

CrossSum: Beyond English-Centric Cross-Lingual Summarization for 1,500+ Language Pairs

Anonymous ACL submission

Abstract

We present CrossSum, a large-scale cross-lingual summarization dataset comprising 1.68 million article-summary samples in 1,500+ language pairs. We create CrossSum by aligning identical articles written in different languages via cross-lingual retrieval from a multilingual summarization dataset and perform a controlled human evaluation to validate its quality. We propose a multistage data sampling algorithm to effectively train a cross-lingual summarization model capable of summarizing an article in any target language. We also introduce LaSE, an embedding-based metric for automatically evaluating model-generated summaries. LaSE is strongly correlated with ROUGE and, unlike ROUGE, can be reliably measured even in the absence of references in the target language. Performance on ROUGE and LaSE indicate that pretrained models fine-tuned on CrossSum consistently outperform baseline models. To the best of our knowledge, CrossSum is the largest cross-lingual summarization dataset and the first-ever that is not centered around English. We will release the dataset, alignment and training scripts, and the models to spur future research on cross-lingual summarization.

1 Introduction

Cross-lingual summarization (hereinafter XLS) 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 XLS thus employed pipeline methods such as translate-then-summarize (Leuski et al., 2003) and summarize-then-translate (Wan et al., 2010). Not only are they computationally expensive, having to use multiple models, but 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 additionally translate the texts to English for better understanding. Words 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) have aided in the emergence of end-to-end methods that can perform XLS with one single model (Zhu et al., 2019; Cao et al., 2020b). The availability of XLS datasets (Ladhak et al., 2020; Perez-Beltrachini and Lapata, 2021) has also helped this task gain popularity in recent times. However, they cover only a few languages, contain a small number of samples for training and evaluation, or use English as the pivot language (i.e., the target language always remains English), thereby limiting their applicability to a great extent.

To democratize XLS beyond high-resource languages, in this work, we introduce **CrossSum**, a large-scale XLS dataset containing 1.68 million

058 article-summary samples in 1,500+ language pairs. 108
059 We align identical articles written in different languages 109
060 via cross-lingual retrieval from the multi-lingual 110
061 XL-Sum (Hasan et al., 2021) dataset. We 111
062 perform a controlled human evaluation of Cross- 112
063 Sum spanning nine languages from high-resource 113
064 to low-resource and show that the alignments are 114
065 highly accurate. We design a multistage sampling 115
066 algorithm for successfully training models that can 116
067 generate a summary in any target language for an 117
068 article in any source language. For the first time, 118
069 we perform XLS with CrossSum on a broad and 119
070 diverse set of languages without relying on English 120
071 as the standalone pivot language, consistently out- 121
072 performing many-to-one and one-to-many models, 122
073 as well as summarize-then-translate baselines. 123

074 We propose **LaSE**, an embedding-based metric 124
075 for evaluating summaries when reference sum- 125
076 maries may not be available in the target language 126
077 but another language, potentially opening new 127
078 doors for evaluating low-resource languages. Fur- 128
079 thermore, we demonstrate the reliability of LaSE 129
080 by its high correlation with ROUGE (Lin, 2004), 130
081 the de-facto metric for summarization evaluation. 131

082 To the best of our knowledge, CrossSum is the 132
083 first publicly available XLS dataset for a large num- 133
084 ber of language pairs. We are releasing the dataset, 134
085 alignment and training scripts, and models, hoping 135
086 that these resources will encourage the community 136
087 to push the boundaries of XLS beyond English and 137
088 other high-resource languages. 138

089 2 The CrossSum Dataset

090 The most straightforward way of curating a high- 140
091 quality XLS dataset is via crowd-sourcing (Nguyen 141
092 and Daumé III, 2019). However, it may be dif- 142
093 ficult to find crowd workers having professional 143
094 command over low-resource languages or distant 144
095 language pairs. Moreover, scalability issues might 145
096 arise due to the time and budget constraints for 146
097 crowd-sourcing. Therefore, synthetic (Zhu et al., 147
098 2019) and automatic methods (Ladhak et al., 2020; 148
099 Perez-Beltrachini and Lapata, 2021) have gained 149
100 traction over crowd-sourcing. 150

101 Automatic curation of an XLS dataset is sim- 151
102 ply to pair an article A in a source language with 152
103 the summary of an identical article B written in a 153
104 different target language (Figure 1), assuming the 154
105 availability of a multilingual dataset having identi-
106 cal contents in different languages. Two contempo-
107 rary works have compiled large-scale multilingual

summarization datasets, namely XL-Sum (Hasan 108
et al., 2021) (1.35M samples in 45 languages) and 109
MassiveSumm (Varab and Schluter, 2021) (28.8M 110
samples in 92 languages). Though substantially 111
larger than the other, MassiveSumm is not publicly 112
available. Since public availability is crucial for 113
promoting open research, we opted for XL-Sum, 114
distributed under a non-commercial license. Addi- 115
tionally, all articles of XL-Sum are crawled from 116
a single source, BBC News. We observed that 117
BBC publishes similar news content in different 118
languages and follow similar summarization strate- 119
gies. Hence adopting XL-Sum would increase the 120
quality and quantity of the article-summary pairs. 121

Unlike previous automatic methods, there are 122
no explicit links between identical articles in XL- 123
Sum. Fortunately, language-agnostic sentence rep- 124
resentations (Artetxe and Schwenk, 2019a; Feng 125
et al., 2022) have achieved state-of-the-art results 126
in cross-lingual text mining (Zweigenbaum et al., 127
2017; Artetxe and Schwenk, 2019b), and hence, 128
we use them to search identical contents across 129
languages. For simplicity¹, we perform the search 130
over summaries only. To ensure maximum quality, 131
we set two conditions for a summary S_A in lan- 132
guage A to be matched with another summary S_B 133
in language B : 134

1. S_B must be the nearest neighbor of S_A among 135
all summaries in B , and vice-versa. 136
2. The similarity between S_A and S_B must be 137
above the threshold, τ . 138

The similarity of a summary pair is measured by 139
the inner product of their Language-agnostic BERT 140
Sentence Representations (LaBSE) (Feng et al., 141
2022) (a unit vector for an input text sequence). 142
We empirically set the similarity threshold as the 143
average over all languages that maximized their 144
respective F_1 score ($\tau = 0.7437$) in the BUCC 145
mining tasks (Zweigenbaum et al., 2017).² 146

Induced Pairs We observed that many summary 147
pairs, despite being nearest neighbors in their lan- 148
guage pairs, were filtered out because of the thresh- 149
old τ . Although interestingly, both were matched 150
with the same summary in a different language. 151
Moreover, these pairs are prevalent if their lan- 152
guages are distant or low-resource. LaBSE uses 153
contrastive learning (Guo et al., 2018; Yang et al., 154

¹The entire procedure is described in Appendix A.

²Around 90% F_1 is achieved using LaBSE in BUCC, hence not all CrossSum alignments will be correct. Therefore, in the following section, we further assess the quality of the alignments using human evaluation.

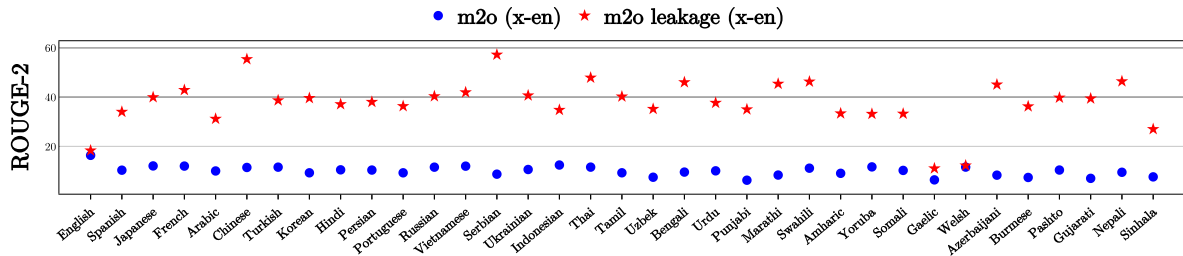


Figure 2: Training on the dataset respecting the original XL-Sum splits causes unusually high ROUGE scores (marked red) in many-to-one models due to implicit data leakage. Therefore, we redid the splits taking the issue into account, and consequently, models trained on the new set (marked blue) do not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

2019) to rank parallel sentences over non-parallel. 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.

To include these pairs into CrossSum, we introduce the notion ‘*induced pairs*.’ Formally, two summaries S_A, S_B in languages A, B are induced pairs if they are nearest neighbors of each other in A, B , their similarity score is below τ , and both are matched with S_C in language C as valid pairs $(S_A, S_C), (S_C, S_B)$, or through a chain of valid pairs $(S_A, S_C), (S_C, S_D), \dots, (S_Y, S_Z), (S_Z, S_B)$ in languages $\{C, D, \dots, Y, Z\}$.

We thus incorporate the induced pairs into CrossSum through a simple graph-based algorithm. First, we represent all summaries as vertices in a graph and draw an edge 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 their minimum acceptance threshold to $\tau' = (\tau - 0.10)$.

We finally assembled the original matched pairs and induced pairs to create the CrossSum dataset. Figure 6 (Appendix) shows the article-summary statistics for all language pairs in CrossSum. As evident from the figure, CrossSum is not centered only around the English language but rather distributed across multiple languages.

Implicit Leakage We initially made the train-dev-test splits respecting the original XL-Sum

splits and performed an initial assessment of CrossSum by training a many-to-one model (articles written in any source language being summarized into one target language). Upon evaluation, we found very high ROUGE-2 scores (around 40) for many language pairs, even reaching as high as 60 for some (Figure 2). In contrast, Hasan et al. (2021) reported ROUGE-2 in the 10-20 range for the multilingual summarization task.

We inspected the model outputs and found that many summaries were precisely 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 model was able to align the representations of identical articles (albeit written in different languages) and generate the 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-Sum dataset³ 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 address the leakage and assign all article-summary pairs originating from a single component in the training (dev/test) set of CrossSum, creating an 80%-10%-10% split for all language pairs. Since identical articles no longer appear in the train set of one and the dev/test set

³XL-Sum has been deduplicated using lexical overlap methods only. But due to the risk of implicit leakage, which is not lexical, we further perform semantic deduplication.

of another, the leakage is not observed anymore (Figure 2). We further validated this by inspecting the model outputs and found no exact copies.

3 Human Evaluation of CrossSum

To establish the validity of our automatic alignment pipeline, we conducted a human evaluation to study the quality of the cross-lingual alignments.

We selected all possible combinations of language pairs from a list of nine languages ranging from high-resource to low-resource to assess the alignment quality in different pair configurations (e.g., high-high, low-high, low-low) as per the language diversity categorization by Joshi et al. (2020). We chose three high-resource languages, English, Arabic, and (simplified) Chinese (category 4 and 5), three mid-resource languages, Indonesian, Bengali, and Urdu (category 3), and three low-resource languages, Punjabi, Swahili, and Pashto (category 1 and 2), as representative languages and randomly sampled fifty cross-lingual summary alignments from each language pair for annotation. As a direct evaluation of these pairs would require bilingually-proficient annotators for both languages, which are practically intractable for distantly related languages (e.g., Bengali-Swahili), we resorted to a pivoting approach during annotation for language pairs that do not contain English. For a language pair $(l_1 - l_2)$, where $l_1 \neq en$ and $l_2 \neq en$, we sampled alignments (x, y) such that $\exists(x, e) \in (l_1 - en)$ and $\exists(y, e) \in (l_2 - en)$, for an English article e . In other words, we ensure that both the articles of the sampled cross-lingual pair have a corresponding cross-lingual pair with an English article. An alignment (x, y) would be deemed correct if both (x, e) and (y, e) are correct. This formulation thus reduced the original problem to annotating samples from language pairs $(l_1 - en)$ and $(l_2 - en)$, where l_1 and l_2 are from the previously selected languages that are not English.

We hired bilingually proficient expert annotators adept in the language of interest and English. Two annotators labeled each language pair where one language is English. We presented them with corresponding summaries of the cross-lingual pairs (and optionally the articles themselves) and elicited yes/no answers to the question:

“Can the provided sequences be considered summaries for the same article?”⁴

⁴We do not explicitly evaluate article-summary correctness as this has already been studied in work on XL-Sum. This was

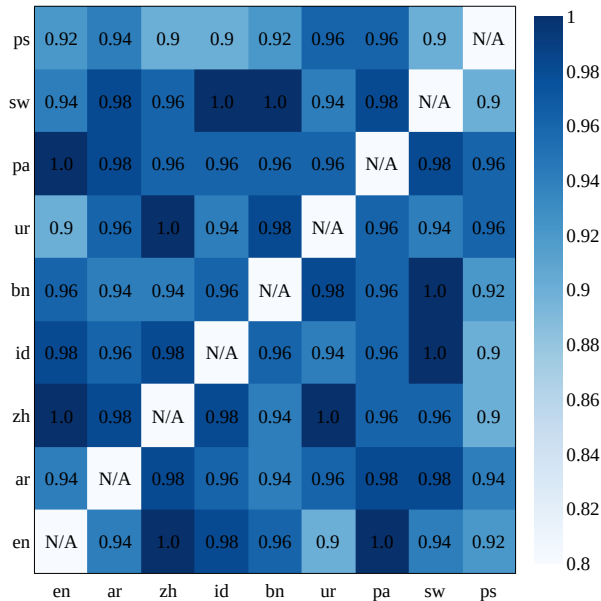


Figure 3: A heatmap showing alignment accuracies of different language pairs obtained by human evaluation.

We deem a sequence pair accurate if both annotators judge it as valid. We show the accuracy of the language pairs in Figure 3.

As evident from the figure, the annotators judge the aligned summaries to be highly accurate, with an average accuracy of 95.67%. We used Cohen’s Kappa (Cohen, 1960) to establish the inter-annotator agreement and show the corresponding statistics in Table 2 in the Appendix.

4 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.

4.1 Multistage Language Sampling

From Figure 6, it can be observed that CrossSum is heavily imbalanced. Thus, training directly without upsampling low-resource languages may result in their degraded performance. Conneau et al. (2020) used probability smoothing for upsampling in multilingual pretraining and sampled all data points of a batch from one language. However, applying this technique to the language pairs in CrossSum would result in many batches having duplicate samples as many language pairs do not have enough examples. At the same time, many would not be sampled during training for lack of enough training steps (due

also done to reduce annotation costs.

Algorithm 1: Multistage sampling

Input: $D_{ij} \forall i, j \in \{1, 2, \dots, n\}$: training data with tgt/src 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  $q_i, q_{j|i}$  using  $c_{ij}$ 
2 while (Model Not Converged) do
3   batch  $\leftarrow \phi$ 
4   Sample  $L_i \sim q_i$ 
5   for  $k \leftarrow 1$  to  $m$  do
6     Sample  $L_j \sim q_{j|i}$ 
7     Create mini-batch  $mb$  from  $D_{ij}$ 
8     batch  $\leftarrow$  batch  $\cup$   $\{mb\}$ 
9   Update model parameters using batch
```

to constraints on computational resources). To address this, we adapt their algorithm to introduce a multistage upsampling method to ensure that the target summaries 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, and c_{ij} be the number of training samples where the target is from L_i and source from L_j . We compute the probability p_i of each target language L_i by

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

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

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

Given the target language L_i , we now compute the probability of a source language L_j , represented by $p_{j|i}$.

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

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

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

Using the probabilities, we describe the training process with multistage 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.

4.2 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, even the training sets are small, let alone the test sets (the median size is only 33). For instance, the Japanese-Bengali language pair only has 34 test samples, which is too few for reliable evaluation. But the in-language test samples for Japanese and Bengali are nearly 1k. Being able to evaluate against reference summaries written in the source language would thus alleviate this insufficiency problem by leveraging the in-language test set of the source language.

For this purpose, cross-lingual similarity metrics that do not rely on lexical overlap (i.e., unlike ROUGE) are required. Embedding-based similarity metrics (Zhang et al., 2020; Zhao et al., 2019) have recently gained popularity. We draw inspiration from them and design a similarity metric that can effectively measure similarity across languages in a language-independent manner. We consider three essential factors:

1. Meaning Similarity: The generated summary and the reference summary should convey the same meaning irrespective of their languages. 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 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 language, it is unclear 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:

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

s_{gen} and s_{ref} may not be of the same language, and identical texts may vary in length across languages. Hence, we use 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.

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

5 Experiments & Analysis

One model capable of generating summaries in any target language for an input article from any source language is highly desirable. However, it may not be the case that such a ‘many-to-many’ model (m2m in brief) would outperform many-to-one (m2o) or one-to-many (o2m) models⁵, which are widely-used practices for XLS (Ladhak et al., 2020; Perez-Beltrachini and Lapata, 2021). In this section, we establish that the m2m model, in the presence of training samples from all possible language pairs, consistently outperforms m2o, o2m, and summarize-then-translate (s.+t.) baselines given equal training steps.

In addition to the proposed m2m model, we train five different m2o and o2m models using five highly spoken⁶ and typologically diverse pivot (i.e., the ‘one’ in m2o and o2m) languages: English, Chinese (simplified), Hindi, Arabic, and Russian. As another baseline, we use a summarize-then-translate pipeline. As fine-tuning pretrained language models (Devlin et al., 2019; Xue et al., 2021a) have shown state-of-the-art results on monolingual and multilingual text summarization (Rothe et al., 2020; Hasan et al., 2021), we fine-tune each model using a pretrained mT5 (Xue et al., 2021a) by providing explicit cross-lingual supervision⁷. We show the results on ROUGE-2 F1 and LaSE in Figures 4 and 5⁸.

⁵Discussed in detail in Appendix D.

⁶<https://w.wiki/Pss>

⁷Zero-shot cross-lingual transfer discussed in Appendix E.3.

⁸A detailed description of the training procedures and hyperparameter choices are detailed in Appendix E.1.

Results indicate that the m2m model consistently outperforms m2o, o2m, and s.+t., with an average ROUGE-2 (LaSE) score of 8.15 (57.15) over all languages tested, 3.12 (9.02) above s.+t. Moreover, compared to the o2m models on language pairs where the pivots are the targets, the m2m model scores 1.80 (5.84) over m2os, and on those where the pivots are the sources, 6.52 (51.80) over o2ms. We additionally perform a significance test of the m2m model’s performance in Appendix 3 and show it to be statistically superior to others.

Upon inspection, we found the m2o models to be able to generate non-trivial summaries, while the o2m models completely failed to produce cross-lingual summaries, performing in-language summarization for all targets⁹. s.+t. performed well on high-resource languages but poorly on low-resource ones. Inspection revealed this to be a limitation of the translation model used in the pipeline.

Target Lang.	ROUGE-2 vs.	LaSE-in-lang vs.
	LaSE-in-lang.	LaSE-out-lang.
	Pearson/Spearman	Pearson/Spearman
English	0.976/0.939	0.993/1.000
Arabic	0.903/0.987	0.968/0.942
Chinese	0.983/1.000	0.996/1.000
Indonesian	0.992/0.975	0.872/0.828
Bengali	0.947/0.902	0.819/0.771
Urdu	0.997/0.951	0.774/0.828
Punjabi	0.988/0.963	0.881/0.885
Swahili	0.990/0.951	0.979/0.885
Pashto	0.994/0.987	0.883/0.885

Table 1: Correlation analysis of ROUGE-2 and LaSE.

How reliable is LaSE? At first, we validated the reliability of LaSE by showing its correlation with ROUGE-2. We took different checkpoints of the in-language summarization model used in s.+t. and computed ROUGE-2 and LaSE for the nine languages in Section 3 for each checkpoint. The correlation coefficients of the computed scores are shown in the second column of Table 1. For all languages (from high- to low-resource), LaSE has a near-perfect correlation with ROUGE-2.

However, the purpose of LaSE is to show that it is language-agnostic and can even be computed in the absence of references in the target language.

⁹We hypothesize that varying the target language in a batch hampers the decoder’s ability to generate from a specific language, possibly because of the vast diversity of target languages in the batch (discussed further in Appendix F).

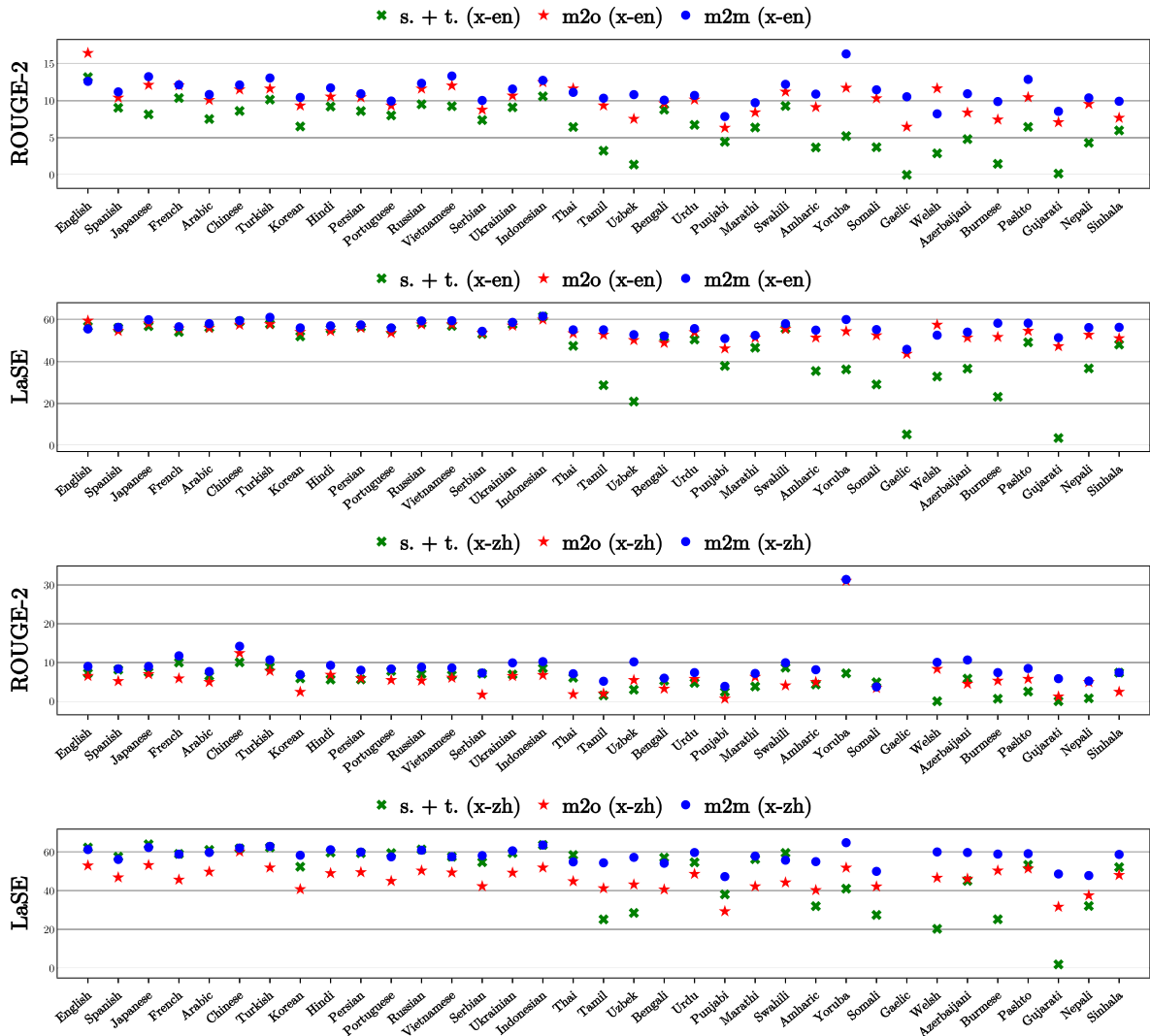


Figure 4: ROUGE-2 and LaSE scores for English and Chinese as target pivots as the source languages vary. The m2m model significantly outperforms the m2o models and summarize-then-translate baseline in most languages. The comparisons with other pivots are shown in the appendix (Figure 8) due to space limitations.

Therefore, we evaluate the summaries with references in a different language from the target using the m2m model. For each target language, we first compute the standard LaSE for different source languages (denoted as LaSE-in-lang). We again compute LaSE after swapping the reference texts with the references in the language of the input text¹⁰ (denoted as LaSE-out-lang). We then show the correlation between the two variants of LaSE in the third column of Table 1¹¹ for each target language. Results show a substantial correlation between the two variants of LaSE for all languages.

¹⁰Our curation method ensures that such summaries always exist in the corresponding test sets.

¹¹Since many test sets of the language pairs from Section 3 have too few samples for reliable evaluation (e.g., Punjabi-Pashto), for each target language, we use only the top-5 source languages by the number of their test set samples.

From these two experiments, we can conclude that LaSE is ideal for summary evaluation and can be computed in a language-independent manner.

6 Related Works

Pipeline-based methods were popular at the beginning stages of XLS research (Leuski et al., 2003; Orasan and Chiorean, 2008; Wan et al., 2010), breaking it into two sequential summarization and translation tasks. End-to-end methods that performed XLS with a single model gained popularity with the emergence of neural models. Ayana et al. (2018) used knowledge distillation (Hinton et al., 2015) to train a student XLS model from two summarization and translation teacher models. Using a synthetic dataset, Zhu et al. (2019); Cao et al. (2020a) performed XLS with a dual Transformer

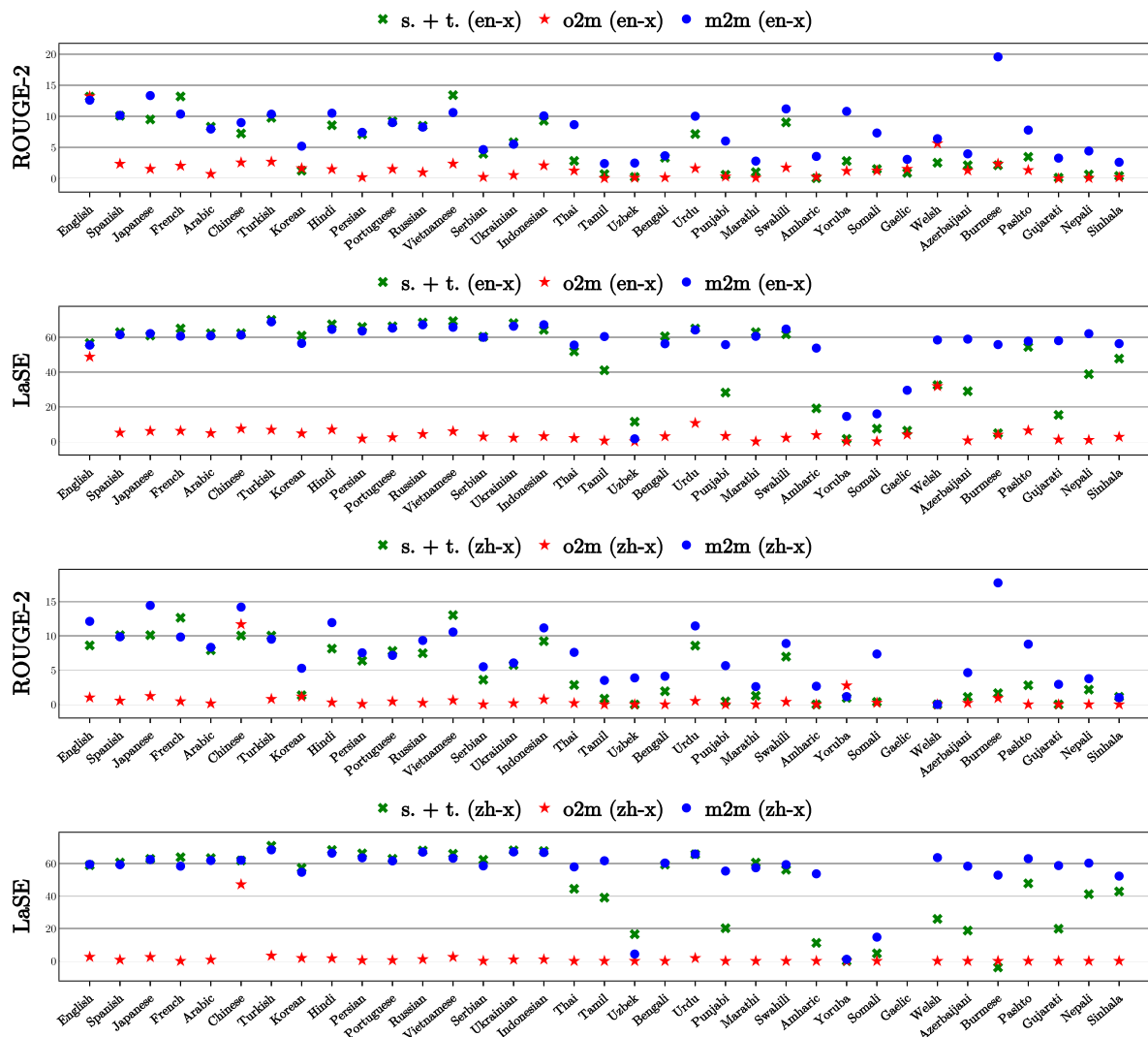


Figure 5: ROUGE-2 and LaSE scores for English and Chinese as source pivots as the target languages vary. The m2m model significantly outperforms the o2m models and summarize-then-translate baseline in most languages. The comparisons with other pivots are shown in the appendix (Figure 9) due to space limitations.

(Waswani et al., 2017) architecture in a multitask framework, while Bai et al. (2021) proposed a single encoder-decoder for better transfer across tasks. Until recently, XLS was limited to English-Chinese only due to the lack of benchmark datasets. To promote the task beyond, 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 languages in 12 directions.

7 Conclusion & Future Works

In this paper, we presented CrossSum, a large-scale, non-English-centric XLS dataset containing 1.68 million samples across 1500+ language pairs. CrossSum provides the first publicly avail-

able XLS dataset for many of these pairs. Performing a limited-scale human evaluation of CrossSum, we introduced a multistage sampling algorithm for general-purpose cross-lingual generation and a language-agnostic metric for evaluating summaries when references in the target languages may not be available. Additionally, we demonstrated that training one multilingual model can help towards better XLS than baselines. We also shed some light on the potential to perform zero-/few-shot XLS with CrossSum.

In the future, we will investigate the use of CrossSum for other summarization tasks, e.g., multi-document (Fabbri et al., 2019) and multi-modal summarization (Zhu et al., 2018). We would also like to explore better techniques for m2m, zero-shot, and few-shot summarization.

513 Limitations

514 Though we believe that our work has many merits,
515 some of its limitations must be acknowledged. De-
516 spite exhaustive human annotation being the most
517 reliable means of ensuring the maximum quality of
518 a dataset, we had to resort to automatic curation of
519 CrossSum due to the enormous scale of the dataset.
520 As identified in the human evaluation, not all of the
521 alignments made by LaBSE are correct. They are
522 primarily summaries describing similar (i.e., hav-
523 ing a substantial degree of syntactic or semantic
524 similarity) but non-identical events. LaBSE also
525 fails to penalize numerical mismatches, especially
526 if the summaries depict the same event.

527 Consequently, any mistake made by LaBSE in
528 the curation phase may propagate to the models
529 trained using CrossSum. And since LaBSE is a
530 component of the proposed LaSE metric, these bi-
531 ases may remain unidentified by LaSE in the evalu-
532 ation stage. However, no matter which automatic
533 method we use, there will be such frailties in these
534 extreme cases. Since the objective of this paper is
535 not to scrutinize the pitfalls of LaBSE but rather
536 to use it as a means of curation and evaluation, we
537 deem LaBSE the best choice due to its extensive
538 language coverage and empirical performance in
539 cross-lingual mining among existing alternatives.

540 Ethical Considerations

541 **License** CrossSum is a derivative of the XL-Sum
542 dataset. XL-Sum has been released under the
543 Creative Commons Attribution-NonCommercial-
544 ShareAlike 4.0 International License (CC BY-NC-
545 SA 4.0), allowing modifications and distributions
546 for non-commercial research purposes. We are
547 adhering to the terms of the license and will also
548 release CrossSum under the same license.

549 **Generated Text** All of our models use the mT5
550 model as the backbone, which is pretrained on a
551 large multilingual text corpus. For a text gener-
552 ation model, even small amounts of offensive or
553 harmful texts in pretraining could lead to danger-
554 ous biases in generated text (Luccioni and Viviano,
555 2021). Therefore, our models can potentially gener-
556 ate offensive or biased content learned during
557 the pretraining phase, which is beyond our control.
558 Text summarization systems have also been shown
559 to generate unfaithful and factually incorrect (albeit
560 fluent) (Maynez et al., 2020) texts. Thus, we sug-
561 gest carefully examining the potential biases before
562 considering them in any real-world deployment.

Human Evaluation Annotators were hired from
the graduates of an institute that provides profes-
sional language training for many languages, in-
cluding the ones evaluated in Section 3. Each an-
notator was given around 200-250 sequence pairs
to evaluate. Each annotation took around one and a
half minutes on average, with a total of approx-
imately 5-6 hours for annotating the whole set.
Annotators were paid hourly as per the standard
remuneration of bilingual professionals in local
currency.

Environmental Impact About 25 models were
trained in total in this paper. Each model was
trained for about 3 days on a 4-GPU Tesla P100
server. Assuming 0.08 kg/kWh carbon emission¹²,
less than 175kg of carbon was released into the
environment in this work, which is orders of mag-
nitude below the most computationally demanding
models.

582 References

- Judit Ács. 2019. Exploring bert’s vocabulary. *Blog Post*.
- 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.
- Ayana, Shi-qi Shen, Yun Chen, Cheng Yang, Zhiyuan Liu, and Mao-song Sun. 2018. [Zero-shot cross-lingual neural headline generation](#). *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(12):2319–2327.
- 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. 2020a. [Jointly learning to align and summarize for neural cross-lingual summarization](#). In *Proceedings of the 58th*

¹²<https://blog.google/technology/ai/minimizing-carbon-footprint/>

612		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. <i>Journal of Machine Learning Research</i> , 22(107):1–48.	668
613			669
614			670
615	Yue Cao, Xiaojun Wan, Jinge Yao, and Dian Yu. 2020b.		671
616	Multisumm: Towards a unified model for multi-		672
617	lingual abstractive summarization . In <i>Proceedings of</i>		673
618	<i>Thirty-Fourth AAAI Conference on Artificial Intelli-</i>		
619	<i>gence, AAAI 2020</i> , pages 11–18. AAAI Press.		
620	Kyunghyun Cho, Bart van Merriënboer, Caglar Gul-	Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Ari-	674
621	cehre, Dzmitry Bahdanau, Fethi Bougares, Holger	vazhagan, and Wei Wang. 2022. Language-agnostic	675
622	Schwenk, and Yoshua Bengio. 2014. Learning	BERT sentence embedding . In <i>Proceedings of the</i>	676
623	phrase representations using RNN encoder–decoder	<i>60th Annual Meeting of the Association for Compu-</i>	677
624	for statistical machine translation . In <i>Proceedings</i>	<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	678
625	<i>of the 2014 Conference on Empirical Methods in</i>	878–891, Dublin, Ireland. Association for Computa-	679
626	<i>Natural Language Processing (EMNLP)</i> , pages 1724–	tional Linguistics.	680
627	1734, Doha, Qatar. Association for Computational		
628	Linguistics.		
629	Jacob Cohen. 1960. A coefficient of agreement for	Mandy Guo, Qinlan Shen, Yinfei Yang, Heming	681
630	nominal scales. <i>Educational and psychological mea-</i>	Ge, Daniel Cer, Gustavo Hernandez Abrego, Keith	682
631	<i>surement</i> , 20(1):37–46.	Stevens, Noah Constant, Yun-Hsuan Sung, Brian	683
632		Strope, and Ray Kurzweil. 2018. Effective parallel	684
633	Alexis Conneau, Kartikay Khandelwal, Naman Goyal,	corpus mining using bilingual sentence embeddings .	685
634	Vishrav Chaudhary, Guillaume Wenzek, Francisco	In <i>Proceedings of the Third Conference on Machine</i>	686
635	Guzmán, Edouard Grave, Myle Ott, Luke Zettle-	<i>Translation: Research Papers</i> , pages 165–176, Brus-	687
636	moyer, and Veselin Stoyanov. 2020. Unsupervised	sels, Belgium. Association for Computational Lin-	688
637	cross-lingual representation learning at scale . In <i>Pro-</i>	guistics.	689
638	<i>ceedings of the 58th Annual Meeting of the Associa-</i>		
639	<i>tion for Computational Linguistics</i> , pages 8440–	Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Is-	690
640	8451, Online. Association for Computational Lin-	lam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang,	691
641	guistics.	M. Sohel Rahman, and Rifat Shahriyar. 2021. XL-	692
642	George Bernard Dantzig and Delbert Ray Fulkerson.	sum: Large-scale multilingual abstractive summariza-	693
643	1955. On the max flow min cut theorem of networks.	tion for 44 languages . In <i>Findings of the Association</i>	694
644	Technical report, RAND CORP SANTA MONICA	<i>for Computational Linguistics: ACL-IJCNLP 2021</i> ,	695
645	CA.	pages 4693–4703, Online. Association for Computa-	696
646	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and	tional Linguistics.	697
647	Kristina Toutanova. 2019. BERT: Pre-training of		
648	deep bidirectional transformers for language under-	Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. 2015.	698
649	standing . In <i>Proceedings of the 2019 Conference of</i>	Distilling the knowledge in a neural network . In	699
650	<i>the North American Chapter of the Association for</i>	<i>NIPS Deep Learning and Representation Learning</i>	700
651	<i>Computational Linguistics: Human Language Tech-</i>	<i>Workshop</i> .	701
652	<i>nologies, Volume 1 (Long and Short Papers)</i> , pages		
653	4171–4186, Minneapolis, Minnesota, USA. Associa-	Melvin Johnson, Mike Schuster, Quoc V Le, Maxim	702
654	tion for Computational Linguistics.	Krikun, Yonghui Wu, Zhifeng Chen, Nikhil Thorat,	703
655	Xiangyu Duan, Mingming Yin, Min Zhang, Boxing	Fernanda Viégas, Martin Wattenberg, Greg Corrado,	704
656	Chen, and Weihua Luo. 2019. Zero-shot cross-	et al. 2017. Google’s multilingual neural machine	705
657	lingual abstractive sentence summarization through	translation system: Enabling zero-shot translation.	706
658	teaching generation and attention . In <i>Proceedings of</i>	<i>Transactions of the Association for Computational</i>	707
659	<i>the 57th Annual Meeting of the Association for Com-</i>	<i>Linguistics</i> , 5:339–351.	708
660	<i>putational Linguistics</i> , pages 3162–3172, Florence,		
661	Italy. Association for Computational Linguistics.	Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika	709
662	Alexander Fabbri, Irene Li, Tianwei She, Suyi Li, and	Bali, and Monojit Choudhury. 2020. The state and	710
663	Dragomir Radev. 2019. Multi-news: A large-scale	fate of linguistic diversity and inclusion in the NLP	711
664	multi-document summarization dataset and abstrac-	world . In <i>Proceedings of the 58th Annual Meeting of</i>	712
665	tive hierarchical model . In <i>Proceedings of the 57th</i>	<i>the Association for Computational Linguistics</i> , pages	713
666	<i>Annual Meeting of the Association for Computational</i>	6282–6293, Online. Association for Computational	714
667	<i>Linguistics</i> , pages 1074–1084, Florence, Italy. Asso-	Linguistics.	715
	ciation for Computational Linguistics.		
		Armand Joulin, Edouard Grave, Piotr Bojanowski, and	716
		Tomas Mikolov. 2017. Bag of tricks for efficient	717
		text classification . In <i>Proceedings of the 15th Con-</i>	718
		<i>ference of the European Chapter of the Association</i>	719
		<i>for Computational Linguistics: Volume 2, Short Pa-</i>	720
		<i>pers</i> , pages 427–431, Valencia, Spain. Association	721
		for Computational Linguistics.	722
		Philipp Koehn. 2004. Statistical significance tests for	723
		machine translation evaluation . In <i>Proceedings of the</i>	724

725		2004 Conference on Empirical Methods in Natural Language Processing, pages 388–395, Barcelona, Spain. Association for Computational Linguistics.	
726			
727			
728	Faisal Ladhak, Esin Durmus, Claire Cardie, and Kathleen McKeown. 2020. WikiLingua: A new benchmark dataset for cross-lingual abstractive summarization . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 4034–4048, Online. Association for Computational Linguistics.		
729			
730			
731			
732			
733			
734	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. <i>ACM Transactions on Asian Language Information Processing (TALIP)</i> , 2(3):245–269.		
735			
736			
737			
738			
739	Mike Lewis, Marjan Ghazvininejad, Gargi Ghosh, Armen Aghajanyan, Sida Wang, and Luke Zettlemoyer. 2020. Pre-training via paraphrasing. In <i>Proceedings of the 34th International Conference on Neural Information Processing Systems, NIPS’20</i> , Red Hook, NY, USA. Curran Associates Inc.		
740			
741			
742			
743			
744			
745	Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries . In <i>Text Summarization Branches Out</i> , pages 74–81, Barcelona, Spain. Association for Computational Linguistics.		
746			
747			
748			
749	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. <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.		
750			
751			
752			
753			
754			
755	Alexandra Luccioni and Joseph Viviano. 2021. What’s in the box? an analysis of undesirable content in the Common Crawl corpus . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)</i> , pages 182–189, Online. Association for Computational Linguistics.		
756			
757			
758			
759			
760			
761			
762			
763	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.		
764			
765			
766			
767			
768			
769	Mark F. Medress, Franklin S Cooper, Jim W. Forgie, CC Green, Dennis H. Klatt, Michael H. O’Malley, Edward P Neuburg, Allen Newell, DR Reddy, B Ritea, et al. 1977. Speech understanding systems: Report of a steering committee. <i>Artificial Intelligence</i> , 9(3):307–316.		
770			
771			
772			
773			
774			
775	Khanh Nguyen and Hal Daumé III. 2019. Global Voices: Crossing borders in automatic news summarization . In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> , pages 90–97, Hong Kong, China. Association for Computational Linguistics.		
776			
777			
778			
779			
780			
	Constantin Orasan and Oana Andreea Chiorean. 2008. Evaluation of a cross-lingual romanian-english multi-document summariser. In <i>Proceedings of the Sixth International Conference on Language Resources and Evaluation (LREC’08)</i> , Marrakech, Morocco. European Language Resources Association (ELRA). Http://www.lrec-conf.org/proceedings/lrec2008/ .		781 782 783 784 785 786 787
	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation . In <i>Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics</i> , pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.		788 789 790 791 792 793 794
	Laura Perez-Beltrachini and Mirella Lapata. 2021. Models and datasets for cross-lingual summarisation . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 9408–9423, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		795 796 797 798 799 800
	Sascha Rothe, Shashi Narayan, and Aliaksei Severyn. 2020. Leveraging pre-trained checkpoints for sequence generation tasks . <i>Transactions of the Association for Computational Linguistics</i> , 8:264–280.		801 802 803 804
	Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks . In <i>Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS 2014)</i> , pages 3104–3112, Montreal, Canada.		805 806 807 808 809
	Chau Tran, Yuqing Tang, Xian Li, and Jiatao Gu. 2020. Cross-lingual retrieval for iterative self-supervised training . In <i>Advances in Neural Information Processing Systems</i> , volume 33, pages 2207–2219. Curran Associates, Inc.		810 811 812 813 814
	Daniel Varab and Natalie Schluter. 2021. MasiveSumm: a very large-scale, very multilingual, news summarisation dataset . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10150–10161, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		815 816 817 818 819 820 821
	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 <i>Proceedings of the 31st International Conference on Neural Information Processing Systems (NIPS 2017)</i> , page 6000–6010, Long Beach, California, USA.		822 823 824 825 826 827 828
	Xiaojun Wan, Huiying Li, and Jianguo Xiao. 2010. Cross-language document summarization based on machine translation quality prediction . In <i>Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics</i> , pages 917–926, Uppsala, Sweden. Association for Computational Linguistics.		829 830 831 832 833 834
	Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim		835 836

837	Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al.	Pierre Zweigenbaum, Serge Sharoff, and Reinhard Rapp.	895
838	2016. Google’s neural machine translation system: Bridging the gap between human and machine translation . <i>arXiv:1609.08144</i> .	2017. Overview of the second bucc shared task: Spotting parallel sentences in comparable corpora. In <i>Proceedings of the 10th Workshop on Building and Using Comparable Corpora</i> , pages 60–67.	896
839			897
840			898
841	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021a. mT5: A massively multilingual pre-trained text-to-text transformer . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.		899
842			
843			
844			
845			
846			
847			
848			
849	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021b. mT5: A massively multilingual pre-trained text-to-text transformer . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 483–498, Online. Association for Computational Linguistics.		
850			
851			
852			
853			
854			
855			
856			
857	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 bidirectional dual encoder with additive margin softmax . In <i>Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19</i> , pages 5370–5378. International Joint Conferences on Artificial Intelligence Organization.		
858			
859			
860			
861			
862			
863			
864			
865			
866	Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with bert . In <i>International Conference on Learning Representations</i> .		
867			
868			
869			
870	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 <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 563–578, Hong Kong, China. Association for Computational Linguistics.		
871			
872			
873			
874			
875			
876			
877			
878			
879			
880	Junnan Zhu, Haoran Li, Tianshang Liu, Yu Zhou, Jijun Zhang, and Chengqing Zong. 2018. Msmo: Multimodal summarization with multimodal output . In <i>Proceedings of the 2018 conference on empirical methods in natural language processing</i> , pages 4154–4164.		
881			
882			
883			
884			
885			
886	Junnan Zhu, Qian Wang, Yining Wang, Yu Zhou, Jijun Zhang, Shaonan Wang, and Chengqing Zong. 2019. NCLS: Neural cross-lingual summarization . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)</i> , pages 3054–3064, Hong Kong, China. Association for Computational Linguistics.		
887			
888			
889			
890			
891			
892			
893			
894			

Appendix

A Aligning Summaries using LaBSE

In Section 2, we curate CrossSum by aligning identical summaries in different languages. It might be argued why the articles themselves were not used for the alignment process. Initially, we experimented with whole-article embeddings. However, this resulted in many false-negative alignments, where similarity scores between identical articles across languages were relatively low (verified manually between English and the authors’ native languages). This is most likely attributed to the 512-token limit of LaBSE and different sequence lengths of those articles due to different languages having different subword segmentation fertility (Ács, 2019). This would entail that identical articles in different languages might be truncated at different locations, resulting in discrepancies between their embeddings. As observed in the BUCC evaluation, LaBSE is well-suited for sentence-level retrieval. Since summaries are good representatives of entire articles, we finally chose summaries as our candidates for the alignment.

B Inter-annotator Agreement of Human Evaluation

Language Pair	Cohen’s Kappa
Arabic-English	0.82
Chinese-English	0.73
Indonesian-English	0.73
Bengali-English	0.73
Urdu-English	0.76
Punjabi-English	0.71
Swahili-English	0.78
Pashto-English	0.75

Table 2: Language pair-wise kappa scores.

C Statistical Significance

While the scores obtained from the experiments in Section 5 are a telling sign that the proposed m2m model performs better than the others, the differences are very close in many language pairs. Therefore, a statistical significance test is still warranted to support our claim further. As such, for each language pair experimented on, we performed the Bootstrap resampling test (Koehn, 2004) with the m2m model against the best performing model among the others in a one vs. all manner: if m2m

has the best score (ROUGE-2/LaSE), we compare it with the model with the second-best score, and if m2m is not the best, we compare it with the best. Results ($p < 0.05$) reveal that in more than 42% language pairs tested, m2m is significantly better, and in less than 10% pairs, it is considerably worse. Details presented in Table 3.

Pivot	Metric	Better	Worse	Insignificant
x-en	R-2/LaSE	8/18	2/2	25/15
en-x	R-2/LaSE	20/15	3/14	12/6
x-zh	R-2/LaSE	11/13	0/0	23/21
zh-x	R-2/LaSE	17/12	1/2	16/20
x-hi	R-2/LaSE	18/15	1/6	15/13
hi-x	R-2/LaSE	19/15	0/6	15/13
x-ar	R-2/LaSE	6/15	2/3	26/16
ar-x	R-2/LaSE	23/15	1/5	10/14
x-ru	R-2/LaSE	6/11	2/7	26/16
ru-x	R-2/LaSE	19/13	2/7	13/14

Table 3: Significance test on different pivot languages.

D Modeling Details

D.1 Choice of Pretrained Model

Many pretrained multilingual text-to-text models are currently available, e.g., mBART (Liu et al., 2020), CRISS (Tran et al., 2020), MARGE (Lewis et al., 2020), and mT5 (Xue et al., 2021b). While mBART and mT5 are pretrained with multilingual objectives, CRISS and MARGE are pretrained with a cross-lingual one, which better suits our use case. However, we choose mT5 for fine-tuning because of its broad coverage of 101 languages with support for 41 of the 45 languages from CrossSum, in contrast to only 15 languages in mBART or CRISS and 26 in MARGE.

D.2 Summarize-then-translate (s. + t.)

The primary reason for using summarize-then-translate rather than translate-then-summarize is the computational cost between these two. Available translation models only work for short sequences and are unsuitable for long documents. One solution is to segment the documents into sentences and then translate them. But that increases the compute overhead, and translations suffer from loss of context. We use a multilingual summarization model (Hasan et al., 2021) coupled with the multilingual machine translation model, M2M-100 (Fan et al., 2021), for our pipeline.

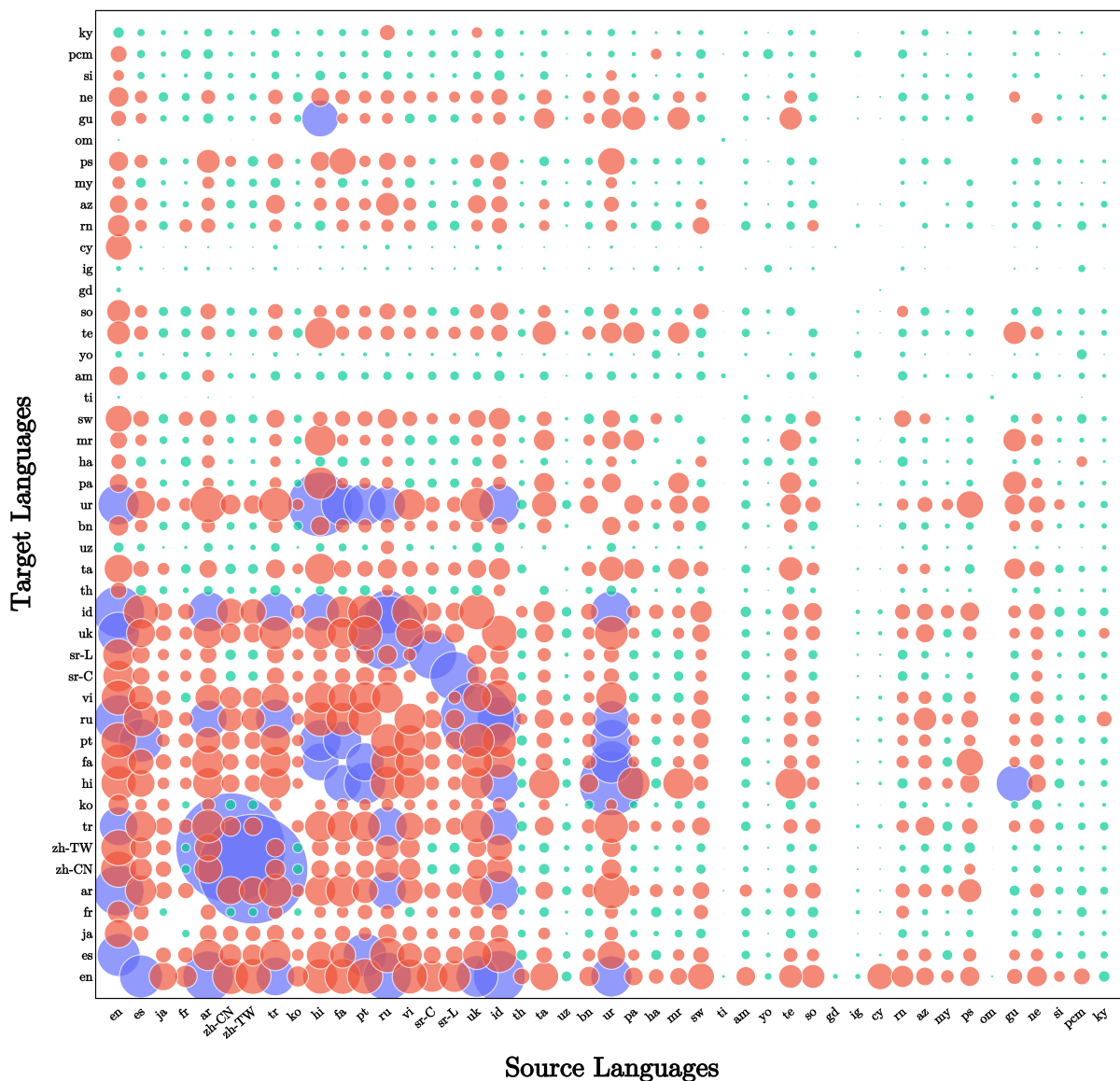


Figure 6: A bubble plot depicting the article-summary frequencies of CrossSum. The radii of the bubbles are proportional to the number of samples for the corresponding language pair (exact numbers are in Table 4). Languages are ordered by the language taxonomy from Joshi et al. (2020). To show better contrast between language pairs, we color a bubble cyan if its frequency is below 500 (1218 pairs), red for 500 to 5000 (688 pairs), and blue for frequencies exceeding 5000 (52 pairs).

D.2.1 Multilingual Summarization

The pipeline first performs in-language summarization (the language of the summary is the same as that of its input article) and then translates the summary into the desired target language. We train our own model for summarization as the model released by Hasan et al. (2021) has been rendered unusable due to the change in the dataset split. We extend our component graphs to curate the in-language dataset splits. We consider articles having no identical counterpart in any other language as

single node components in the component graph. As before, we assign all articles originating from a single component to the training (dev/test) set of the dataset, extending them to the in-language splits too. We then train the multilingual model by fine-tuning mT5 with the in-language splits, sampling each batch of 256 samples from a single language with a sampling factor of $\alpha = 0.5$.

D.2.2 Multilingual Translation

For multilingual translation, we used M2M-100 (Fan et al., 2021) (418M parameters variant), a

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

990

991

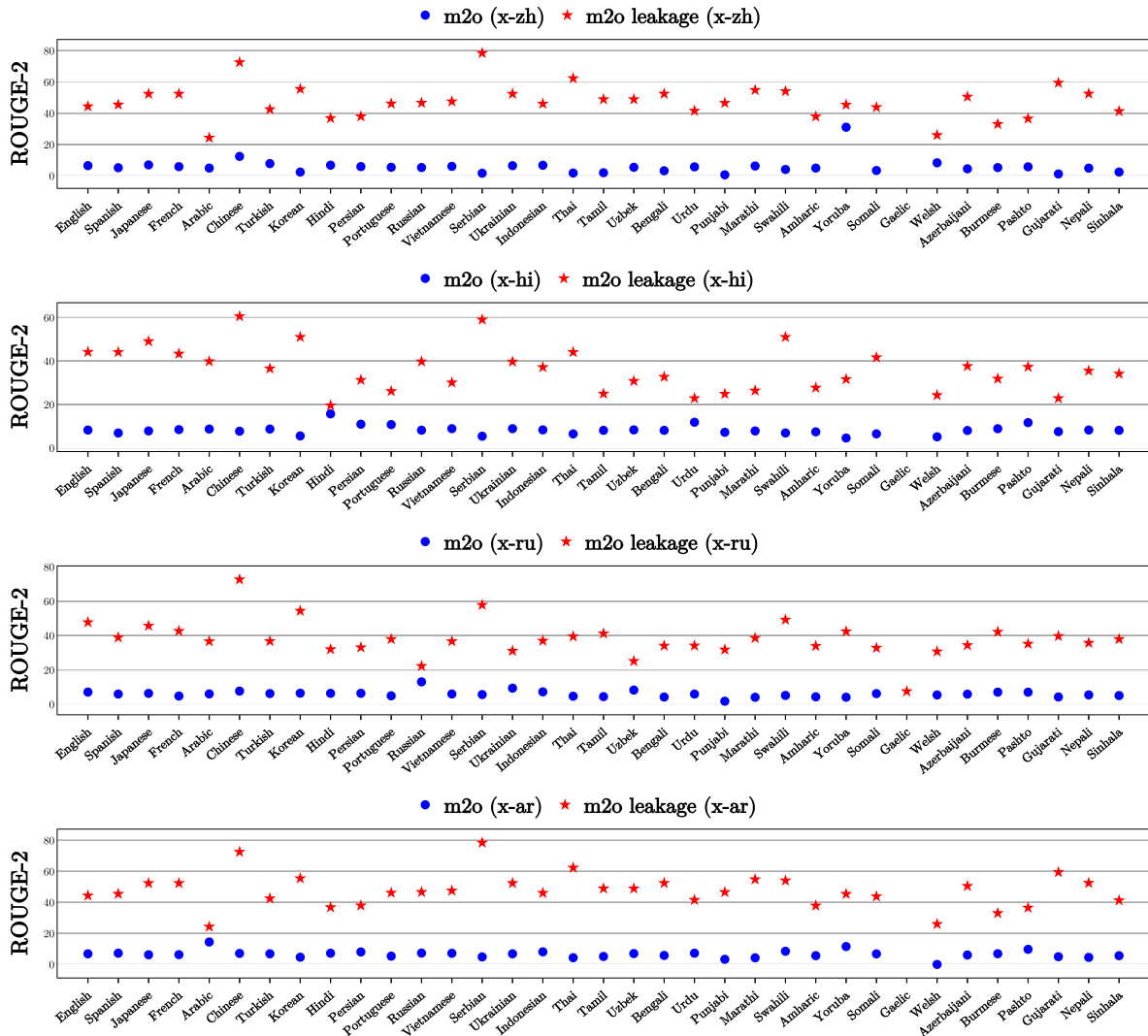


Figure 7: Training on the dataset respecting the original XL-Sum splits causes absurdly high ROUGE scores (marked red) in many-to-one models due to implicit data leakage. Therefore, we split taking the issue into account and consequently, models trained on the new set (marked blue) do not exhibit any unusual spike in ROUGE-2.

many-to-many multilingual translation model, with support for 37 languages from CrossSum.

D.3 Many-to-One (m2o) Model

Many-to-one training is standard for evaluating cross-lingual summarization. In these models, the language of the source text can vary, but the target language remains the same, i.e., as the pivot language. Instead of sampling all samples of a batch from the same language pair, we sample 8 mini-batches of 32 samples using a sampling factor of $\alpha = 0.25$, the source side of each originating from a single language while the target language remains fixed. We then merge the mini-batches into a single batch and update the model parameters. This is to ensure that there are not many duplicates in a single batch (if all 256 samples of a batch are sampled

from a single language pair, there might be many duplicates as many language pairs do not have 256 training samples) and the model still benefits the advantages of low-resource upsampling.

D.4 One-to-many (o2m) Model

o2m models are complementary to m2o models: we train them by keeping the source language fixed and varying the target language. We upsample the low-resource target languages with the same sampling factor of $\alpha = 0.25$ and merge 8 mini-batches of 32 samples each, analogous to m2o models.

D.5 Many-to-many (m2m) Multistage Model

This is the model obtained from the Algorithm 1. In contrast to standard language sampling (Conneau et al., 2020), we sample the target language and

1023	then choose the source based on that decision. We	zero-shot cross-lingual generation (Duan et al.,	1069
1024	use batch size 256, 8 mini-batches with size 32 and	2019) without relying on any labeled examples.	1070
1025	$\alpha = 0.5, \beta = 0.75$.	To this end, we fine-tuned mT5 with the in-	1071
1026	D.6 Many-to-many (m2m) Unistage Model	language (both source and target are in the same	1072
1027	This algorithm is similar to standard language sam-	language) samples only in a multilingual fashion	1073
1028	pling, the difference being that languages are sam-	and, during inference, varied the target language.	1074
1029	pled as pairs from all possible combinations. In-	Unfortunately, the model totally fails at generating	1075
1030	stead of sampling one language pair at each training	cross-lingual summaries and performs in-language	1076
1031	step, we sample 8 pairs, one for each mini-batch	summarization instead.	1077
1032	of size 32. We then merge the mini-batches into	We also fine-tuned m2o models in a zero-shot	1078
1033	a single batch of 256 samples before updating the	setting (with only the in-language samples of the	1079
1034	model parameters. We use a sampling factor of	target language) in a monolingual fashion. Here,	1080
1035	$\alpha = 0.25$.	the models are able to generate non-trivial sum-	1081
1036	In all models, we discarded a language pair from	maries for some language pairs but still lag behind	1082
1037	training if it had fewer than 30 training samples to	fully supervised models by a significant margin	1083
1038	prevent too many duplicates in a mini-batch. The	(Figure 10 and 11).	1084
1039	training was done together with the in-language	Furthermore, we ran inference with the m2m	1085
1040	samples.	model on distant low-resource language pairs that	1086
1041	E Experimental Details	were absent during training. Their LaSE scores	1087
1042	E.1 Training Setups	were substantially below supervised pairs, mean-	1088
1043	Fine-tuning generation models is compute-	ing zero-shot transfer in supervised multilingual	1089
1044	intensive, and due to computational limitations,	models (Johnson et al., 2017) shows weak perfor-	1090
1045	we fine-tune all pretrained models for 25k steps	mance as well.	1091
1046	with an effective batch size of 256, which roughly	We do not perform any few-shot experiments	1092
1047	takes about three days on a 4-GPU NVIDIA P100	and leave them as potential future directions.	1093
1048	server. We use the base variant of mT5, having	F Ablation Studies	1094
1049	250k vocabulary, 768 embedding and dimension	We make several design choices in the multistage	1095
1050	size, 12 attention heads, and 2048 FFN size, with	sampling algorithm. We break them into two main	1096
1051	580M parameters. We limit the input to 512 and	decisions:	1097
1052	output to 84 tokens. All models are trained on the	1. Making mini-batches and sampling the lan-	1098
1053	respective subsets of the CrossSum training set.	guage pair for each mini-batch.	1099
1054	E.2 Inference	2. Keeping either the source or the target lan-	1100
1055	During inference, we jump-start the decoder with	guage fixed for each batch.	1101
1056	language-specific BOS (beginning of sequence) to-	To verify that these choices indeed affect perfor-	1102
1057	kens (Johnson et al., 2017) at the first decoding step	mance positively, we train five different models for	1103
1058	for guiding the decoder to generate summaries in	ablation:	1104
1059	the intended target language. We use beam search	1. Sampling the language pair in mini-batches	1105
1060	(Medress et al., 1977) with the beam size 4 and use	in one stage only and then merging them into	1106
1061	a length penalty (Wu et al., 2016) of 0.6. We limit	large batches before updating model parame-	1107
1062	ourselves only to the languages supported by mT5,	ters: m2m-unistage.	1108
1063	fastText, and M2M-100.	2. Sampling the language pair with large batches	1109
1064	E.3 Zero-shot Cross-lingual Transfer	of 256 samples without mini-batching: m2m-	1110
1065	The previous experiments were done in a fully su-	large.	1111
1066	perervised fashion. However, for many low-resource	3. Multistage sampling keeping only the target	1112
1067	language pairs, samples are not abundantly avail-	language fixed in a batch: m2m-tgt [<i>our pro-</i>	1113
1068	able. Hence, it is attractive to be able to perform	<i>posed model</i>].	1114

- 1115 4. Multistage sampling keeping only the source
1116 language fixed in a batch: m2m-src; i.e., the
1117 complement of our proposed model.
- 1118 5. Multistage sampling keeping either the source
1119 or the target language fixed (with equal proba-
1120 bility) for each batch: m2m-src-tgt.

1121 We benchmark on all the language pairs done
1122 previously and show the mean ROUGE-2 and LaSE
1123 scores in Table 5.

1154 with our hypothesis made in Footnote 9, as m2m-
1155 src and m2m-tgt mimic the training settings of the
1156 o2m and m2o models, respectively, at the batch
1157 level. The m2m-src-tgt is the middle ground be-
1158 tween m2m-src and m2m-tgt and, likewise, scores
1159 between these two. In our opinion, the perfor-
1160 mance dynamics between the m2o (m2m-tgt) and
1161 o2m (m2m-src) models is an interesting finding
1162 and should be studied in depth as a new research
1163 direction in future works.

Model	Scores		Significance		
	R-2/LaSE	Better	Worse	Insignificant	
m2m-large	8.31/57.45	122	59	503	
m2m-unistage	7.51/55.36	191	149	344	
m2m-tgt	8.15/57.15	289	66	329	
m2m-src	4.44/26.75	34	477	173	
m2m-src-tgt	6.47/42.55	89	297	298	

Table 5: ROUGE-2 and LaSE scores for ablation.

1124 As can be seen from the table, m2m-large, the
1125 standard m2m model, has the best average ROUGE-
1126 2/LaSE scores among all m2m variants. This begs
1127 the question of whether our proposed multistage
1128 sampling is, after all, needed or not. But the scores
1129 of the proposed m2m-tgt model does not fall much
1130 below. Therefore, we show statistical significance
1131 test results of all m2m models, comparing them
1132 against m2o, o2m, and s.+t. in one vs. all manner.

1133 Significance results paint a different picture:
1134 m2m-tgt triumphs over all other models, getting
1135 significantly better results on 42% language pairs,
1136 more than double the m2m-large model. We in-
1137 spected the results individually and found that the
1138 results are notably better on language pairs that are
1139 not adequately represented in the training set. m2m-
1140 tgt performs comparatively worse on high-resource
1141 language pairs, which we believe is a fair compro-
1142 mise to uplift low-resource ones. As m2m-large
1143 can sample a pair only once per batch, it fails to
1144 incorporate many language pairs due to them hav-
1145 ing insufficient participation during training. On
1146 the other hand, our proposed multistage sampling
1147 algorithm performs well in this regard by sampling
1148 in two stages.

1149 While m2m-tgt outperforms all the rest, m2m-
1150 src falls behind all other models by a large margin.
1151 This phenomenon also has the same trend as the
1152 results in Section 5, where o2m models failed at
1153 generating cross-lingual summaries. This is in line

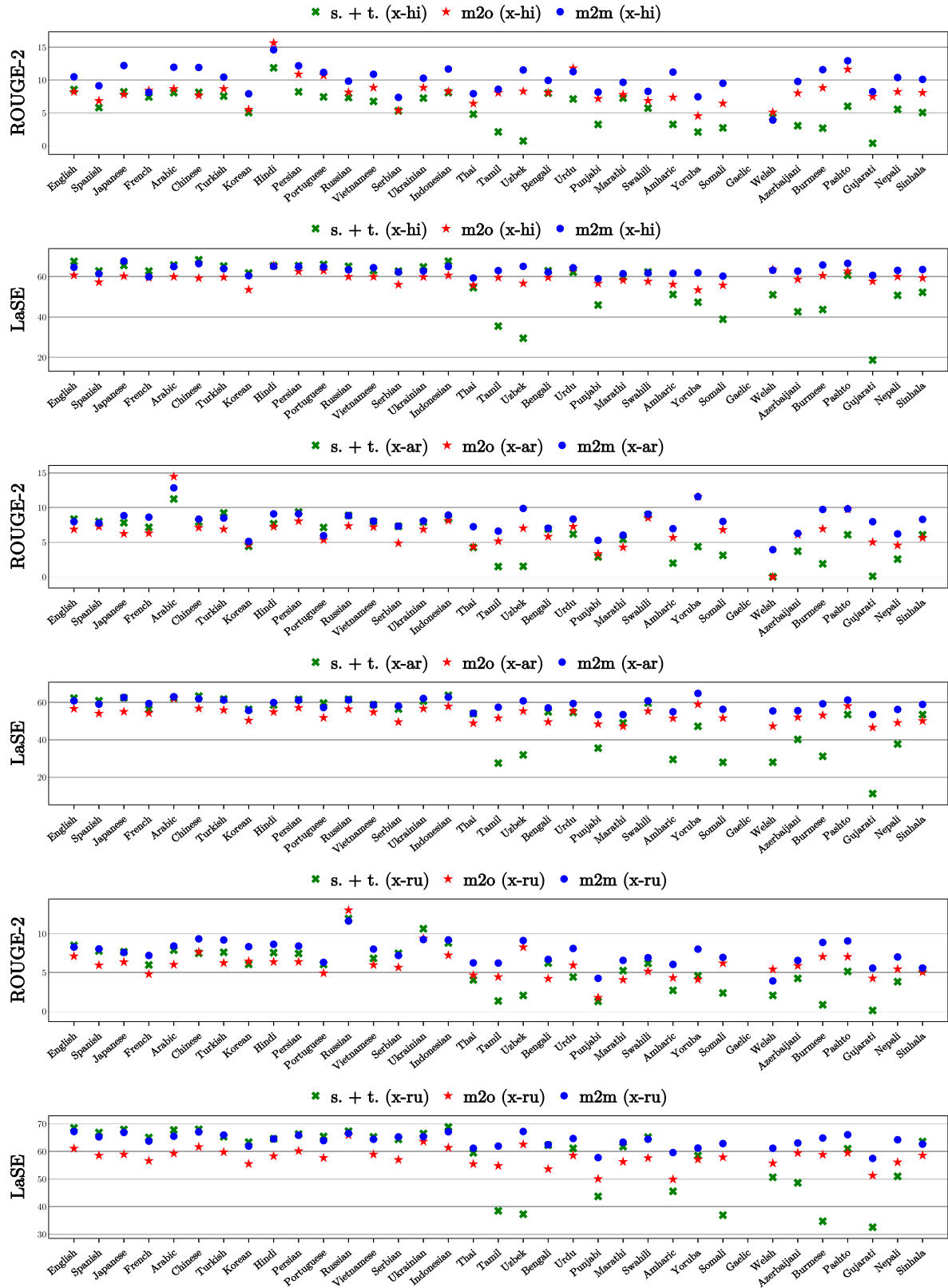


Figure 8: ROUGE-2 and LaSE scores for Hindi, Arabic, and Russian as target pivots as the sources languages vary. Just like Figure 4, the m2m model significantly outperforms the m2o models and s. + t. baseline on most languages.

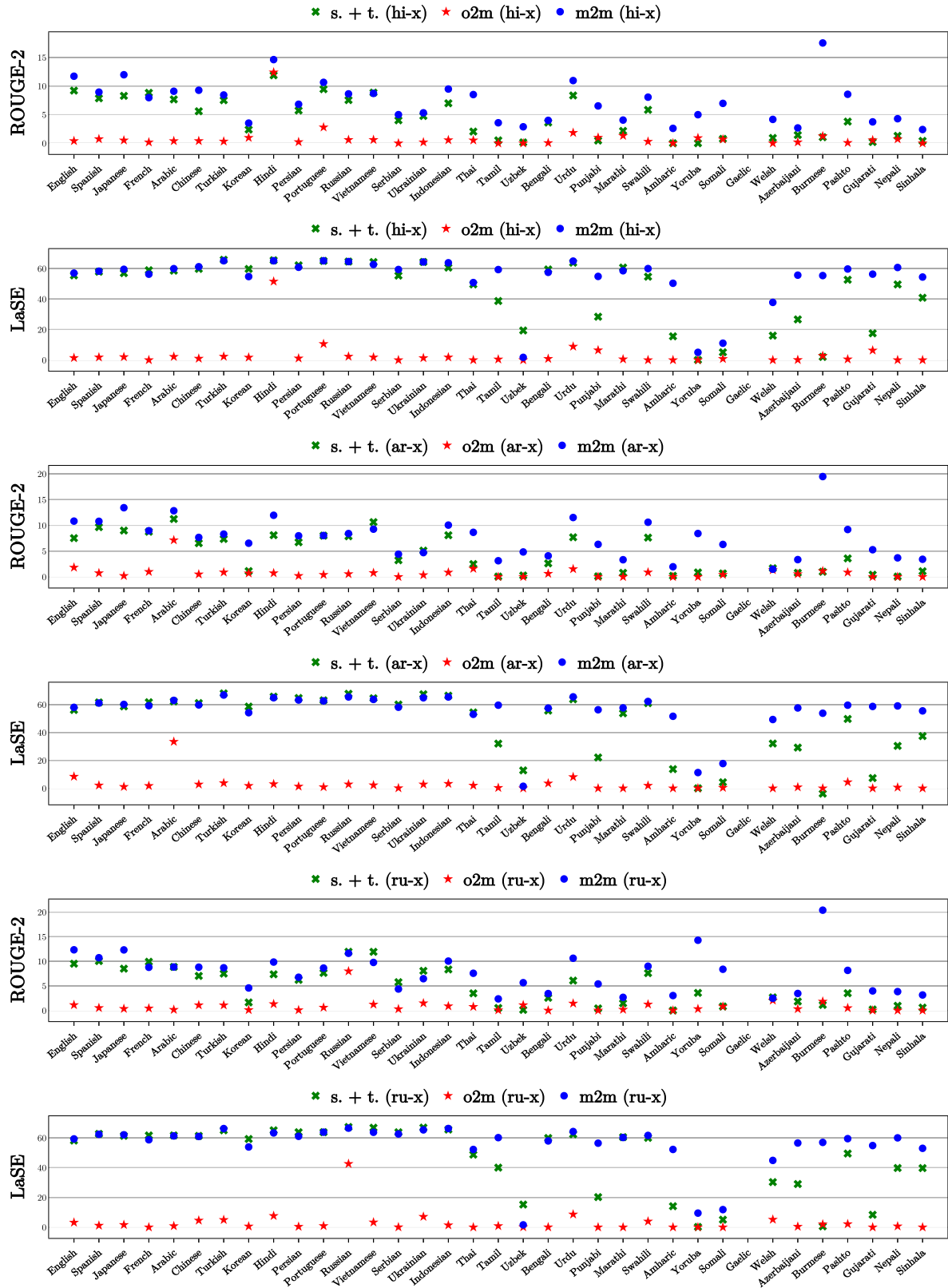


Figure 9: ROUGE-2 and LaSE scores for Hindi, Arabic, and Russian as source pivots as the target languages vary. Just like Figure 5, the m2m model significantly outperforms the o2m models and s. + t. baseline on most languages.

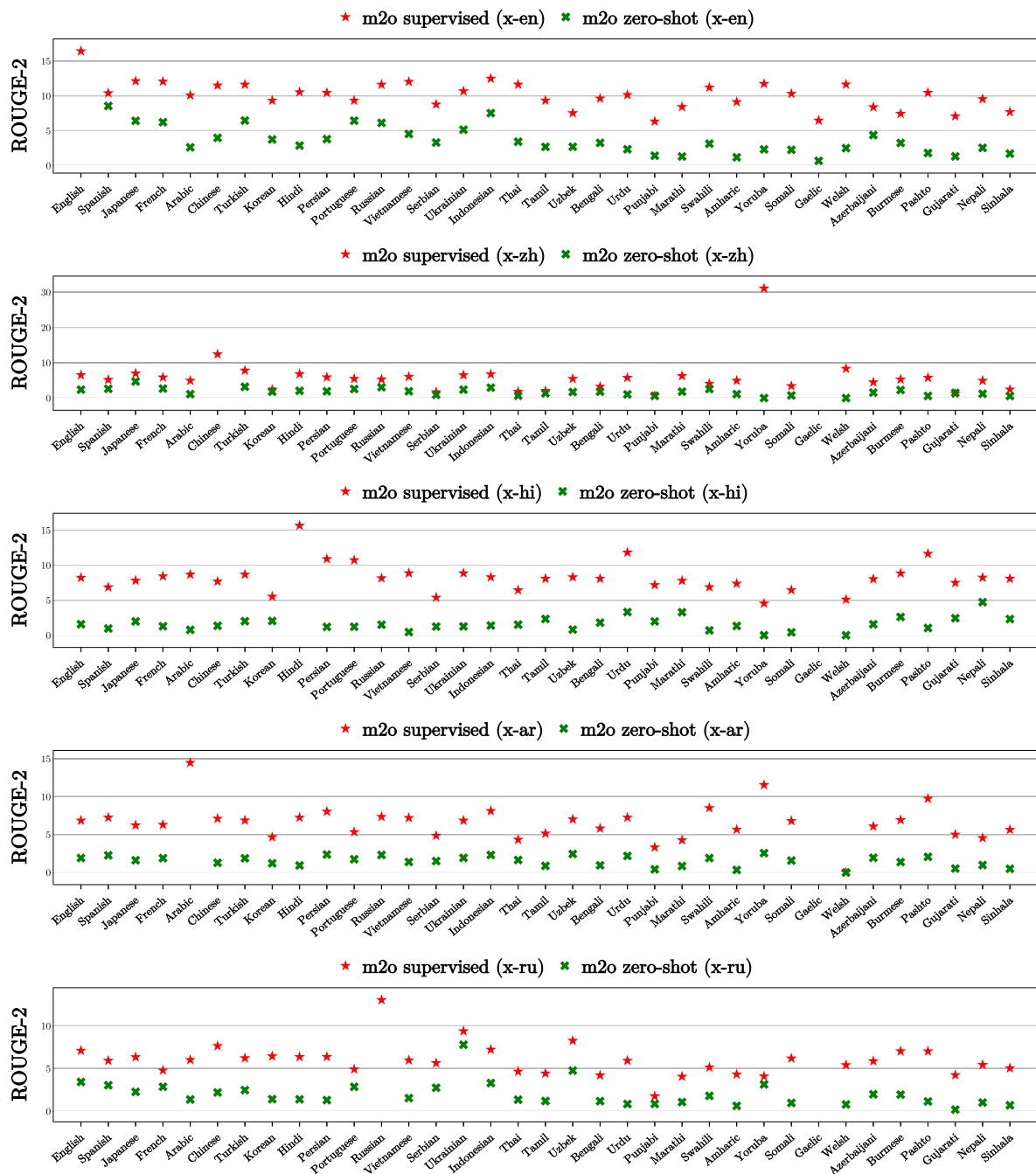


Figure 10: Zero-shot ROUGE-2 scores for the different target languages as the source languages vary. The zero-shot models are trained with only the in-language samples of the pivot. Though their results are clearly behind the fully supervised models, the zero-shot models are able to generate non-trivial summaries for many language pairs.

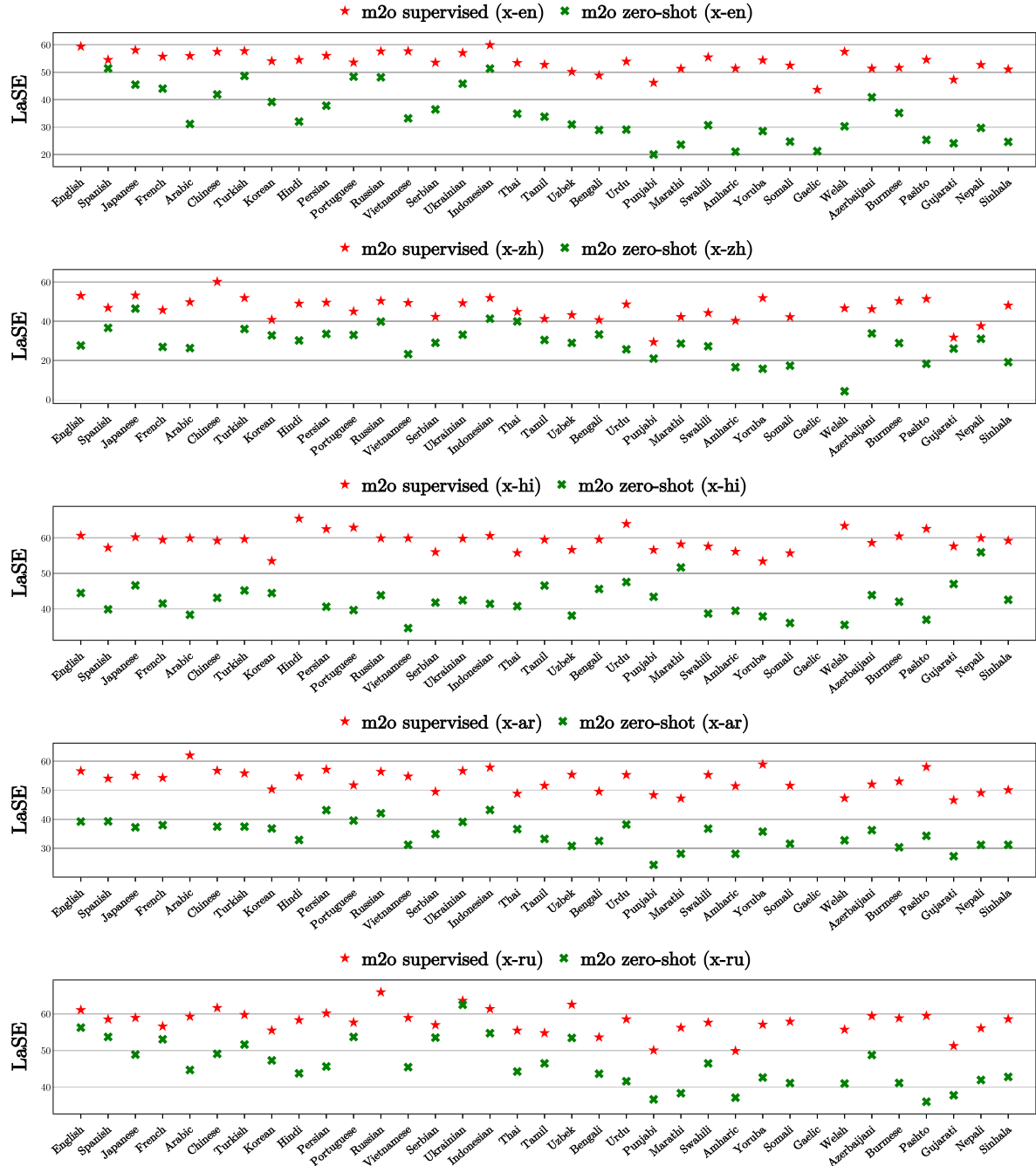


Figure 11: Zero-shot LaSE scores for the different source languages as the target languages vary. The zero-shot models are trained with only the in-language samples of the pivot. Though their results are clearly behind the fully supervised models, the zero-shot models are able to generate non-trivial summaries for many language pairs.

