s neural machine translation system: Bridging the gap between human and machine translation**. *arXiv:1609.08144*.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. **mT5: A massively multilingual pre-trained text-to-text transformer**. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 483–498, Online. Association for Computational Linguistics.
- Yinfei Yang, Gustavo Hernandez Abrego, Steve Yuan, Mandy Guo, Qinlan Shen, Daniel Cer, Yun-hsuan Sung, Brian Strope, and Ray Kurzweil. 2019. **Improving multilingual sentence embedding using bi-directional dual encoder with additive margin softmax**. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 5370–5378. International Conferences on Artificial Intelligence Organization.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. **Bertscore: Evaluating text generation with bert**. In *International Conference on Learning Representations*.
- Wei Zhao, Maxime Peyrard, Fei Liu, Yang Gao, Christian M. Meyer, and Steffen Eger. 2019. **MoverScore: Text generation evaluating with contextualized embeddings and earth mover distance**. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 563–578, Hong Kong, China. Association for Computational Linguistics.
- Junnan Zhu, Haoran Li, Tianshang Liu, Yu Zhou, Jijun Zhang, and Chengqing Zong. 2018. **Msmo: Multimodal summarization with multimodal output**. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pages 4154–4164.

Junnan Zhu, Qian Wang, Yining Wang, Yu Zhou, Jijun Zhang, Shaonan Wang, and Chengqing Zong. 2019. *NCLS: Neural cross-lingual summarization*. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 3054–3064, Hong Kong, China. Association for Computational Linguistics.

Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp. 2017. Overview of the second bucc shared task: Spotting parallel sentences in comparable corpora. In *Proceedings of the 10th Workshop on Building and Using Comparable Corpora*, pages 60–67.

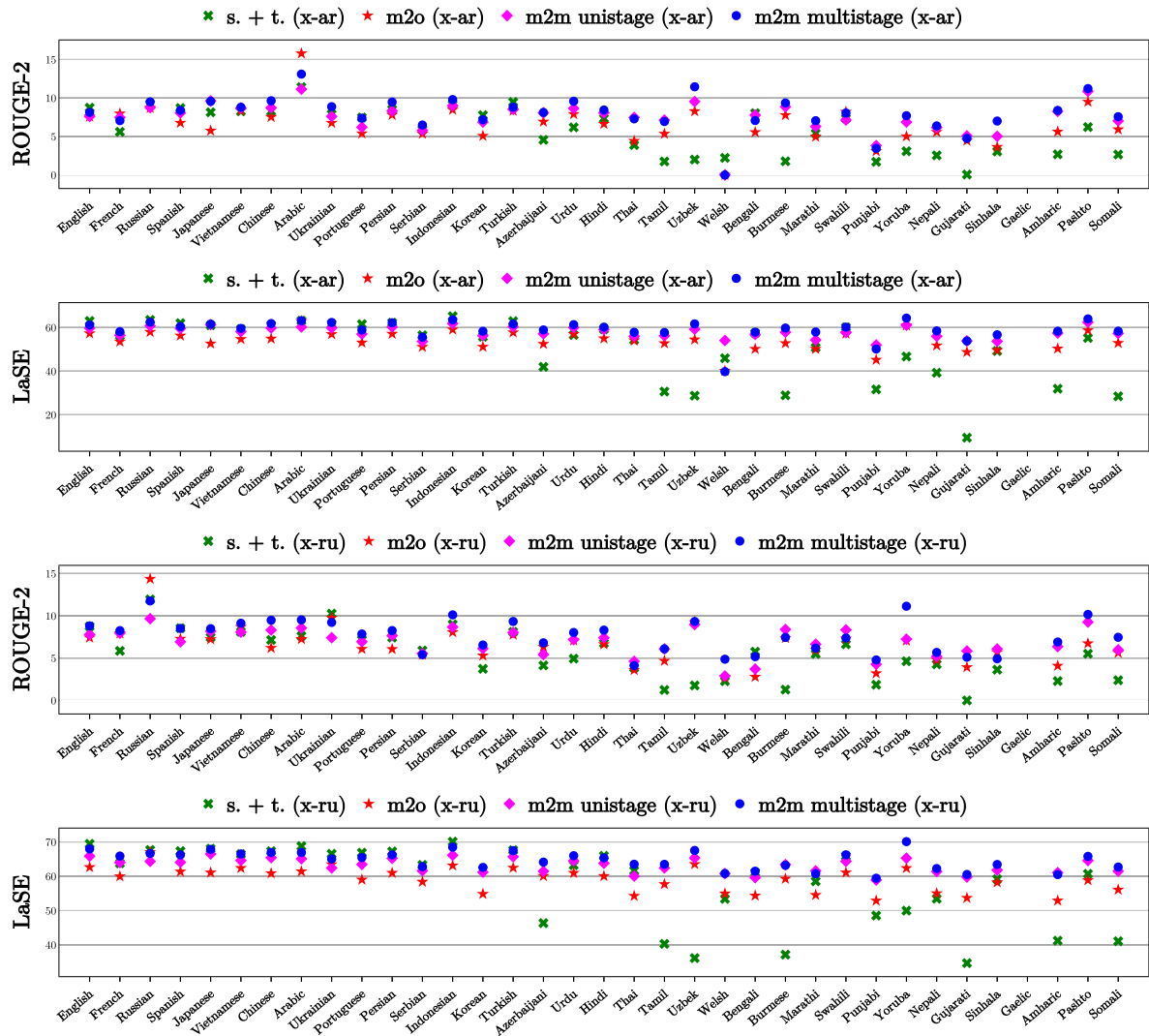


Figure 6: ROUGE-2 and LaSE scores for Arabic and Russian as target pivots as the sources languages vary. Scores indicate that our many-to-many (m2m) model with multistage sampling significantly outperforms the one-to-many models, summarize-then-translate and unistage m2m baselines models on most languages. The comparisons with other pivots are shown in the appendix due to space restrictions.

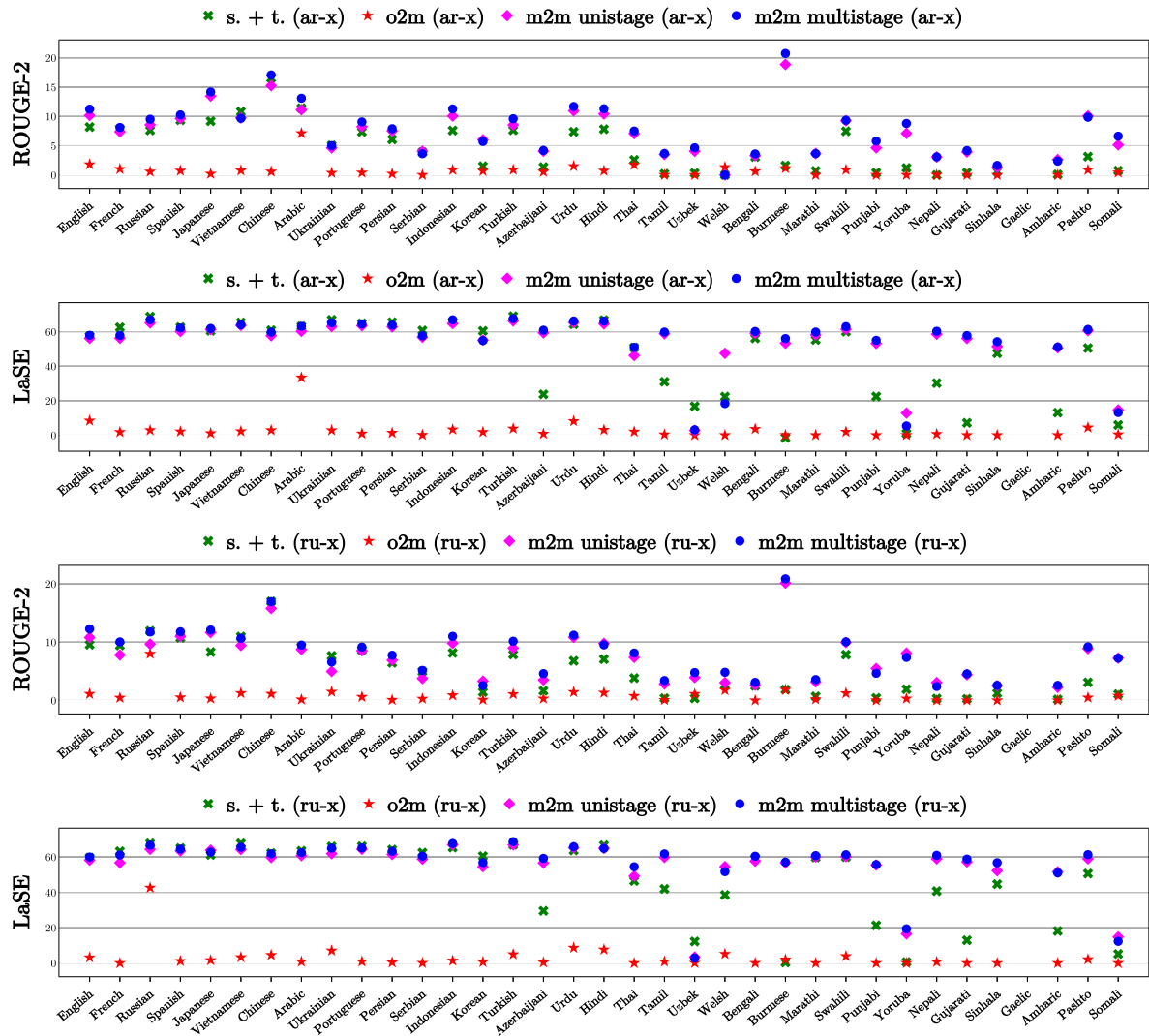


Figure 7: ROUGE-2 and LaSE scores for Arabic and Russian as source pivots as the target languages vary. Scores indicate that our many-to-many (m2m) model with multistage sampling significantly outperforms the one-to-many models, summarize-then-translate and unistage m2m baselines models on most languages. The comparisons with other pivots are shown in the appendix due to space restrictions.

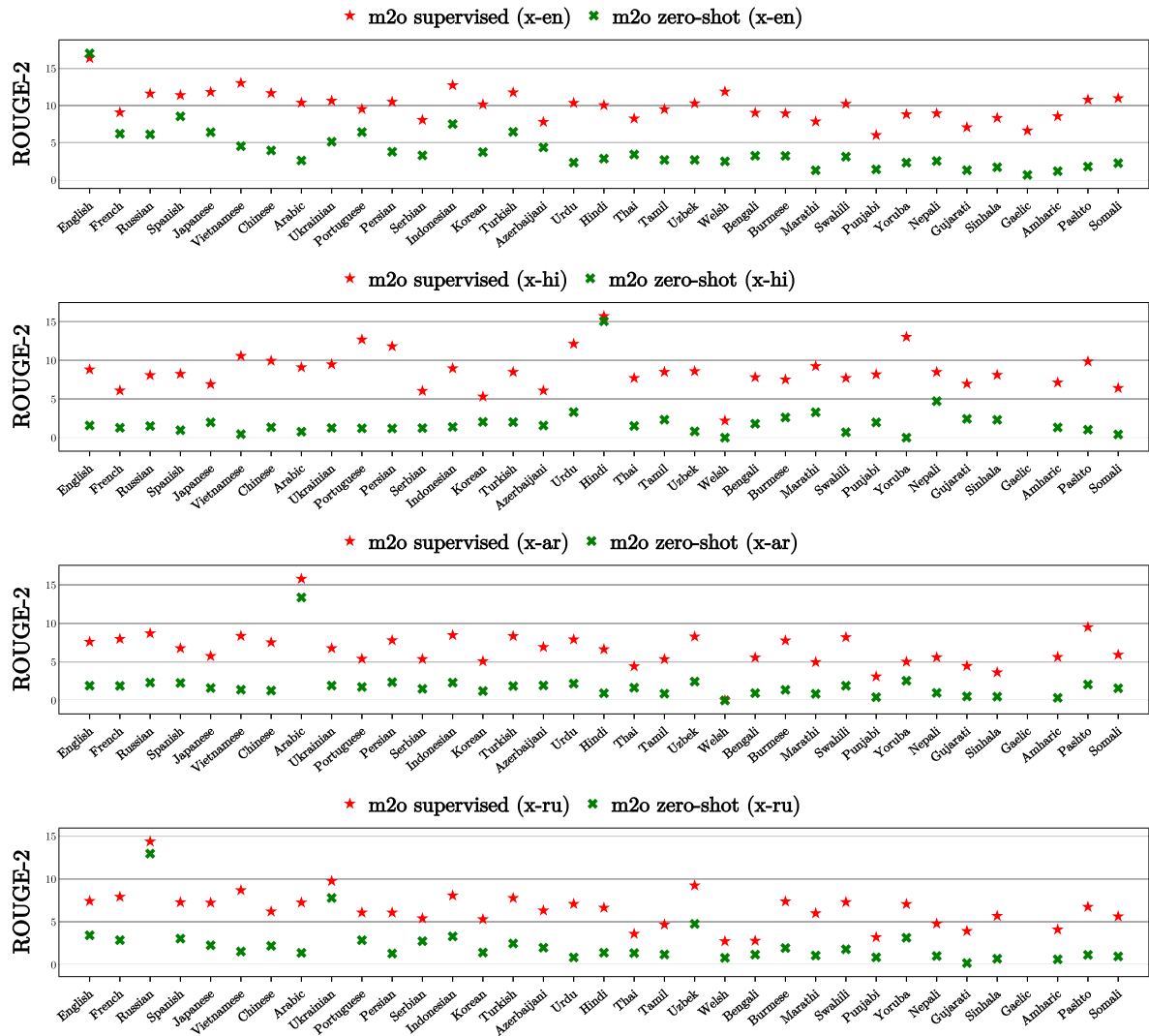


Figure 8: Zero-shot ROUGE-2 scores for the pivot languages as the target languages vary. The zero-shot models are trained with only the in-language samples of the pivot. Though the results are clearly behind the fully supervised model, the model is able to generate non-trivial summaries for many language pairs.

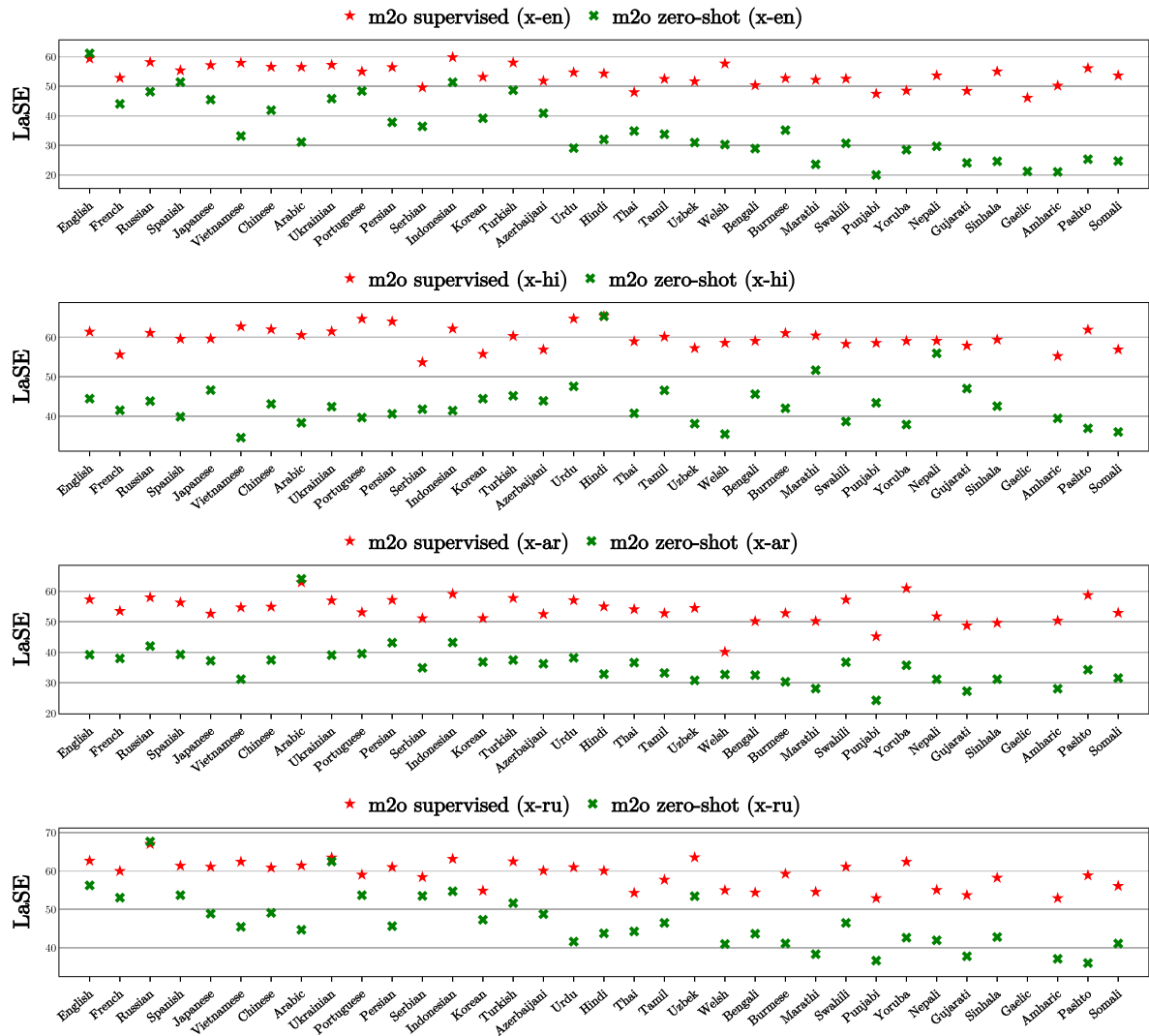


Figure 9: Zero-shot LaSE scores for the pivot languages as the target languages vary. The zero-shot models are trained with only the in-language samples of the pivot. Though the results are clearly behind the fully supervised model, the model is able to generate non-trivial summaries for many language pairs.

