

You Are What You Train: Effects of Data Composition on Training Context-aware Machine Translation Models

Anonymous ACL submission

Abstract

Achieving human-level translations requires leveraging context to ensure coherence and handle complex phenomena like pronoun disambiguation. Sparsity of contextually rich examples in the standard training data has been hypothesized as the reason for the difficulty of context utilization. In this work, we systematically validate this claim in both single- and multilingual settings by constructing training datasets with a controlled proportions of contextually relevant examples. We demonstrate a strong association between training data sparsity and model performance confirming sparsity as a key bottleneck. Importantly, we reveal that improvements in one contextual phenomenon do not generalize to others. While we observe some cross-lingual transfer, it is not significantly higher between languages within the same sub-family. Finally, we propose and empirically evaluate two training strategies designed to leverage the available data. These strategies improve context utilization, resulting in accuracy gains of up to 6 and 8 percentage points on the ctxPro evaluation in single- and multilingual settings respectively.¹

1 Introduction

Context-Aware Machine Translation (MT) models use surrounding sentences (context) to improve translation by maintaining coherence and resolving ambiguities (Agrawal et al., 2018; Bawden et al., 2018; Müller et al., 2018; Voita et al., 2019b). The context can be sentences in the source language and the previously translated sentences in the target language. While many works improved the translation quality of the context-aware MT by applying standard Transformer (Vaswani et al., 2017) model (Sun et al., 2022; Majumde et al., 2022; Gete et al., 2023; Post and Junczys-Dowmunt, 2024; Alves



Figure 1: Composition of the English-to-German training datasets with the Gender phenomenon in Pure IWSLT and IWSLT+OpenSubtitles settings. Annotations are based on ctxPro (Wicks and Post, 2023), and the dashed bars represent the contextually-rich datasets. Note that the horizontal axis starts at 100,000.

et al., 2024; Kocmi et al., 2024), specialized architectures (Tu et al., 2017; Bawden et al., 2018; Miculicich et al., 2018; Maruf et al., 2019; Huo et al., 2020; Zheng et al., 2021), and decoder-only LLMs (Alves et al., 2024; Kocmi et al., 2024), the reason why the context utilization is challenging for the models remain an open question.

The low density of contextually rich (requiring context for correct translation) examples in the training datasets has been suspected as the main reason why MT models have trouble in translating contextual phenomena. For example, Lupo et al. (2022) proposed the two-fold sparsity hypothesis, where the low density of examples in the dataset and the tokens in the examples requiring context increases the difficulty of learning to leverage context. Post and Junczys-Dowmunt (2024) show that sparsity in the evaluation datasets makes it difficult to assess the context utilization of the models. We argue that this also points to the sparsity hypothesis in the training data, as the evaluation datasets are

¹<https://anonymous.4open.science/r/data-composition-0C80>.

often subsets of the training datasets.

In this work, we evaluate how the the proportion of contextually rich examples in the training data of the encoder-decoder context-aware MT models affects the overall translation quality measured by BLEU (Papineni et al., 2002) and COMET (Rei et al., 2020), and performance on the examples requiring context (using generative and contrastive evaluations). To this end, we use ctxPro toolset (Wicks and Post, 2023) to extract the relevant examples containing the following phenomena: Gender, Formality, Auxiliary, Inflection, and Animacy. We refer the reader to the original ctxPro paper for the details of the annotation and phenomena (Wicks and Post, 2023). We constructed training data by mixing contextually rich and poor examples with varying proportions (Figure 1 illustrates this for Gender in English-to-German). Moreover, we evaluate cross-lingual transfer of context utilization in multilingual models on English-to-X and X-to-English where X is {German, French, Polish, Russian, and Spanish}. Finally, we explore several ways to effectively leverage the available data to obtain models that perform well both generally and in context-sensitive settings. The contributions of this work are:

1. We **empirically validate the sparsity hypothesis**, showing strong relation between the density of the contextual phenomena in the training data and the resulting performance of the context-aware MT models.
2. We **reveal limitations in generalization**, showing that the improvement in one linguistic phenomenon does not transfer to others. We observe limited cross-lingual transfer, not substantially higher between languages in the same sub-family.
3. We propose and empirically evaluate **two training strategies** designed to improve context utilization by leveraging the available data. We show a trade-off between improving context utilization and general translation metrics such as BLEU.

2 Related Work

Through years many dedicated architectures has been proposed for context-aware MT (Miculicich et al., 2018; Voita et al., 2019b,a; Bao et al., 2021; Chen et al., 2022; Feng et al., 2022; Bulatov et al., 2022; Maka et al., 2024) including popular multi-encoder (where a separate encoder is responsible for processing the context sentences; Jean et al.,

2017; Miculicich et al., 2018; Maruf et al., 2019; Huo et al., 2020; Zheng et al., 2021), but the standard Transformer model (Vaswani et al., 2017) with the sentences being concatenated (single-encoder; Tiedemann and Scherrer, 2017; Ma et al., 2020; Zhang et al., 2020) exhibited high performance despite its relative simplicity (Majumde et al., 2022; Sun et al., 2022; Gete et al., 2023; Post and Junczys-Dowmunt, 2023). While decoder-only LLMs have achieved state-of-the-art results in MT (Alves et al., 2024; Kocmi et al., 2024), they require extensive datasets for training and have a large number of parameters that can limit their usefulness in computationally constrained environments (e.g., edge devices) and low-resource settings. In recent years, research interest in the architectures other than decoder-only has remained relevant (Mohammed and Niculae, 2024; Warner et al., 2024; Alastruey et al., 2024; Azeemi et al., 2025; Marashian et al., 2025; Breton et al., 2025; Clavié et al., 2025). Therefore, we focus this paper on encoder-decoder models.

The standard sentence-level metrics (e.g., BLEU (Papineni et al., 2002) do not capture the contextual utilization by the models (Hardmeier, 2012; Wong and Kit, 2012). To address this, several evaluation datasets have been proposed including contrastive (Müller et al., 2018; Bawden et al., 2018; Voita et al., 2019b; Lopes et al., 2020) and generative such as ctxPro (Wicks and Post, 2023) used in this study. Moreover, metrics like CXMI (Fernandes et al., 2021) and PCXMI (Fernandes et al., 2023) can measure how much the model relies on context during translation.

The effects of the training dataset on the final model has also been studied extensively (Kaplan et al., 2020; Hoffmann et al., 2022) in different domains (Alabdulmohsin et al., 2023), including document-level MT (Zhuocheng et al., 2023). The studies mostly concentrated on the scale of the training dataset. We, instead, investigate the composition of the dataset and its effect on the context-aware MT models.

Several works proposed methods increasing contextual capabilities of the models by training the models on annotated data (Jwalapuram et al., 2020; Yin et al., 2021; Mağa et al., 2025) but they target only pronoun disambiguation. Fine-tuning in this case can be seen as similar to domain adaptation (Luong and Manning, 2015; Chu et al., 2017) where loss weighting (similar to one of our methods) is an effective strategy (Wang et al., 2017).

3 Effects of Data Composition

We first measured how the presence of contextually rich examples in the training data affects both translation quality and the models’ ability to leverage context. To that end, we trained models on datasets whose composition we systematically varied. Specifically, we identified contextual examples (containing relevant phenomena) from the available datasets using ctxPro toolset (Wicks and Post, 2023) and constructed a series of datasets with varying densities of different phenomena. This setup allowed us to assess inter-phenomena as well as cross-lingual effects of the composition of the training datasets. We used two settings: single language pair (English-to-German) and multilingual. For the multilingual setting, we used English-to-X and X-to-English language directions, where X is {German, French, Polish, Russian, and Spanish} - a subset of directions covered by the ctxPro. We utilized two Germanic, Romance, and Slavic languages.

3.1 Datasets

We base our research on two document-level translation datasets: IWSLT 2017 English-to-German (Cettolo et al., 2017) and OpenSubtitles 2018 (Lison et al., 2018). For the English-to-German direction, we employ both datasets, and for the multilingual setting, we only use OpenSubtitles. We extract contextual annotations from the training subset of the IWSLT dataset using the ctxPro toolset. The annotated (containing contextually-rich examples) subset forms **IWSLT-dense** dataset, which can be further divided based on the target phenomenon: Gender, Formality, Auxiliary, Inflection, and Animacy. From the remaining examples we form **IWSLT-sparse** dataset of size 123,000, containing no examples annotated with contextual phenomena. CtxPro released annotations extracted from the OpenSubtitles 2018 dataset divided into *dev*, *devtest*, and *test* subsets. We set aside the *test* subset for the evaluation and used the combined *dev* and *devtest* subsets for training, forming **OS-dense** dataset. The released ctxPro dataset is not exhaustive; therefore, we do not create the sparse version of the OpenSubtitles dataset. Instead, we randomly sample the OpenSubtitles dataset to the desired size (referred to as **OS-random**). In Appendix A we present the sizes of the dense component datasets.

To create the training datasets with varying densities of contextually rich examples, we sample and

concatenate examples from both dense and sparse datasets to form a training dataset. For English-to-German, we study two settings: *Pure IWSLT* (only IWSLT-sparse and IWSLT-dense datasets) and *IWSLT + OS* (using IWSLT-sparse, IWSLT-dense, and English-to-German OS-rand and OS-dense datasets). These allow us to study two regimes: extremely low sparsity with the first setting, and very dense with the second one. We progressively replace examples from sparse and random datasets with the examples sampled from dense datasets. In the multilingual experiments, we formed the baseline training dataset by sampling 50,000 examples from OS-rand for all language directions we considered. For each phenomenon in a language direction, we formed the enriched datasets by replacing n examples with the examples sampled from the OS-dense dataset corresponding to the phenomenon and language direction. We chose n to be the minimum number of examples (rounded) for a particular phenomenon across language directions maximizing the resulting density of the training datasets while making the results comparable between language directions. We present the illustration of the composition of the datasets in Figure 1 for Gender on English-to-German and further details in Appendix A.

3.2 Training

We employed a two-stage training process where first the sentence-level model is trained on more abundant sentence-aligned datasets, followed by the context-aware training on the document-level dataset. Following Mała et al. (2025), we rely on the publicly available pre-trained sentence-level models, namely *OPUS-MT en-de* and *No Language Left Behind* (NLLB-200) with 600M parameters (NLLB Team et al., 2022). We concatenate consecutive sentences separated by the special [SEP] token on both the source and target sides. Similar to Sun et al. (2022), we create examples with all context sizes (number of previous sentences to concatenate) from zero to the maximum context size. We set the maximum context size to three as further increases have shown diminishing returns regarding context utilization (Post and Junczys-Dowmunt, 2023). During inference, the models receive only the source-side context and generate the target-side context before the current sentence. We obtained the translation of the current sentence by splitting the output on the separator token. The training hyper-parameters and additional details can be seen

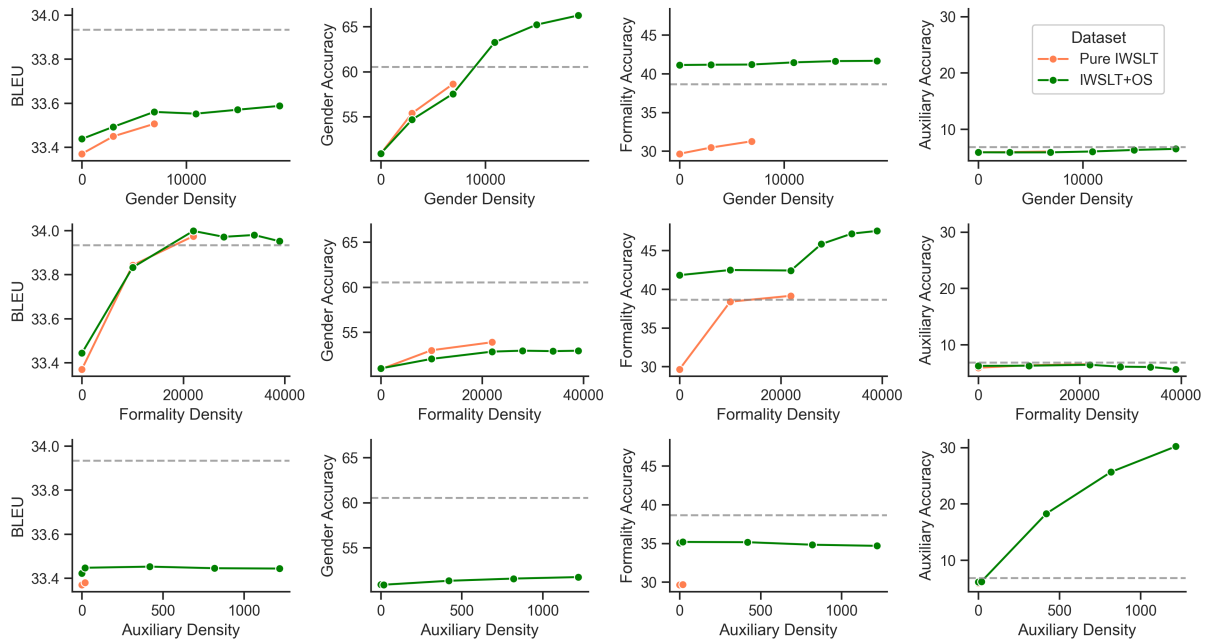


Figure 2: Measured metrics of BLEU on IWSLT 2017 testset, and ctxPro accuracy on Gender, Formality, and Auxiliary phenomena (in columns) of the OpusMT en-de models trained on the datasets with varying amounts of contextually-rich examples of Gender, Formality, and Auxiliary phenomena (in rows). Shows two experimental settings: Pure IWSLT and combined IWSLT+OS.

in Appendix B.

3.3 Single Language Pair Results

For the models in the English-to-German experiments, we trained 5 models with different seeds and averaged the results. Apart from the constructed datasets, we also trained a baseline model on the unmodified IWSLT training dataset. To measure the general translation quality, we translated the IWSLT 2017 English-to-German test subset (with BEAM search of 5) and measured BLEU (Papineni et al., 2002). Additionally, we translated test subsets of the ctxPro dataset (based on OpenSubtitles) and measured the accuracy of matching the expected word in the translation (using the scripts provided with the dataset). The results can be seen in Figure 2. Extended results including COMET and ContraPro (Müller et al., 2018) accuracy can be found in Appendix C.

We observed a drop in BLEU for the models trained on the sparse datasets, even for the datasets with mixed OpenSubtitles examples. While the reduction was relatively small (less than 2%), it returned to the baseline value only when Formality IWSLT-dense examples were added to the dataset. This could mean that the examples from the IWSLT dataset annotated with Formality were particularly influential for the model’s general translation abil-

ity, and mixing in the random examples from OpenSubtitles did not help.

For Gender and Formality, increasing their density in the training dataset improved the ctxPro accuracy for the corresponding phenomenon. Notably, Formality in the IWSLT+OS setting only improved when OS-dense examples were added, but exceeded the accuracy of the baseline model even with the most sparse dataset. Adding OS-dense examples improved the accuracy significantly above the baseline (up to 30%). Interestingly, adding dense examples in one phenomenon had minimal effect on the accuracy of the other phenomena, with only a very small increase of Formality for the Gender-enriched dataset and vice versa. Those results show that the generalizability of the models’ ability to handle contextual phenomena is very limited.

3.4 Multilingual Results

For the multilingual experiments, we trained models (with a single seed due to the computational cost of training and evaluation) on the composed datasets and measured ctxPro accuracy for all available phenomena and language directions included in the experiments. Note that Inflection applies only to English-to-Polish and English-to-Russian, and Animacy only to X-to-English. The results are

Enriched Dataset	Baseline	61.3	45.9	51.4	38.6	45.9	45.6	36.8	50.2	35.5	56.7	16.4	25.3	16.2	34.3	33.5	45.1	39.7	62.0	86.9	69.7	71.7	64.8	
En-De Gender		+12.8	+0.7	+1.8	+1.7	+3.4	+0.3	-0.1	-0.4	+0.0	+0.0	+0.2	-0.3	-0.7	-0.2	-0.3	-0.2	-0.2	-0.2	-0.3	-0.2	-0.5	-0.2	
En-Es Gender		+1.5	+23.9	+1.6	+1.5	+1.0	+0.1	+0.6	-0.4	+0.1	-0.2	+0.8	+0.3	-0.5	-0.2	-0.4	-0.1	-0.2	-0.3	-0.5	-0.4	-0.6	-0.4	
En-Fr Gender		+2.7	+1.2	+10.6	+1.2	+2.4	+0.2	-0.0	-0.2	+0.2	+0.1	+0.2	+0.2	-0.3	-0.2	-0.6	-0.0	-0.1	-0.5	-0.4	-0.2	-0.3	-0.4	
En-Pl Gender		+2.4	+0.9	+1.3	+20.1	+2.8	+0.1	-0.1	-0.3	+0.3	-0.1	+0.2	+0.2	-0.4	+0.2	-0.5	+0.0	-0.1	-0.5	-0.4	-0.4	-0.2	-0.4	
En-Ru Gender		+2.9	+0.4	+1.6	+2.3	+21.9	+0.1	-0.0	-0.3	-0.1	-0.2	+1.0	+0.3	-0.4	+0.1	-0.4	-0.1	-0.1	-0.1	-0.1	-0.3	-0.2	-0.5	+0.1
En-De Formality		-0.1	-0.2	+0.0	+0.0	-0.0	+3.2	+0.0	-0.1	+0.5	+0.1	+0.2	+0.0	-0.1	-0.2	-0.1	-0.2	-0.1	-0.6	-0.5	-0.6	-0.7	-0.5	
En-Es Formality		+0.7	+0.8	+0.1	+0.0	-0.0	+0.4	+9.1	-0.3	+0.7	+0.2	+0.3	-2.0	-0.6	-0.4	-0.9	-0.1	-0.0	-0.8	-0.7	-0.6	-0.9	-1.0	
En-Fr Formality		+0.2	-0.1	-0.1	+0.0	+0.1	+0.1	+0.2	+0.8	+0.2	-0.0	+0.4	+0.1	-1.2	+0.2	-0.5	-0.1	-0.1	-0.3	-0.6	-0.5	-0.2	-0.3	
En-Pl Formality		+0.1	-0.3	+0.1	-0.3	-0.3	+0.4	+0.4	-0.1	+8.6	+0.3	+0.4	-0.2	-0.5	+0.4	-0.5	-0.1	+0.0	-0.6	-0.6	-0.7	-0.8	-0.8	
En-Ru Formality		+0.4	-0.2	+0.1	+0.2	-0.4	+0.2	-0.0	-0.4	+0.3	+3.1	+0.8	+0.6	-0.4	-0.3	-0.9	-0.1	-0.3	-0.6	-0.6	-0.5	-0.5	-0.5	
En-De Auxiliary		+0.4	-0.2	+0.0	-0.0	+0.1	+0.3	+0.0	-0.4	+0.1	-0.0	+18.6	+5.4	+5.0	+4.6	+3.9	-0.1	-0.1	-0.3	-0.4	-0.3	-0.5	-0.3	
En-Es Auxiliary		+0.3	+0.2	+0.4	-0.1	+0.1	+0.2	-0.0	-0.5	+0.1	-0.1	+5.5	+25.4	+6.3	+5.7	+4.3	+0.1	-0.0	-0.6	-0.6	-0.6	-0.7	-0.6	
En-Fr Auxiliary		+0.5	-0.2	+0.1	+0.2	+0.3	+0.2	-0.0	-0.3	+0.1	+0.1	+5.1	+8.3	+15.8	+4.3	+4.8	-0.1	-0.1	-0.4	-0.5	-0.5	-0.4	-0.5	
En-Pl Auxiliary		+0.3	-0.1	+0.1	-0.1	-0.1	+0.2	-0.1	-0.4	+0.0	-0.2	+5.0	+7.3	+3.6	+13.6	+6.4	-0.0	+0.0	-0.2	-0.4	-0.3	-0.2	-0.4	
En-Ru Auxiliary		+0.4	-0.1	+0.2	+0.2	+0.6	+0.1	-0.0	-0.5	+0.1	-0.2	+4.0	+6.0	+3.3	+7.1	+10.2	-0.1	-0.2	-0.2	-0.6	-0.4	-0.5	-0.3	
En-Pl Inflection		+0.4	-0.2	+0.1	-0.6	-0.2	+0.1	+0.0	-0.4	-0.2	-0.1	+0.5	+0.1	-0.4	-0.1	-0.5	+9.7	+1.6	-0.5	-0.7	-0.5	-0.6	-0.6	
En-Ru Inflection		+0.4	-0.1	+0.2	+0.0	-0.2	+0.2	+0.1	-0.5	+0.0	-0.2	+0.5	+0.2	-0.4	-0.2	-1.2	+1.1	+5.6	-0.5	-0.6	-0.5	-1.0	-0.7	
De-En Animacy		+0.5	-0.1	+0.3	+0.3	+0.5	+0.3	-0.1	-0.4	+0.3	-0.0	+0.1	-0.4	-0.6	-0.3	-0.7	-0.0	-0.0	+11.6	+0.8	+1.8	+3.1	+2.6	
Es-En Animacy		+0.3	-0.3	+0.2	+0.1	+0.2	+0.2	+0.0	-0.5	+0.2	-0.1	+0.5	+0.5	-0.4	-0.2	-0.7	-0.1	-0.1	+1.1	+5.4	+1.5	+1.9	+1.1	
Fr-En Animacy		+0.2	-0.4	+0.2	+0.0	+0.4	+0.0	-0.1	-0.5	+0.1	-0.2	+0.3	-0.1	-0.5	-0.1	-0.6	-0.2	-0.1	+2.5	+2.0	+6.3	+2.0	+2.1	
Pl-En Animacy		+0.3	-0.1	+0.1	+0.1	+0.2	+0.1	+0.0	-0.5	-0.0	-0.2	+0.5	+0.1	-0.6	+0.1	-0.3	-0.2	-0.1	+1.7	+0.6	+0.9	+13.8	+2.4	
Ru-En Animacy		+0.3	-0.1	+0.2	+0.2	+0.4	+0.2	-0.0	-0.5	-0.0	-0.2	+0.4	+0.2	-0.6	-0.1	-0.3	+0.0	+0.0	+2.6	+0.6	+1.4	+3.7	+7.7	

Figure 3: Measured ctxPro accuracy on all phenomena for each of the relevant language directions (in columns) of the models trained on the OpenSubtitles datasets with varying amounts of contextually-rich examples for each phenomena and language directions (in rows). We show the differences from the baseline model (top row).

presented in Figure 3. Results in terms of BLEU and COMET on the testsets sampled from OpenSubtitles for each language direction can be seen in Appendix C.

For each model, the highest improvement in accuracy was observed for the phenomenon and language direction that was added to the training dataset (values on the diagonal in the figure). In line with the results on the single language pair, we did not observe any intra-lingual transfer between phenomena. Interestingly, there was some transfer between language directions for the same phenomenon, which was the strongest for Auxiliary, moderate for Gender, Inflection, and Animacy, and no transfer for Formality. Contrary to our expectations, we did not observe notably stronger transferability between languages in the same linguistic sub-family, with the exception of Auxiliary, where the increase in accuracy is slightly higher inside Romance and Slavic languages than for other languages. Surprisingly, the transferability between Gender and Animacy was not observed, even though the phenomena in question are a reflection of each other.

3.5 Discussion

We experimentally confirmed the dataset sparsity hypothesis by showing that the models trained on datasets sparse in contextually rich examples exhibit poor context utilization, and increasing the density leads to large improvements for the tested phenomena. Both experiments showed that the models do not generalize context utilization between phenomena. This finding calls for caution when interpreting the results of evaluations targeting a single phenomenon (Müller et al., 2018; Lopes et al., 2020). While document-level training datasets typically include a representative (for a particular domain) mixture of contextual phenomena, we found that models can develop strong capabilities for some phenomena, while remaining weak on others. Mařka et al. (2025) found attention heads in context-aware MT models responsible for pronoun disambiguation with some cross-lingual behavior, which is in line with the observed transferability between language directions. We hypothesize that the poor transfer between phenomena can be explained by the models developing separate heads for each of them.

4 Methods Exploiting Contextual Data

Inspired by the fact that increased density in contextually-relevant examples of the training dataset leads to improvement in context utilization, we tested several techniques that could leverage the available data more efficiently. We broadly divide them into annotation-based and annotation-free. Annotations can inform the training process but require an external tool (e.g., ctxPro) to mark the relevant examples. A straightforward method is to simply extract the annotated examples from the training dataset and use them to fine-tune the model. Annotation-free methods do not rely on an external tool and have the advantage of generalizability beyond the phenomena covered by any tool. Crucially, the presented methods aim to improve contextual capabilities without the need for any additional data beyond the standard training datasets.

4.1 Token-level Loss Weighting

We adapted the weighting of the loss elements (Wang et al., 2017), which increases the error signal coming from selected examples. Instead of weighting the whole examples, we apply a token-level approach as phenomena annotations contain an expected word or phrase that requires context for successful translation. We train the models using the weighted negative log-likelihood loss function:

$$\mathcal{L} = -\frac{1}{|D_a|} \sum_{(x_i, y_i, a_i) \sim D_a} \sum_{j=1}^{|y_i|} w(a_{i,j}) \log(\hat{y}_{i,j}), \quad (1)$$

where $\hat{y}_{i,j}$ is the probability of the j -th token in i -th example, D_a is the annotated training dataset with examples containing input and output sequences (x_i and y_i respectively), as well as the token-level annotations a_i marking the contextually-dependent tokens, and $w_{i,j}$ is defined as:

$$w_{i,j} = \begin{cases} 1 + \lambda, & \text{if contextually dependent,} \\ 1, & \text{otherwise,} \end{cases} \quad (2)$$

for each token j in the i -th output sequence, where λ is the hyper-parameter.

4.2 Metric-based Example Selection

A major issue with using annotations is that, according to our experiments on data composition, the model will improve only on the included phenomena. To mitigate this, we propose to utilize the

model itself to mark contextually-rich examples. Fernandes et al. (2023) proposed the Point-wise Cross-Mutual Information (PCXMI) metric to measure the context reliance of the translations, which is based on the output probabilities of the context-aware MT model. For a particular example it is calculated as:

$$PCXMI = \sum_{j=1}^{|y|} \log \frac{q(y_j | y_{t < j}, x, C)}{q(y_j | y_{t < j}, x)}, \quad (3)$$

where C is the context, and q represents the context-aware MT model (returning token probabilities, noted as $q(y_j | y_{t < j}, x, C)$) that is trained to also be used as a sentence-level model (noted as $q(y_j | y_{t < j}, x)$). We introduce a slightly modified metric that computes the maximum token-level PCXMI for a given example:

$$MaxPCXMI = \max_j \left(\log \frac{q(y_j | y_{t < j}, x, C)}{q(y_j | y_{t < j}, x)} \right). \quad (4)$$

We motivate it by the fact that an example with even a single token being dependent on context can be considered a contextually-rich example (certainly the case for pronouns), which is better captured by our metric. The proposed method consists of the following steps:

1. **train** the model on context-aware data,
2. **calculate** the metric using the trained model for the examples in the training dataset,
3. **select** top k examples (a hyper-parameter),
4. **fine-tune** the model on the selected subset.

While the method can be seen as similar to curriculum learning (Zhang et al., 2018), we select the examples that the model is already competent at translating using context. Intuitively, this is a positive feedback where the model learns to generalize to the difficult examples by becoming better at what it already knows.

5 Experiments

We experimentally evaluated Token-level Loss Weighting and Metric-based Example Selection for fine-tuning and compared them to the following baselines:

- **Fine-tuning** (annotation-based) - simply fine-tuning the model on the annotated data after the context-aware training.
- **CoWord Dropout** (annotation-free; Fernandes et al., 2021) - masking random tokens in the current

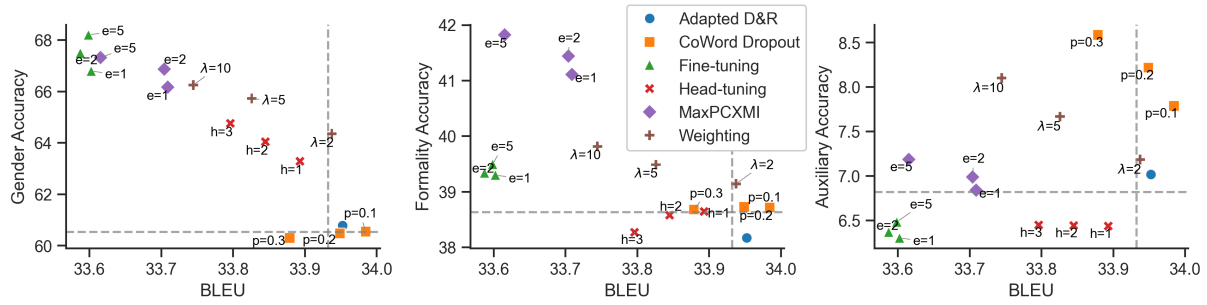


Figure 4: Accuracy of ctxPro English-to-German phenomena (Gender, Formality, and Auxiliary) against BLEU on the IWSLT 2017 en-de testset of the trained models. Labels show: the number of epochs ("e"), CoWord Dropout probability ("p"), number of tuned heads("h"), and weighting strength ("λ") hyper-parameters.

source sentence to force the model to use context for translation, the probability of masking a token is controlled by the hyper-parameter p .

- **Adapted Divide and Rule** (annotation-free; [Lupo et al., 2022](#)) - splitting the current source and target sentences in the middle and appending the first parts to the context. Notably, this method was introduced for the multi-encoder architecture and trained only the contextual parameters. We adapt it to the single-encoder architectures we use in this study and do not freeze the parameters of the model.

- **Head-tuning** (annotation-based; [Mařka et al., 2025](#)) - training selected attention heads to attend the context cue, therefore, available only for Gender.

We evaluated all methods in the single language pair (English-to-German) setting and annotation-free methods in the multilingual setting (due to the lack of exhaustive annotations for the dataset). We used the same base sentence-level models: OpusMT en-de and NLLB-200 600M, respectively. For English-to-German, we trained on the full IWSLT 2017 en-de dataset with ctxPro annotations, and for multilingual, we sampled 50,000 examples for each language direction from the OS-rand dataset. We used the same hyper-parameters shared by all methods as in previous experiments (see Appendix B for more details) for both training and fine-tuning with the exception of Head-tuning where we applied the hyper-parameters from the original paper. In the English-to-German setting, we repeated the training 5 times with different seeds and averaged the results. In the multi-lingual setting, we performed a single training run for all models with the same seed. Fine-tuning used the base model trained with the corresponding seed.

5.1 Single Language Pair Results

We tested several parameters for most methods. For fine-tuning-based models, we trained for $e \in \{1, 2, 5\}$ epochs and utilized only the examples with the maximum context size. For Weighting we set the λ parameter to 2, 5, and 10. In addition to the values of p for CoWord Dropout recommended by the authors (0.1, 0.2), we also included the value of 0.3. For Metric-based example selection, we set $k=30,000$ based on the number of annotated examples in the dataset, and used the MaxPCXMI metric (in Appendix D we present the comparison to the PCXMI metric). For Head-tuning we selected top $h \in \{1, 2, 3\}$ heads from [Mařka et al. \(2025\)](#). Results in terms of accuracy on the ctxPro dataset and BLEU on the IWSLT testset can be seen in Figure 4. Extended results are presented in Appendix D.

It can be seen that with four metrics, the models' performance varies, and improvement in one metric comes at a cost of a reduction in another. In particular, we observe a negative relation between ctxPro accuracies and BLEU for all methods with the increase of the hyper-parameters. Metric-based example selection achieved highest improvement in Formality and outperformed the annotation-based selection for fine-tuning in Formality and Auxiliary, and achieved similar results for Gender, with a smaller decrease in BLEU. Head-tuning showed improvement only on Gender but with smaller drop in BLEU. Methods applied during training (Weighting, CoWord Dropout, and Divide and Rule) showed a smaller reduction in BLEU compared to fine-tuning. We attribute this to the smaller discrepancy in the dataset distribution between training and evaluation. Weighting outperformed CoWord Dropout on Gender and Auxiliary. Conversely, CoWord Dropout achieved the highest accuracy on Auxiliary (with Weighting being the

Model	BLEU	Gender	Formality	Auxiliary	Inflection	Animacy
Adapted D&R	-0.05	-0.06	-0.19	-0.16	-0.03	+0.07
CoWord p=0.1	-0.09	+0.02	-0.16	+0.35	-0.10	+0.16
CoWord p=0.2	-0.11	+0.07	-0.28	+0.65	-0.21	+0.01
CoWord p=0.3	-0.08	+0.01	-0.42	+0.97	-0.29	-0.27
MaxPCXMI e=1	-0.42	+1.13	+0.05	+3.41	+0.44	+1.08
MaxPCXMI e=2	-0.45	+1.42	+0.05	+4.25	+0.57	+1.10
MaxPCXMI e=5	-0.50	+1.93	+0.11	+5.80	+0.76	+1.64

Table 1: The averaged (over language directions) difference from the baseline in terms of BLEU on OpenSubtitles 2018 testsets and ctxPro phenomena accuracies for the tested methods. Number of epochs is noted as "e", and CoWord Dropout probability as "p".

second-best) but did not show any improvement for Gender and Formality. Notably, the highest reduction in BLEU was around 1% compared to the baseline. Lack of improvement exhibited by Adapted Divide and Rule can be attributed to our adaptation implementation, which did not utilize parameter freezing as in the original paper. Among all methods, metric-based example selection achieved the highest average ctxPro accuracy across phenomena, while token-level loss weighting was the most effective among annotation-based approaches, demonstrating that both proposed techniques can substantially improve context utilization.

5.2 Multilingual Results

We trained models based on NLLB-200 600M on all relevant language-directions using annotation-free methods to assess their performance in the multilingual setting. For CoWord Dropout, we used the same values of p (0.1, 0.2, and 0.3), and for Metric-based example selection, we set $k=10,000$ per language direction and the number of epochs equal to 1, 2, and 5. The results aggregated over language directions can be seen in Table 1 and extended results in Appendix D.

Fine-tuning on examples selected by MaxPCXMI outperformed all baselines in terms of ctxPro accuracy across phenomena, with the highest improvement of 5.8, 1.9, and 1.6 percentage points (on average) for Auxiliary, Gender, and Animacy, respectively. Contrary to the English-to-German experiments, no improvement (on average) was observed for Formality. This was caused by a drop of up to 1 percentage point in the English-to-French direction, which offset small gains in other language directions. These accuracy improvements came at the cost of a greater reduction in BLEU compared to other methods, and both trends—accuracy gains and BLEU

drops—intensified with more fine-tuning epochs, mirroring the patterns seen in the single-language-direction experiments.

6 Conclusions

This work provided a systematic empirical evaluation of the influence of training data composition, in terms of contextually rich examples, on the context utilization capabilities for encoder-decoder MT models. By systematically adapting the proportion of contextually rich examples in the training data, we demonstrated that such data sparsity is the key bottleneck in learning to leverage context efficiently in MT models. Crucially, we found that (1) models do not generalize well across different contextual phenomena (e.g. gender or formality) and (2) while there is some cross-lingual transfer, it was not significantly higher between languages in the same linguistic sub-family.

Motivated by these findings, we proposed two methods designed to mitigate the effect of data sparsity in context-aware MT: token-level loss weighting (based on token-level annotations of context-dependent words) and metric-based instance selection (fine-tuning on most contextually important examples). Both methods significantly improved context utilization without the need for extensive architectural changes or additional annotated data. Notably, the metric based method showed strong gains across multiple phenomena and language directions.

In practical terms, data composition and targeted training should be considered as potential solutions to developing strong context-aware MT models. In future work, we will evaluate our findings for decoder-only models and combine the strengths of weighting and metric-based example selection.

7 Limitations

While we investigate many language directions and three sub-families, all of them come from the Indo-European family. This limitation was imposed by the language directions covered by ctxPro toolset. Additionally, for the single language pair setting, we only tested English-to-German direction. We suspect that the uncovered effects of data composition go beyond the tested language pairs, but this claim has not been tested experimentally.

Similarly, we only tested encoder-decoder architectures with the single encoder (standard Transformer). Both multi-encoder and decoder-only models lay beyond the scope of this study. Furthermore, we leveraged only the publicly available pre-trained sentence-level models as the basis for context-aware training. We argue that this increases reproducibility as the models are freely available for further investigation by other researchers. Nevertheless, we did not experiment in this study on the models trained from randomly initialized weights.

References

Ruchit Agrawal, Marco Turchi, and Matteo Negri. 2018. Contextual handling in neural machine translation: Look behind, ahead and on both sides. In *Proceedings of the 21st Annual Conference of the European Association for Machine Translation*, pages 31–40.

Ibrahim M Alabdulmohsin, Xiaohua Zhai, Alexander Kolesnikov, and Lucas Beyer. 2023. Getting vit in shape: Scaling laws for compute-optimal model design. *Advances in Neural Information Processing Systems*, 36:16406–16425.

Belen Alastruey, Gerard I. Gállego, and Marta R. Costajussà. 2024. [Unveiling the role of pretraining in direct speech translation](#). In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 11259–11265, Miami, Florida, USA. Association for Computational Linguistics.

Duarte M Alves, José Pombal, Nuno M Guerreiro, Pedro H Martins, João Alves, Amin Farajian, Ben Peters, Ricardo Rei, Patrick Fernandes, Sweta Agrawal, and 1 others. 2024. Tower: An open multilingual large language model for translation-related tasks. *arXiv preprint arXiv:2402.17733*.

Abdul Hameed Azeemi, Ihsan Ayyub Qazi, and Agha Ali Raza. 2025. [To label or not to label: Hybrid active learning for neural machine translation](#). In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 3071–3082, Abu Dhabi, UAE. Association for Computational Linguistics.

Guangsheng Bao, Yue Zhang, Zhiyang Teng, Boxing Chen, and Weihua Luo. 2021. [G-transformer for document-level machine translation](#). 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 3442–3455, Online. Association for Computational Linguistics.

Rachel Bawden, Rico Sennrich, Alexandra Birch, and Barry Haddow. 2018. [Evaluating discourse phenomena in neural machine translation](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1304–1313, New Orleans, Louisiana. Association for Computational Linguistics.

Lola Le Breton, Quentin Fournier, Mariam El Mezouar, and Sarath Chandar. 2025. Neobert: A next-generation bert. *arXiv preprint arXiv:2502.19587*.

Aydar Bulatov, Yury Kuratov, and Mikhail Burtsev. 2022. Recurrent memory transformer. *Advances in Neural Information Processing Systems*, 35:11079–11091.

Mauro Cettolo, Marcello Federico, Luisa Bentivogli, Jan Niehues, Sebastian Stüker, Katsuhito Sudoh, Koichiro Yoshino, and Christian Federmann. 2017. [Overview of the IWSLT 2017 evaluation campaign](#). In *Proceedings of the 14th International Conference on Spoken Language Translation*, pages 2–14, Tokyo, Japan. International Workshop on Spoken Language Translation.

Linqing Chen, Junhui Li, Zhengxian Gong, Min Zhang, and Guodong Zhou. 2022. [One type context is not enough: Global context-aware neural machine translation](#). *ACM Trans. Asian Low-Resour. Lang. Inf. Process.*, 21(6).

Chenhui Chu, Raj Dabre, and Sadao Kurohashi. 2017. [An empirical comparison of domain adaptation methods for neural machine translation](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 385–391, Vancouver, Canada. Association for Computational Linguistics.

Benjamin Clavié, Nathan Cooper, and Benjamin Warner. 2025. It’s all in the [mask]: Simple instruction-tuning enables bert-like masked language models as generative classifiers. *arXiv preprint arXiv:2502.03793*.

Yukun Feng, Feng Li, Ziang Song, Boyuan Zheng, and Philipp Koehn. 2022. [Learn to remember: Transformer with recurrent memory for document-level machine translation](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 1409–1420, Seattle, United States. Association for Computational Linguistics.

Patrick Fernandes, Kayo Yin, Emmy Liu, André Martins, and Graham Neubig. 2023. [When does translation require context? a data-driven, multilingual](#)

823	heads for pronoun disambiguation in context-aware machine translation models . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 6348–6377, Abu Dhabi, UAE. Association for Computational Linguistics.	<i>Machine Translation: Research Papers</i> , pages 186–191, Brussels, Belgium. Association for Computational Linguistics.	879 880 881
824			
825			
826			
827			
828	Ali Marashian, Enora Rice, Luke Gessler, Alexis Palmer, and Katharina von der Wense. 2025. From priest to doctor: Domain adaptation for low-resource neural machine translation . In <i>Proceedings of the 31st International Conference on Computational Linguistics</i> , pages 7087–7098, Abu Dhabi, UAE. Association for Computational Linguistics.	Matt Post and Marcin Junczys-Dowmunt. 2023. Escaping the sentence-level paradigm in machine translation . <i>arXiv preprint arXiv:2304.12959</i> .	882 883 884
829			
830			
831			
832			
833			
834			
835	Sameen Maruf, André F. T. Martins, and Gholamreza Haffari. 2019. Selective attention for context-aware neural machine translation . In <i>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)</i> , pages 3092–3102, Minneapolis, Minnesota. Association for Computational Linguistics.	Matt Post and Marcin Junczys-Dowmunt. 2024. Evaluation and large-scale training for contextual machine translation . In <i>Proceedings of the Ninth Conference on Machine Translation</i> , pages 1125–1139, Miami, Florida, USA. Association for Computational Linguistics.	885 886 887 888 889 890
836			
837			
838			
839			
840			
841			
842			
843	Lesly Miculicich, Dhananjay Ram, Nikolaos Pappas, and James Henderson. 2018. Document-level neural machine translation with hierarchical attention networks . In <i>Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing</i> , pages 2947–2954, Brussels, Belgium. Association for Computational Linguistics.	Ricardo Rei, Craig Stewart, Ana C Farinha, and Alon Lavie. 2020. COMET: A neural framework for MT evaluation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 2685–2702, Online. Association for Computational Linguistics.	891 892 893 894 895 896
844			
845			
846			
847			
848			
849			
850	Wafaa Mohammed and Vlad Niculae. 2024. On measuring context utilization in document-level MT systems . In <i>Findings of the Association for Computational Linguistics: EACL 2024</i> , pages 1633–1643, St. Julian’s, Malta. Association for Computational Linguistics.	Noam M. Shazeer and Mitchell Stern. 2018. Adafactor: Adaptive learning rates with sublinear memory cost . <i>ArXiv</i> , abs/1804.04235.	897 898 899
851			
852			
853			
854			
855	Mathias Müller, Annette Rios, Elena Voita, and Rico Sennrich. 2018. A large-scale test set for the evaluation of context-aware pronoun translation in neural machine translation . In <i>Proceedings of the Third Conference on Machine Translation: Research Papers</i> , pages 61–72, Brussels, Belgium. Association for Computational Linguistics.	Zewei Sun, Mingxuan Wang, Hao Zhou, Chengqi Zhao, Shujian Huang, Jiajun Chen, and Lei Li. 2022. Re-thinking document-level neural machine translation . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 3537–3548, Dublin, Ireland. Association for Computational Linguistics.	900 901 902 903 904 905
856			
857			
858			
859			
860			
861			
862	NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia-Gonzalez, Prangthip Hansanti, and 20 others. 2022. No language left behind: Scaling human-centered machine translation .	Jörg Tiedemann and Yves Scherrer. 2017. Neural machine translation with extended context . In <i>Proceedings of the Third Workshop on Discourse in Machine Translation</i> , pages 82–92, Copenhagen, Denmark. Association for Computational Linguistics.	906 907 908 909 910
863			
864			
865			
866			
867			
868			
869			
870	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.	Zhaopeng Tu, Yang Liu, Zhengdong Lu, Xiaohua Liu, and Hang Li. 2017. Context gates for neural machine translation . <i>Transactions of the Association for Computational Linguistics</i> , 5:87–99.	911 912 913 914
871			
872			
873			
874			
875			
876			
877	Matt Post. 2018. A call for clarity in reporting BLEU scores . In <i>Proceedings of the Third Conference on</i>	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need . <i>Advances in neural information processing systems</i> , 30.	915 916 917 918 919
878			
		Elena Voita, Rico Sennrich, and Ivan Titov. 2019a. Context-aware monolingual repair for neural machine translation . 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 877–886, Hong Kong, China. Association for Computational Linguistics.	920 921 922 923 924 925 926 927
		Elena Voita, Rico Sennrich, and Ivan Titov. 2019b. When a good translation is wrong in context: Context-aware machine translation improves on deixis, ellipsis, and lexical cohesion . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 1198–1212, Florence, Italy. Association for Computational Linguistics.	928 929 930 931 932 933 934

935	Rui Wang, Masao Utiyama, Lemao Liu, Kehai Chen, and Eiichiro Sumita. 2017. Instance weighting for neural machine translation domain adaptation . In <i>Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing</i> , pages 1482–1488, Copenhagen, Denmark. Association for Computational Linguistics.	994
936		995
937		996
938		997
939		998
940		999
941		
942	Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, and 1 others. 2024. Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference. <i>arXiv preprint arXiv:2412.13663</i> .	1000
943		1001
944		1002
945		1003
946		1004
947		1005
948		
949	Rachel Wicks and Matt Post. 2023. Identifying context-dependent translations for evaluation set production . In <i>Proceedings of the Eighth Conference on Machine Translation</i> , pages 452–467, Singapore. Association for Computational Linguistics.	1006
950		1007
951		1008
952		1009
953		1010
954	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, and 3 others. 2020. Transformers: State-of-the-art natural language processing . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pages 38–45, Online. Association for Computational Linguistics.	1011
955		1012
956		1013
957		1014
958		1015
959		1016
960		1017
961		1018
962		1019
963		1020
964		1021
965	Billy T. M. Wong and Chunyu Kit. 2012. Extending machine translation evaluation metrics with lexical cohesion to document level . In <i>Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning</i> , pages 1060–1068, Jeju Island, Korea. Association for Computational Linguistics.	1022
966		1023
967		1024
968		1025
969		1026
970		1027
971		1028
972	Kayo Yin, Patrick Fernandes, Danish Pruthi, Aditi Chaudhary, André F. T. Martins, and Graham Neubig. 2021. Do context-aware translation models pay the right attention? 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 1: Long Papers)</i> , pages 788–801, Online. Association for Computational Linguistics.	1029
973		1030
974		1031
975		1032
976		1033
977		1034
978		1035
979		1036
980		1037
981	Pei Zhang, Boxing Chen, Niyu Ge, and Kai Fan. 2020. Long-short term masking transformer: A simple but effective baseline for document-level neural machine translation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 1081–1087, Online. Association for Computational Linguistics.	1038
982		1039
983		1040
984		1041
985		1042
986		1043
987		1044
988	Xuan Zhang, Gaurav Kumar, Huda Khayrallah, Kenton Murray, Jeremy Gwinnup, Marianna J Martindale, Paul McNamee, Kevin Duh, and Marine Carpuat. 2018. An empirical exploration of curriculum learning for neural machine translation. <i>arXiv preprint arXiv:1811.00739</i> .	1045
989		
990		
991		
992		
993		
	Zaixiang Zheng, Xiang Yue, Shujian Huang, Jiajun Chen, and Alexandra Birch. 2021. Towards making the most of context in neural machine translation. In <i>Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence</i> , pages 3983–3989.	
	Zhang Zhuocheng, Shuhao Gu, Min Zhang, and Yang Feng. 2023. Scaling law for document neural machine translation . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 8290–8303, Singapore. Association for Computational Linguistics.	
	A Composition of the Datasets	
	In this section, we describe how the constructed datasets were created. For the <i>Pure IWSLT</i> setting, we start with the IWSLT-sparse (123,000 examples with no annotations) and progressively replace it with the examples sampled from IWSLT-dense. The steps are based on the size of the IWSLT-dense dataset for a particular phenomenon: 3,000 and 6,915 (full size) for Gender, 10,000 and 21,977 (full size) for Formality, and 19 (full size) for Auxiliary. For the <i>IWSLT + OS</i> setting, we start with the datasets formed by combining IWSLT-sparse with examples sampled from OS-rand. To maximize the density of the resulting datasets, we set the number of examples sampled from OS-rand to be dependent on the phenomenon and equal to the (rounded) size of the OS-dense datasets: 12,000 for Gender, 17,000 for Formality, and 1,200 for Auxiliary. We start by replacing examples from IWSLT-sparse (we retain the steps from the Pure IWSLT setting). After reaching the maximum density in the IWSLT portion of the dataset, we start replacing OS-rand with OS-dense in the following steps: 4,000, 8,000, and 12,000 for Gender, 6,000, 12,000, and 17,000 for Formality, and 400, 800, and 1,200 for Auxiliary.	
	Tables 2, 3, and 4 show the composition of the training datasets we used in the experiments for Gender, Formality, and Auxiliary phenomena, respectively. Each example was encoded with the context size ranging from zero to the maximum context size (three in our experiments), increasing the size of the datasets four times.	
	In the multilingual experiments, we formed the baseline training dataset by sampling 50,000 examples from OpenSubtitles (OS-rand) for each language direction we considered. For each phenomenon in a language direction, we replaced examples with the rich ones: 6,900 for Gender, 10,000 for Formality, 1,200 for Auxiliary, 10,000	

1046 for Inflection, and 4,000 for Animacy.
1047 ([NLLB Team et al., 2022](#)).

1048 B Details of Context-aware Training

1049 We implemented all experiments in *Hugging-*
1050 *face transformers* framework ([Wolf et al., 2020](#)).
1051 We trained the models with Adafactor optimizer
1052 ([Shazeer and Stern, 2018](#)) on a single GPU
1053 (NVIDIA GeForce RTX 3090 24GB and NVIDIA
1054 H100 80GB for the models based on OpusMT
1055 en-de and NLLB-200 600M, respectively) for 10
1056 epochs. OpusMT en-de² and NLLB-200 600M³
1057 contain 163M and 615M parameters, respectively.
1058 The hyper-parameters are presented in Table 5. We
1059 tuned the hyper-parameters (learning rate, batch
1060 size, number of epochs) during the preliminary ex-
1061 periments on OpusMT en-de model with context
1062 size of one trained on IWSLT 2017 English-to-
1063 German dataset.

1064 C Extended Data Composition Results

1065 In this section, we present the extended results
1066 of the data composition experiments. For single
1067 language pair setting, we measured COMET ([Rei](#)
1068 [et al., 2020](#)) (based on `Unbabel/wmt22-comet-da`)
1069 on the IWSLT 2017 en-de testset and evaluated
1070 the models on the ContraPro ([Müller et al., 2018](#))
1071 contrastive evaluation. The results for the Pure
1072 IWSLT and IWSLT+OS settings can be found in
1073 Tables 6 and 7, respectively.

1074 For the multilingual setting, we measured BLEU
1075 (we used the `sacreBLEU` library ([Post, 2018](#)) using
1076 the default parameters) and COMET on the testsets
1077 formed by sampling 20,000 examples from Open-
1078 Subtitles 2018 for each language direction. The
1079 results can be seen in Tables 8 and 9 for BLEU and
1080 COMET, respectively.

1081 D Extended Fine-tuning Results

1082 For the English-to-German experiment, apart
1083 from BLEU and `ctxPro` accuracy, we also mea-
1084 sured COMET ([Rei et al., 2020](#)) (based on
1085 `Unbabel/wmt22-comet-da`) on the IWSLT 2017
1086 en-de testset and the accuracy on the ContraPro
1087 contrastive evaluation. The results (including
1088 BLEU and `ctxPro` accuracies) can be seen in Ta-
1089 ble 10.

²<https://huggingface.co/Helsinki-NLP/opus-mt-en-de>

³<https://huggingface.co/facebook/nllb-200-distilled-600M>

1090 Next, we present the results of Metric-based se-
1091 lection of examples for fine-tuning for two metrics:
1092 PCXMI ([Fernandes et al., 2023](#)) and MaxPCXMI
1093 (ours). We fine-tuned the models for 1, 2, and 5
1094 epochs and repeated the experiment 5 times with
1095 different seeds (using the base context-aware model
1096 trained with the corresponding seed). The averaged
1097 results can be seen in Figure 5. Selecting exam-
1098 ples based on MaxPCXMI outperforms PCXMI
1099 in Gender and Formality at a lower reduction in
1100 BLEU. PCXMI achieves a better increase in Auxili-
1101 ary but reduces BLEU even below the level of the
1102 annotation-based method.

1103 The un-aggregated results of the trained mod-
1104 els for each language direction in the multilingual
1105 experiment can be seen in Figure 6 and Tables 11
1106 and 12 for `ctxPro` accuracies, BLEU and COMET,
1107 respectively.

Setting	IWSLT-sparse	IWSLT-dense	OS-rand	OS-dense	Total
Pure IWSLT	123,000	0	0	0	123,000
	120,000	3,000	0	0	123,000
	116,085	6,915	0	0	123,000
IWSLT+OS	123,000	0	12,000	0	135,000
	120,000	3,000	12,000	0	135,000
	116,085	6,915	12,000	0	135,000
	116,085	6,915	8,000	4,000	135,000
	116,085	6,915	4,000	8,000	135,000
	116,085	6,915	0	12,000	135,000

Table 2: Number of examples from datasets that were used to compose training datasets (in rows) for the **Gender** phenomenon in the single language direction (English-to-German) setting.

Setting	IWSLT-sparse	IWSLT-dense	OS-rand	OS-dense	Total
Pure IWSLT	123,000	0	0	0	123,000
	113,000	10,000	0	0	123,000
	101,023	21,977	0	0	123,000
IWSLT+OS	123,000	0	17,000	0	140,000
	113,000	10,000	17,000	0	140,000
	101,023	21,977	17,000	0	140,000
	101,023	21,977	11,000	6,000	140,000
	101,023	21,977	5,000	12,000	140,000
	101,023	21,977	0	17,000	140,000

Table 3: Number of examples from datasets that were used to compose training datasets (in rows) for the **Formality** phenomenon in the single language direction (English-to-German) setting.

Setting	IWSLT-sparse	IWSLT-dense	OS-rand	OS-dense	Total
Pure IWSLT	123,000	0	0	0	123,000
	122,981	19	0	0	123,000
IWSLT+OS	123,000	0	1,200	0	124,200
	122,981	19	1,200	0	124,200
	122,981	19	800	400	124,200
	122,981	19	400	800	124,200
	122,981	19	0	1,200	124,200

Table 4: Number of examples from datasets that were used to compose training datasets (in rows) for the **Auxiliary** phenomenon in the single language direction (English-to-German) setting.

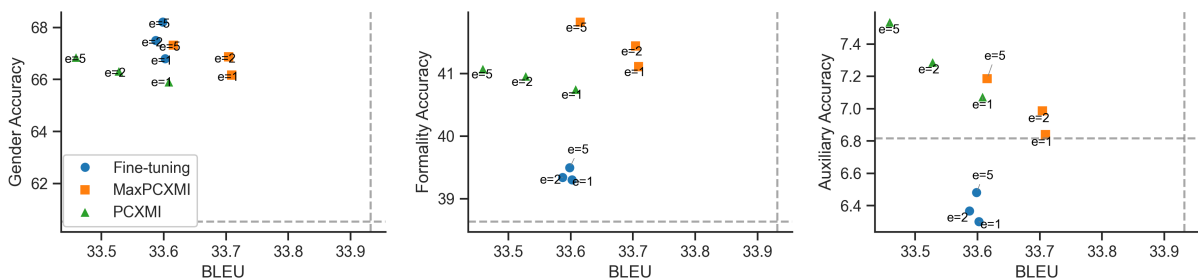


Figure 5: Accuracy of ctxPro English-to-German phenomena (Gender, Formality, and Auxiliary) against BLEU on the IWSLT 2017 en-de testset of the fine-tuned models with Metric-based (PCXMI and MaxPCXMI) and annotation-based (for comparison) selection of examples. Labels show the number of epochs ("e").

Baseline	61.3	45.9	51.4	38.6	45.9	45.6	36.8	50.2	35.5	56.7	16.4	25.3	16.2	34.3	33.5	45.1	39.7	62.0	86.9	69.7	71.7	64.8
Adapted D&R	+0.3	-0.5	-0.0	+0.0	-0.1	-0.0	-0.1	-0.5	-0.2	-0.2	+0.5	-0.4	-0.4	+0.1	-0.6	+0.0	-0.1	+0.3	-0.1	-0.0	-0.0	+0.2
CoWord p=0.1	+0.4	-0.0	+0.1	-0.1	-0.3	+0.0	-0.3	-0.3	-0.0	-0.2	+0.5	+0.8	+0.7	-0.3	+0.1	-0.1	-0.1	+0.2	+0.2	+0.1	+0.1	+0.2
CoWord p=0.2	+0.5	-0.1	+0.4	-0.2	-0.1	-0.0	-0.4	-0.4	-0.2	-0.3	+0.9	+1.4	+1.0	-0.3	+0.2	-0.3	-0.2	-0.2	+0.3	-0.3	+0.2	+0.1
CoWord p=0.3	+0.4	-0.1	+0.4	-0.3	-0.4	-0.2	-0.6	-0.4	-0.4	-0.4	+1.3	+1.9	+1.2	-0.1	+0.5	-0.3	-0.3	-0.7	+0.1	-0.5	-0.2	-0.2
maxPCXMI e=1	+1.7	+0.6	+1.2	+1.0	+1.3	+0.5	-0.1	-0.8	+0.4	+0.2	+3.3	+5.0	+2.8	+2.9	+3.0	+0.1	+0.8	+1.8	+0.3	+0.8	+1.2	+1.1
maxPCXMI e=2	+2.1	+0.8	+1.5	+1.2	+1.5	+0.5	+0.0	-1.1	+0.6	+0.2	+4.3	+6.3	+3.3	+3.8	+3.5	+0.1	+1.0	+2.1	+0.1	+0.8	+1.2	+1.2
maxPCXMI e=5	+2.8	+1.0	+1.9	+1.7	+2.2	+0.7	+0.1	-1.1	+0.7	+0.2	+5.1	+8.8	+4.9	+5.5	+4.7	+0.3	+1.2	+3.0	+0.3	+1.2	+1.7	+2.0
	En-De Gender	En-Es Gender	En-Fr Gender	En-Pl Gender	En-Ru Gender	En-De Formality	En-Es Formality	En-Fr Formality	En-Pl Formality	En-Ru Formality	En-De Auxiliary	En-Es Auxiliary	En-Fr Auxiliary	En-Pl Auxiliary	En-Ru Auxiliary	En-Pl Inflection	En-Ru Inflection	De-En Animacy	Es-En Animacy	Fr-En Animacy	Pl-En Animacy	Ru-En Animacy
	ctxPro Accuracy Difference																					

Figure 6: Measured ctxPro accuracy on all phenomena for each of the relevant language directions (in columns) of tested methods (in rows).

Hyper-parameter	Value
Optimizer	Adafactor
Learning Rate	1e-5
LR Scheduler	Inverse Sqrt
LR Warmup Ratio	0.1
Weight Decay	0.01
Batch Size	32 ^a
Gradient Accumulation Steps	16 ^a
Num Epoch	10
Precision	fp16
Seeds	1,2,3,4,5 ^b
Max Length	512 / 1024 ^c
Max Context Size	3
Beam size	5

Table 5: The hyper-parameters of the context-aware training and fine-tuning.

^a For the cases where the CUDA Out Of Memory error occurred, we reduced the batch size to 16 and increased the Gradient Accumulation Steps to 32, keeping the same effective size of the batch.

^b For the multilingual setting, we used only one seed of 1.

^c For models based on OpusMT en-de and NLLB-200 600M respectively.

Dataset	Count	COMET	ContraPro
Sparse	0	0.8415	69.23
Gender	3,000	0.8417	74.70
	6,915	0.8417	78.45
Formality	10,000	0.8429	69.55
	21,977	0.8430	70.02
Auxiliary	19	0.8413	69.14

Table 6: Performance in terms of COMET on IWSLT 2017 en-de testset and ContraPro accuracy for the models in the **Pure IWSLT** setting trained on datasets with different numbers of examples annotated with different phenomena.

Dataset	Count	COMET	ContraPro
Gender	0	0.8417	70.28
	3,000	0.8417	75.03
	6,915	0.8420	78.52
	10,915	0.8419	83.58
	14,915	0.8418	84.77
	18,915	0.8420	85.24
Formality	0	0.8416	70.15
	10,000	0.8426	70.59
	21,977	0.8428	71.12
	27,977	0.8429	71.04
	33,977	0.8429	70.85
	38,977	0.8430	71.03
Auxiliary	0	0.8414	69.47
	19	0.8415	69.39
	419	0.8415	69.60
	819	0.8415	69.75
	1,219	0.8416	69.79

Table 7: Performance in terms of COMET on IWSLT 2017 en-de testset and ContraPro accuracy for the models in the **IWSLT+OS** setting trained on datasets with different numbers of examples annotated with different phenomena.

Model	En-De	En-Es	En-Fr	En-Pl	En-Ru	De-En	Es-En	Fr-En	Pl-En	Ru-En
Baseline	26.50	37.68	29.39	21.98	24.49	32.04	41.99	32.84	29.42	31.35
Gender										
En-De	26.66	37.61	29.27	21.85	24.46	31.98	42.03	32.87	29.54	31.32
En-Es	26.88	37.60	29.33	22.12	24.52	32.03	41.96	32.86	29.46	31.36
En-Fr	26.75	37.53	29.16	22.05	24.41	32.01	41.97	32.87	29.50	31.33
En-Pl	26.80	37.57	29.21	21.54	24.48	32.05	42.00	32.86	29.53	31.34
En-Ru	26.78	37.60	29.56	21.91	24.45	32.01	42.04	32.81	29.52	31.41
Formality										
En-De	26.61	37.27	29.29	21.75	24.44	31.98	42.05	32.85	29.52	31.31
En-Es	26.58	37.29	29.43	21.65	24.57	32.01	42.04	32.84	29.49	31.39
En-Fr	26.70	37.63	29.67	21.89	24.48	32.02	41.99	32.92	29.52	31.37
En-Pl	26.62	37.38	29.44	21.83	24.35	32.03	42.00	32.88	29.44	31.23
En-Ru	26.88	37.53	29.36	22.05	24.22	32.04	42.03	32.91	29.50	31.39
Auxiliary										
En-De	26.86	37.57	29.26	21.77	24.48	32.01	42.08	32.91	29.51	31.42
En-Es	26.88	37.44	29.38	22.01	24.46	32.09	41.98	32.85	29.46	31.40
En-Fr	26.94	37.53	29.56	21.97	24.42	32.01	41.99	32.83	29.51	31.28
En-Pl	26.65	37.69	29.33	21.70	24.47	32.04	42.05	32.82	29.47	31.26
En-Ru	26.73	37.50	29.35	22.03	24.55	32.08	41.95	32.84	29.51	31.36
Inflection										
En-Pl	26.95	37.58	29.41	21.68	24.59	32.07	41.98	32.87	29.49	31.40
En-Ru	26.80	37.63	29.31	21.90	24.43	32.06	42.04	32.85	29.51	31.30
Animacy										
De-En	26.80	37.43	29.32	21.84	24.65	32.05	42.05	32.84	29.48	31.41
Es-En	26.83	37.59	29.39	22.20	24.50	32.02	41.97	32.81	29.51	31.27
Fr-En	26.93	37.70	29.23	21.85	24.55	32.04	42.02	32.88	29.46	31.27
Pl-En	26.71	37.55	29.35	21.89	24.46	32.09	42.01	32.88	29.44	31.35
Ru-En	26.83	37.51	29.35	21.73	24.48	32.00	41.95	32.86	29.48	31.35

Table 8: BLEU scores for the models trained on datasets with different densities of annotated examples in the multilingual setting on the test subsets of the OpenSubtitles 2018 datasets for all relevant language pairs.

Model	En-De	En-Es	En-Fr	En-Pl	En-Ru	De-En	Es-En	Fr-En	Pl-En	Ru-En
Baseline	0.8023	0.8459	0.8005	0.8171	0.8321	0.8182	0.8522	0.8192	0.8009	0.8086
Gender										
En-De	0.8025	0.8456	0.8001	0.8171	0.8325	0.8182	0.8522	0.8189	0.8011	0.8085
En-Es	0.8025	0.8462	0.8004	0.8172	0.8326	0.8181	0.8521	0.8193	0.8011	0.8086
En-Fr	0.8023	0.8456	0.8000	0.8172	0.8322	0.8182	0.8521	0.8192	0.8011	0.8085
En-Pl	0.8025	0.8458	0.8004	0.8176	0.8324	0.8182	0.8522	0.8193	0.8011	0.8084
En-Ru	0.8021	0.8456	0.7999	0.8168	0.8321	0.8182	0.8523	0.8189	0.8009	0.8086
Formality										
En-De	0.8023	0.8456	0.8002	0.8168	0.8324	0.8181	0.8522	0.8190	0.8010	0.8084
En-Es	0.8026	0.8455	0.8003	0.8171	0.8325	0.8182	0.8523	0.8191	0.8011	0.8087
En-Fr	0.8024	0.8458	0.8008	0.8173	0.8321	0.8183	0.8522	0.8192	0.8011	0.8087
En-Pl	0.8024	0.8456	0.8005	0.8176	0.8325	0.8185	0.8523	0.8192	0.8009	0.8085
En-Ru	0.8022	0.8456	0.8001	0.8171	0.8318	0.8183	0.8524	0.8190	0.8009	0.8085
Auxiliary										
En-De	0.8023	0.8458	0.8001	0.8171	0.8321	0.8185	0.8524	0.8190	0.8011	0.8085
En-Es	0.8024	0.8458	0.8006	0.8174	0.8327	0.8185	0.8522	0.8191	0.8011	0.8085
En-Fr	0.8025	0.8455	0.7999	0.8165	0.8322	0.8181	0.8521	0.8189	0.8010	0.8085
En-Pl	0.8026	0.8458	0.8001	0.8170	0.8321	0.8183	0.8522	0.8191	0.8009	0.8083
En-Ru	0.8024	0.8457	0.8001	0.8169	0.8326	0.8183	0.8520	0.8190	0.8009	0.8085
Inflection										
En-Pl	0.8025	0.8458	0.8004	0.8162	0.8323	0.8184	0.8522	0.8191	0.8010	0.8087
En-Ru	0.8021	0.8457	0.7999	0.8168	0.8309	0.8184	0.8523	0.8190	0.8010	0.8084
Animacy										
De-En	0.8026	0.8458	0.8003	0.8174	0.8324	0.8184	0.8524	0.8188	0.8010	0.8085
Es-En	0.8025	0.8459	0.8005	0.8171	0.8328	0.8184	0.8522	0.8191	0.8009	0.8086
Fr-En	0.8021	0.8458	0.8000	0.8168	0.8325	0.8181	0.8523	0.8191	0.8008	0.8083
Pl-En	0.8021	0.8456	0.8004	0.8171	0.8322	0.8183	0.8522	0.8192	0.8008	0.8085
Ru-En	0.8022	0.8455	0.8003	0.8172	0.8321	0.8182	0.8521	0.8189	0.8008	0.8083

Table 9: COMET scores for the models trained on datasets with different densities of annotated examples in the multilingual setting on the test subsets of the OpenSubtitles 2018 datasets for all relevant language pairs.

Model	BLEU	COMET	Gender	Formality	Auxiliary	ContraPro
Baseline	33.93	0.8431	60.52%	38.63%	6.81%	78.88%
Fine-tuning e=1	33.60	0.8416	66.79%	39.30%	6.30%	83.02%
Fine-tuning e=2	33.59	0.8416	67.49%	39.34%	6.37%	83.78%
Fine-tuning e=5	33.60	0.8415	68.20%	39.49%	6.48%	84.50%
Head-tuning h=1	33.89	0.8428	63.28%	38.64%	6.43%	82.61%
Head-tuning h=2	33.85	0.8427	64.04%	38.58%	6.44%	83.40%
Head-tuning h=3	33.80	0.8425	64.75%	38.27%	6.45%	84.36%
Weighting $\lambda=2$	33.94	0.8430	64.35%	39.14%	7.18%	83.10%
Weighting $\lambda=5$	33.83	0.8430	65.72%	39.48%	7.67%	84.63%
Weighting $\lambda=10$	33.74	0.8426	66.24%	39.81%	8.10%	85.11%
Adapted D&R None	33.95	0.8429	60.77%	38.17%	7.01%	78.66%
CoWord p=0.1	33.98	0.8435	60.54%	38.72%	7.79%	78.65%
CoWord p=0.2	33.95	0.8436	60.47%	38.72%	8.22%	78.52%
CoWord p=0.3	33.88	0.8433	60.29%	38.68%	8.59%	78.39%
MaxPCXMI e=1	33.71	0.8420	66.16%	41.11%	6.84%	82.95%
MaxPCXMI e=2	33.70	0.8418	66.86%	41.44%	6.99%	83.79%
MaxPCXMI e=5	33.62	0.8414	67.31%	41.82%	7.18%	84.39%

Table 10: Performance in terms of BLEU and COMET on IWSLT 2017 en-de testset and ctxPro and ContraPro accuracy for the different methods. Number of epochs is noted as "e", and CoWord Dropout probability as "p", number of tuned heads as "h", and weighting strength as " λ ".

Model	En-De	En-Es	En-Fr	En-Pl	En-Ru	De-En	Es-En	Fr-En	Pl-En	Ru-En
Baseline	26.50	37.68	29.39	21.98	24.49	32.04	41.99	32.84	29.42	31.35
Adapted D&R	26.50	37.00	29.48	22.00	24.44	32.05	42.01	32.88	29.50	31.30
CoWord p=0.1	26.72	37.48	28.86	21.89	24.27	32.10	41.97	32.77	29.41	31.31
CoWord p=0.2	26.45	37.31	29.27	22.01	24.25	32.05	41.88	32.75	29.35	31.30
CoWord p=0.3	26.58	37.61	29.48	21.95	24.15	32.11	41.82	32.68	29.28	31.22
MaxPCXMI e=1	26.00	37.04	28.71	21.23	24.02	31.89	41.78	32.73	29.35	30.76
MaxPCXMI e=2	26.04	37.02	28.59	21.34	23.90	31.81	41.81	32.71	29.31	30.68
MaxPCXMI e=5	26.09	36.93	28.74	21.29	23.85	31.78	41.65	32.62	29.22	30.46

Table 11: BLEU scores for the methods in the multilingual setting on the test subsets of the OpenSubtitles 2018 datasets for all relevant language pairs.

Model	En-De	En-Es	En-Fr	En-Pl	En-Ru	De-En	Es-En	Fr-En	Pl-En	Ru-En
Baseline	0.8023	0.8459	0.8005	0.8171	0.8321	0.8182	0.8522	0.8192	0.8009	0.8086
Adapted D&R	0.8026	0.8456	0.8000	0.8175	0.8322	0.8183	0.8522	0.8191	0.8011	0.8085
CoWord p=0.1	0.8023	0.8454	0.7994	0.8167	0.8317	0.8182	0.8521	0.8188	0.8006	0.8086
CoWord p=0.2	0.8015	0.8453	0.7994	0.8166	0.8316	0.8178	0.8518	0.8187	0.8002	0.8083
CoWord p=0.3	0.8014	0.8453	0.7990	0.8164	0.8313	0.8176	0.8516	0.8183	0.7996	0.8083
MaxPCXMI e=1	0.7990	0.8433	0.7963	0.8125	0.8296	0.8155	0.8501	0.8170	0.7988	0.8057
MaxPCXMI e=2	0.7987	0.8431	0.7958	0.8123	0.8296	0.8150	0.8499	0.8167	0.7982	0.8053
MaxPCXMI e=5	0.7974	0.8427	0.7947	0.8109	0.8285	0.8137	0.8490	0.8158	0.7970	0.8043

Table 12: COMET scores for the methods in the multilingual setting on the test subsets of the OpenSubtitles 2018 datasets for all relevant language pairs.