

Key ingredients for effective zero-shot cross-lingual knowledge transfer in generative tasks: learning rate is (almost) all you need

Anonymous ACL submission

Abstract

Zero-shot cross-lingual generation implies finetuning of the multilingual pretrained language model (mPLM) on a generation task in one language and then using it to make predictions for this task in other languages. Previous works notice a frequent problem of generation in a wrong language and propose approaches to address it, usually using mT5 as a backbone model. In this work we compare various approaches proposed from the literature in unified settings, also including alternative backbone models, namely mBART and NLLB-200. We first underline the importance of tuning learning rate used for finetuning, which helps to substantially alleviate the problem of generation in the wrong language. Then, we show that with careful learning rate tuning, the simple full finetuning of the model acts as a very strong baseline and alternative approaches bring only marginal improvements. Finally, we find that mBART performs similarly to mT5 of the same size, and NLLB-200 can be competitive in some cases. Our final models reach the performance of the approach based on data translation which is usually considered as an upper baseline for zero-shot cross-lingual generation.

1 Introduction

Multilingual pretrained language models (mPLMs) such as mBERT (Devlin et al., 2019), mBART (Liu et al., 2020), and mT5 (Xue et al., 2021) provide high-quality representations for texts in various languages and serve as a universal backbone for finetuning on language-specific task data. The latter, however, is not always available for a language of interest, providing motivation for studying *zero-shot cross-lingual* capabilities of mPLMs. In this setting, the model is finetuned on the task data in one *source* language, usually English, and then applied in a zero-shot manner to make predictions in another *target* language, seen only at the pretraining stage.

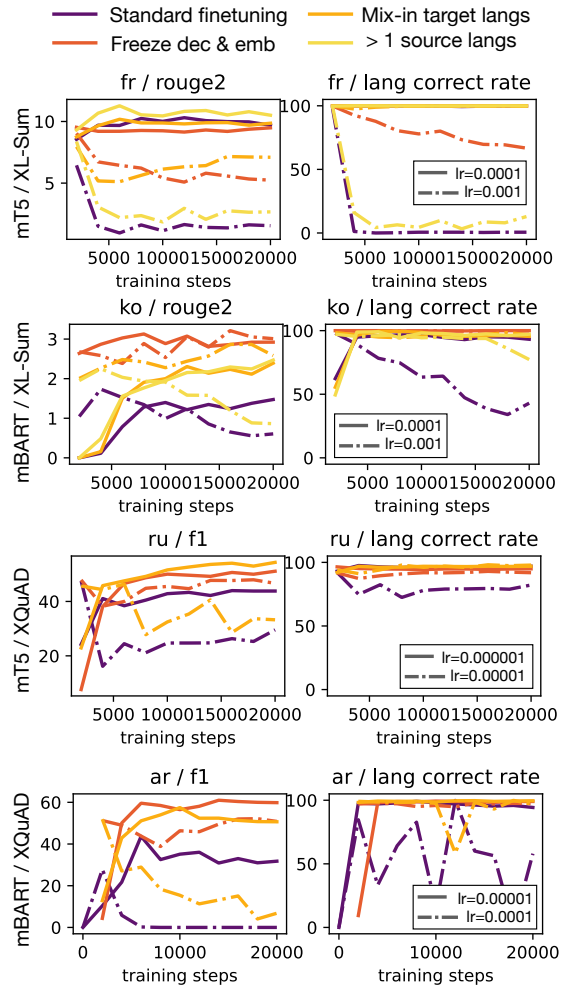


Figure 1: Learning rate plays a key role in cross-lingual transfer: decreasing LR almost completely eliminates generation in the wrong language with standard full finetuning, and often brings larger improvements that using complex adaptation methods developed to overcome this problem. Full results in Fig. 8–11 in Appendix.

While the described setting was broadly studied for natural language understanding tasks (Xue et al., 2021; Conneau et al., 2020; Artetxe et al., 2020a; Pires et al., 2019; Wu and Dredze, 2019; Pfeiffer et al., 2020), work on zero-shot cross-lingual gen-

048 *eration* is more limited (Vu et al., 2022; Pfeiffer
049 et al., 2023; Maurya et al., 2021; Li and Murray,
050 2023). Previous work highlight two main problems
051 arising in this scenario: producing incoherent or
052 irrelevant answers, and generating text in a wrong
053 language. A series of potential solutions were pro-
054 posed, such as freezing parts of the weights during
055 finetuning, utilizing parameter-efficient finetuning
056 methods, mixing-in unsupervised target language
057 data together with supervised source language data,
058 or using more than one source language. A com-
059 mon strategy is also to perform an intermediate
060 tuning of the model on the language generation
061 task in a self-supervised manner (as opposed to
062 denoising tasks used for pretraining).

063 However, despite listed efforts, the state of zero-
064 shot cross-lingual generation still remains unclear
065 and poses open questions:

- 066 • *Which adaptation method is most effective?*
067 Methods proposed for mitigating generation in
068 the wrong language, were all tested on different
069 tasks and benchmarks, and not compared to meth-
070 ods from other works, making it hard to establish
071 the best performing one.
- 072 • *What makes a better mPLM for zero-shot cross-*
073 *lingual transfer?* Different models have different
074 pretraining objectives, training and architectural
075 choices. How do those factors impact the quality
076 of the cross-lingual transfer in generation?
- 077 • *Importance of hyperparameters in downstream*
078 *task adaptation.* None of the previous work stud-
079 ied an impact of hyper-parameters used during
080 downstream task adaptation for zero-shot cross-
081 lingual generation.
- 082 • *Finally, if we pick the best solutions from all*
083 *of the three listed dimensions, how far in per-*
084 *formance can we get?* Can we reach the per-
085 formance of a strong baseline, data translation,
086 consisting in translating train data into target lan-
087 guage? Previous studies either did not reach its
088 performance or did not compare to this baseline.

089 The contribution of this work is conducting a
090 deep empirical study addressing the listed ques-
091 tions. We consider most commonly used multilin-
092 gual encoder-decoder mPLMs, namely mT5 and
093 mBART, as well as the translation model NLLB-
094 200. We systematically study six adaptation meth-
095 ods, investigate the effect of intermediate tuning,
096 pay attention to adaptation hyperparameters, and
097 compare models and adaptation methods *in a uni-*

fied setting. We consider two tasks: summarization
and questions answering (QA). Our main findings
are as follows:

- 098 • Hyperparameter tuning plays a very important
099 role in cross-lingual transfer: while most of the
100 works report severe problems with generation in
101 wrong language for mT5 with full finetuning, we
102 find that simply reducing learning rate helps to
103 alleviate this problem almost completely, without
104 hurting performance.
- 105 • Intermediate tuning substantially improves per-
106 formance in the majority of cases;
- 107 • With carefully chosen learning rates and interme-
108 diate tuning when necessary, simple full finetun-
109 ing is a very strong baseline in zero-shot cross-
110 lingual generation. Improvements brought by
111 more advanced methods are quite modest, and
112 none of the methods consistently outperform full
113 finetuning in all cases. The notable methods are
114 freezing model decoder and embeddings, which
115 performs consistently well with mBART (but not
116 with mT5), and using more than one source lan-
117 guage, which performs consistently well with
118 mT5 (but not with mBART).
- 119 • mBART and mT5 of similar size lead to com-
120 parable performance. Qualitatively, due to
121 the specifics of masking pretraining objective,
122 mBART is better suited for tasks with long out-
123 puts while mT5 is for tasks with short outputs.
- 124 • NLLB-200 is surprisingly competitive in sum-
125 marization, reaching performance of mT5 and
126 mBART for high-resource Latin-alphabet lan-
127 guages, but lags behind in QA.
- 128 • The final performance of cross-lingual genera-
129 tion reaches or outperforms the data translation
130 approach, often considered as an upper bound
131 for zero-shot cross-lingual generation. Notably,
132 careful learning rate tuning coupled with inter-
133 mediate tuning allow mT5 closely approach the
134 performance of data translation simply with full
135 finetuning adaptation.

136 2 Related Work 137

138 All works on zero-shot cross-lingual generation
139 underline (and try to address) the severe problem
140 of generating in a wrong language at the test time.
141 This problem is also referred to under terms cata-
142 strophic forgetting (of languages not participating
143 in finetuning, Vu et al., 2022), source language
144 hallucination (Pfeiffer et al., 2023), or accidental
145

translation problem (Li and Murray, 2023). Vu et al. (2022) propose to overcome generation in a wrong language by using parameter-efficient finetuning instantiated by prompt-tuning (Lester et al., 2021). They also mix-in the unsupervised target language task together with the supervised source language task, and factorize learnable prompts into language and task components.

Pfeiffer et al. (2023) propose mmT5 (modular mT5), allocating a small amount of language-specific parameters in the model during pretraining and freezing them during task-specific finetuning. To alleviate generation in a wrong language, they freeze some additional mmT5 parameters during finetuning, e. g. embedding layer and feed forward layers in Transformer decoder. Li and Murray (2023) argue that learning language-invariant representations during finetuning is harmful for cross-lingual generation and propose finetuning on data from more than one source language to avoid generation in a wrong language, with mT5 as a backbone model. ZMBART (Maurya et al., 2021) is the only work which considers other backbone model than mT5: they perform an intermediate tuning of mBART on an auxiliary unsupervised task on Hindi, Japanese and English. To avoid generation in a wrong language, they freeze embeddings and Transformer decoder, and mix-in data from auxiliary pretraining during finetuning.

In our work we are interested to compare all previously proposed approaches in the unified settings to better assess the impact of different factors on the zero-shot cross-lingual transfer for generation.

Alternative approaches to zero-shot cross-lingual transfer include data translation approaches, often referred as *translate-train* and *translate-test* paradigms. The former one implies translating train task data to the target language and finetuning the model on this translated data, and the latter one assumes translating test input examples into the source language, generating outputs in the source language and translating them back into the target language. The drawbacks of these approaches include a high computational cost either at training or testing time, lack of high-quality translation models for low-resource languages, and potential inconsistencies between sentences in translation (Vu et al., 2022). Despite its computational cost, data translation is a strong baseline which is usually considered as an upper bound on cross-lingual generation. Another related field is *few-shot cross-lingual generation* which assumes access to a

small amount of labeled examples in the target language (Schmidt et al., 2022; Lauscher et al., 2020; Zhao et al., 2021). This setting is out of scope of this study, but could be considered in the future work.

3 Methodology and experimental setup

Adaptation methods. We investigate the following adaptation methods:

- *Full finetuning*: all weights of the model are finetuned on the source language data;
- *Prompt tuning* (Vu et al., 2022): comprises prepending several learnable vectors ("prompt") to the list of embeddings of text input and freezing all other model weights during finetuning. Parameter-efficient approaches were shown in the literature to be better suited for transfer learning than full finetuning.
- *Adapters* (Houlsby et al., 2019; Bapna and Firat, 2019): lightweight tuned modules inserted after each fully-connected and attention block of Transformer, when the rest of (pretrained) model weights are frozen. We consider adapters as the most widely used parameter-efficient adaptation approach in the literature;
- *Freezing of encoder and embeddings* (Maurya et al., 2021): only weights in the encoder are finetuned. The motivation behind this approach is that the decoder should keep capabilities of generating in various languages while the encoder will adapt the model to the task;
- *Mixing-in self-supervised data for target languages* (Lester et al., 2021; Maurya et al., 2021): during finetuning, task data instances in source language will be alternated with self-supervised data instances in target language. The motivation is that such a mixing will preserve model's capability of generation in target languages;
- *Using several source languages* (Li and Murray, 2023): performing finetuning on more than one source language to better decouple task knowledge from language knowledge.

In the rest of the text term "full finetuning" refers to the finetuning of all weights on the English task data, even though two last described methods also finetune all weights. We do not consider mmT5 as it was not publicly released and requires substantial resources for pretraining.

We also experiment with *intermediate tuning* (IT) of the model, used in several works and per-

formed before finetuning on the task data. Standard encoder-decoder mPLMs rely on a self-supervised denoising training, where often input corresponds to corrupted text (eg. masked tokens), and output can follow some very specific structure (eg. unmasked span rather than full sentence, output containing special tokens, etc.). Therefore, in their raw form, these mPLMs are not necessarily well suited to receive well-formed text as an input and generate clean text as an output. IT performs finetuning on language modeling-like tasks, e.g. predicting the continuation of a paragraph based on its beginning, to compensate for this gap. IT was shown to be necessary in [Vu et al. \(2022\)](#) with prompt tuning of mT5 and in [Maurya et al. \(2021\)](#) with full or partial finetuning of mBART. We systematically test the necessity of IT for all methods and models.

Models. We focus on encoder-decoder mPLMs as they are well suited for generation purposes, as opposed to encoder-only mPLMs such as mBERT or XLM-R. We leave the investigation of decoder-only mPLMs such as BLOOM ([Scao et al., 2022](#)) for future work. We consider mT5 and mBART as two most widely used mPLMs and NLLB-200 as a high-quality translation model:

- *mT5*: pretrained using the masked language modeling objective where parts of the input sequence are masked and the missing fragments act as targets¹. mT5 is pretrained on the mC4 corpora, supports 101 languages, and does not use any language codes. Among released sizes from 300M to 13B we experiment with mT5-base (580M, most of the experiments) and mT5-Large (1.2B, additional experiment).
- *mBART (pt)*: pretrained using the denoising objective where parts of the input sequence are masked and the entire original sequence acts as a target ([Liu et al., 2020](#); [Tang et al., 2021](#)). mBART is pretrained on Common Crawl ([Conneau et al., 2020](#)) corpora, supports 50 languages, has 680M parameters in total and uses language codes in both encoder and decoder sides. Both input sequence X and target sequence Y are prepended with the language code: [lang_code, X] and [lang_code, Y], and at the inference time lang_code is forced as a first generated token. Our hypothesis is that the use of the language

code in the decoder can help to alleviate the problem of generation in a wrong language.

- *mBART (tr)*: In addition to the *pretrained* version, we also consider mBART finetuned for *translation* ([Tang et al., 2021](#)).
- *NLLB-200*: translation model supporting 200 languages, pretrained on sentence-level data mined from the web and automatically paired using multilingual embeddings. NLLB-200 uses the same language code scheme as mBART and is released in various sizes from 600M to 54.5B, among them we consider 600M (distilled version). Our hypothesis is that translation-based pretraining may provide good representations for cross-lingual transfer as suggested by ([Reid and Artetxe, 2023](#)).

Evaluation. We select two generative tasks to evaluate cross-lingual zero-shot knowledge transfer:

- *XL-Sum*: news summarization on the XL-Sum dataset ([Hasan et al., 2021](#)). The model needs to generate a 1–2 sentences summary based on a 1–2 news paragraphs. The evaluation is performed with ROUGE-2 metric ([Lin, 2004](#)) computed on the test sets (first 2k examples per language).
- *XQuAD*: question answering dataset ([Artetxe et al., 2020b](#)), the model needs to generate a short phrase answer based on a paragraph and question about it appended in the end of the paragraph. The evaluation is performed with F-measure comparing tokens in the gold answer and model-generated answer computed on publicly available development sets. For better metrics interpretability, we only consider questions for which groundtruth answers do not contain numbers and are correctly identified to be written in the target language.

We select a representative subset of languages for each task², covering Latin- and non-Latin scripts, and report how do task-specific metrics evolve during adaptation. For better interpretability, in addition to task metrics, we also consider (1) *lang. correct rate* metric (the percentage of outputs generated in the correct target language) and (2) *average sequence length* metric that allow to spot some edge behaviour of the models.

¹In contrast to English-centric T5, mT5 did not include supervised tasks in pretraining.

²XL-Sum: Chinese, French, Korean, Russian, and Spanish. XQuAD: Arabic, Chinese, German, Russian, and Spanish

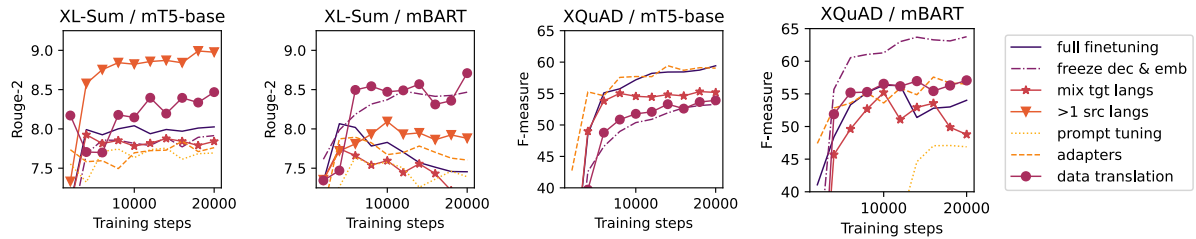


Figure 2: Comparison of adaptation methods, with tuned learning rates and intermediate tuning when it is needed. Results averaged across target languages and 2 runs. Language correct rate is close to 100% in almost all cases, due to hyperparameter tuning. The exception is prompt tuning of mT5 in the XQuAD task which is not shown because of too low performance.

Adaptation settings. For all adaptation methods we train models on English data for 20k steps with batch size of 4000 tokens on a single A100 GPU, and run evaluation each 2k steps. We crop input sequences to the maximum length supported by models, which equals to 512 (mT5, NLLB-200) or 1024 tokens (mBART). We grid search the learning rate (LR) for each task-model-adaptation method combination, details are given below.

For *Intermediate tuning* we finetune models for 100k steps on the CommonCrawl data with the batch size of 5k tokens and the LR chosen to optimize fluency of model generations, inspected manually. We use PrefixLM-inspired self-supervision from (Vu et al., 2022), where the continuation of the text needs to be predicted based on its beginning. It has shown more promising results in our preliminary experiments compared to self-supervised objective from (Maurya et al., 2021) (see details in Appendix B).

- *Prompt tuning*: we use the prompt dimension of 100 and initialize the prompt with randomly selected rows of the embedding matrix, following Vu et al. (2022).
- *Adapters*: we use the adapter dimension of 64 and insert adapters after each attention and fully-connected layer, following Bapna and Firat (2019).
- *Mixing-in target languages*: we use the same self-supervised objective as in IT and sample the corresponding data with probability 1% (all languages represented uniformly within this 1%), following Vu et al. (2022). We experimented with higher portions in Appendix C, as well as with mixing-in the pretraining task of the base model, and found that they lead to worse results.

- *Using several source languages*: we test this approach only on XL-Sum, because for XQuAD only English training data is available; for XL-Sum we use English, Japanese and Arabic, selecting them uniformly when forming mini-batches. More details on the experimental setting are given in Appendix A.

Hyperparameter tuning. We tune LR and decide on the necessity of IT, for each considered task-model-adaptation method combination. We initially grid searched LR for full finetuning, adapters and prompt tuning, for each task and model, without IT. The result of this step is the preliminary LR (PLR), and we utilize the PLR of full finetuning for other adaptation methods since they are also based on full finetuning. PLR usually corresponds to the highest LR which still enables generation in the correct language. After finding PLR, for each task-model-adaptation method combination, we select the best of four hyperparameter combinations: two options for LR (PLR and PLR $\times 10$) and two options for IT (used or not). Our intuition is that the use of advanced adaptation methodology or IT could potentially increase the LR which still does not lead to generation in the wrong language. In practice, this happened only once, for freezing of mBART in the summarization task. For XL-Sum, we perform the described tuning on the validation sets, looking at the performance averaged over considered target languages. For XQuAD, we use held-out languages (Thai, Romanian, and Vietnamese), since publicly available validation sets are used for the main evaluation. Results are usually consistent between languages.

We report the resulting optimal setting in Table 4 in Appendix. We could not find information on the used LR in (Pfeiffer et al., 2023) and (Vu et al., 2022), to compare our chosen LR with theirs. Maurya et al. (2021) and Li and Murray (2023) use

a constant LR for all tasks, which are hard to compare to ours because of different data³.

4 Experiments

First, we investigate the effect of learning rate, intermediate tuning and adaptation method for two most commonly used models, mT5 and mBART. Second, we compare them with other models and consider larger models. Finally, we present some qualitative examples and observations from manual inspection of predictions. In general, model predictions reaching highest metric values in our plots, form quite meaningful and reasonable responses to the considered tasks; more details in Section 5.

Effect of learning rate. We begin our study with analysing the effect of LR on the full finetuning on the English task data. With too small or too large LR the model does not learn even the English task because of too short steps or divergence. For the range of LRs when the English task is learned well, we observe that larger LRs lead to the effect reported in other works, when the model overfits to the source English language and generates answers in English when applied in cross-lingual setting. However, *with the reduced LR, this effect almost completely eliminates and the model mostly generates in the target language.* This effect is demonstrated in Figure 1 on a subset of languages and in Fig. 8–11 in Appendix on all considered languages.

Figure 1 also shows a comparison of enhancements of full finetuning proposed in the literature, such as mixing-in target language or freezing the decoder and the embedding. Even though these enhancements improve performance and percentage of outputs in the correct language, with fixed LR, we find that *reduced LR in full finetuning settings often brings larger improvements.* Reducing LR for other methods makes them even stronger.

We note that performance in English is usually a little higher with larger LR. This may raise a hypothesis that for non-English languages, outputs generated with larger LR in English may be of higher semantic quality than the ones generated in the correct target language with smaller LR. In Appendix D we test this hypothesis and demonstrate that this is not the case.

Effect of intermediate tuning. For each combination of a task and adaptation method, we com-

³Maurya et al. (2021) use LR=3e-5 larger than ours 1e-6, Li and Murray (2023) use LR=7e-5 close to ours 1e-4.

Method	XL-Sum		XQuAD	
	mT5	mBART	mT5	mBART
Full finetuning	+0.1	+2.5	+6.3	+9.0
Ft + mix tgt langs	0	+0.6	+3.1	-8.3
Ft + >1 src langs	0	+1	n/a	n/a
Freeze emb & dec	+4.3	+4.1	+11.2	+1.3
Adapters	0	0	+1.0	+3.9
Prompt tuning	+7.5	+7.2	+26.8	+25.1

Table 1: Difference in performance between task adaptation with and without intermediate tuning, for various methods. Rouge-2 for XL-Sum, F-measure for XQuAD.

pare the mT5-base/mBART task adaptation with and without intermediate tuning (IT).

We choose the best LR between PLR and PLR $\times 10$ (section 3). Results are presented in Table 1. We observe that *intermediate tuning substantially increases performance in the majority of cases.* In particular, IT appears to be essential for mBART with almost all adaptation methods and in all tasks, and important for mT5 in question answering. For mT5 in summarization, the use of IT does not increase performance, except with prompt tuning and freezing methods. We believe that this is because these two approaches do not modify the decoder, which was trained only on masked spans during mT5 pretraining and never was exposed to realistic text, and IT closes this gap. This result is consistent with (Vu et al., 2022) and (Maurya et al., 2021).

Comparison of adaptation methods. Figure 2 shows results (averaged over target languages) comparing adaptation methods for mT5-base and mBART models. Detailed per-language results are presented in Figure 7 in Appendix.

We observe that *with carefully chosen learning rates and intermediate tuning, simple full finetuning is a very strong baseline in zero-shot cross-lingual generation.* Improvements brought by the use of more advanced adaptation methods are rather modest, and *none of the adaptation methods consistently outperform full finetuning in all cases.* The notable approach for mBART is freezing the decoder and embeddings, proposed by Maurya et al. (2021) for this base model: freezing consistently outperforms full finetuning in all target languages in both tasks. However, this approach does not show such improvements for mT5. For XL-Sum, using more than one source language proposed in (Li and Murray, 2023) brings consistent improvement over target languages for mT5. For mBART this approach performs on par with using

one source language. The obvious drawback of this approach is that multi-lingual data may be not available, e.g. this is the case for XQuAD.

Mixing-in unsupervised tasks for target languages often degrades performance and increases the length of predictions, see Appendix C. Prompt tuning often has difficulties learning an English task and substantially underperforms other adaptation methods on XQuAD. Adapters usually perform on par or slightly worse than full finetuning.

Comparison of models. Figure 2 allows us to compare mT5-base and mBART after tuning of hyperparameters and adaptation methods. These models incorporate comparable numbers of parameters. We observe that *mT5 and mBART reach the similar level of performance in both tasks*. The same conclusion holds if we simply compare full finetuning runs of both models.

In Figure 3 we compare all four models we consider, adapted using full finetuning. We compare models without intermediate tuning, to avoid hindering model capabilities behind this additional step. We find that *translation-pretrained NLLB-200 performs well in summarization*, achieving performance of mT5 and mBART in Latin-language high-resource languages, French and Spanish, and performing on par with mBART without intermediate tuning in other languages⁴. We selectively inspected the predictions of NLLB and found that they indeed form meaningful summaries. However, in QA, NLLB-200 performs poorly, often (but not always) generating non-relevant answers. Translation-finetuned version of mBART performs poorly in all tasks, generating a lot of wrong language predictions.

Comparison versus data translation. Figure 2 also shows comparison versus the data translation⁵ approach, when English training data is translated into target languages using the NLLB-3.3B model. We translate data sentence-by-sentence and grid search the LR for finetuning. The results show that after careful tuning, *zero-shot cross-lingual generation reaches or outperforms the data translation approach in both considered tasks*. If we consider a simpler setting when only LR and the use of IT are tuned, i.e. comparing full finetuning and data translation runs in Figure 2, we observe that zero-

⁴Expect Chinese, for which NLLB-200 generates a lot of empty predictions. NLLB-200 was noticed previously in the literature to have issues with processing Chinese.

⁵Data translation is often referred as translate-train method.

Method	XL-Sum		XQuAD	
	R2	LCR	F1	LCR
Large / IT + ft	9.9	99.8%	69.8	94.7%
Large / IT + ft >1 src lg	10.9	99.8%	n/a	n/a
Large / Data translation	10.8	99.8%	63.6	96.7%
Base / IT + ft	8.0	99.7%	59.4	92.9%
Base / IT + ft >1 src lg	9.0	99.8%	n/a	n/a
Base / Data translation	8.5	99.6%	53.9	95.3%

Table 2: Results for mT5-large model, averaged over target languages. Metrics: Rouge-2 for XL-Sum, F-measure for XQuAD, LCR: language correct rate. LCR is lower than 100% on XQuAD (partly) because of language identification errors for short sequences.

shot cross-lingual generation closely approaches the data translation approach in summarization and performs the same in question answering. The XQuAD dataset is harder to automatically translate than XL-Sum, e.g. single words often present in targets may be translated into short full sentences.

Experiments with larger models. Table 2 reports results for the mT5-large model where we compare performance achieved with full finetuning after intermediate tuning versus training on translated data. We also include the leader approach of using several source languages. We consider only mT5 because mBART is released in one size. We reduce LR to 0.00001 for the larger model, as the LR of 0.0001 used for the base model was sometimes producing English outputs. We also list mT5-base results for reference.

We find that the same conclusions hold for the mT5-large model as for mT5-base: reducing LR eliminates generation in the wrong language, and the zero-shot cross-lingual model is on par or better than the data translation approach.

5 Inspection of predictions

We inspected a subset of predictions in the languages we speak and found that models achieving highest scores in both tasks generate fluent, meaningful and reasonable predictions in a lot of cases, but sometimes have issues with factuality, grammaticality or hallucinations. Examples are shown in Figure 4. Analyzing effects of LR, we observe that increasing LR leads first to increase in code switching and then to wrong language generation, while *reducing LR leads to producing rudiments of pretraining in generation*. For example, models sometimes generate extra tokens used in pretraining, such as <extra_id_{N}>

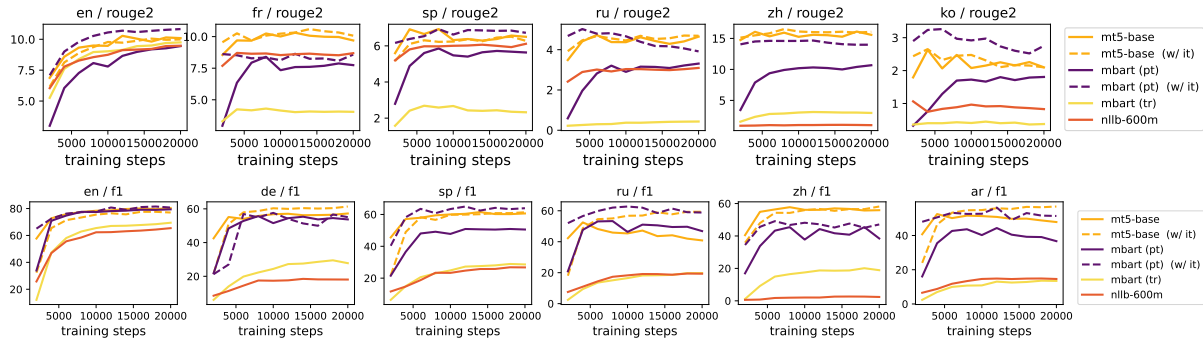


Figure 3: Comparison of base models with full finetuning. Each plot averaged over 3 runs. Correct language rate is close to 100%, due to hyperparameter tuning, in almost all cases except the translation-tuned version of mBART. pt: pretrained version of mBART, tr: translation-finetuned version of mBART.

Model	Model outputs in French	Translation into English	Avg len	
1	Grountruth	Pr Hulda Swai, professeure tanzanienne de sciences de la vie et de bio-ingénierie, a été désignée lauréate 2020 du prix continental Kwame Nkrumah de l'Union africaine (UA) pour l'excellence scientifique.	Pr. Hulda Swai, a Tanzanian female professor of life sciences and bioengineering, has been named the 2020 winner of the African Union (AU) Kwame Nkrumah Continental Prize for Scientific Excellence.	151
2	mT5-base, full ft w/o IT	<extra_id_0> femmes sont très motivées à entreprendre une carrière scientifique dans le domaine de la nanotechnologie.	<extra_id_0> women are highly motivated to pursue a scientific career in the field of nanotechnology.	101
3	mT5-base, full ft with IT	La professeure de nanotechnologie de l'Union africaine a réussi à réunir 7,5 millions de dollars.	The nanotechnology female professor from the African Union managed to raise \$7.5 million.	91
4	mT5-base, >1 src langs w/o IT	<extra_id_0> une scientifique africaine a reçu une récompense pour son travail dans le domaine de la nanotechnologie.	<extra_id_0> an African female scientist has received an award for her work in the field of nanotechnology.	107
5	mT5-base, train data translation	Un professeur de nanotechnologie a reçu une récompense continentale pour son travail dans le domaine des nanotechnologies.	A nanotechnology professor has received a continental award for her work in the field of nanotechnology.	112
6	mBART, full ft w/o IT	Ancienne professeure de l'université de Durban a reçu un prix de la part de la Banque mondiale.	Former professor at the University of Durban received an award from the World Bank.	117
7	mBART, full ft with IT	A ne pas manquer sur BBC Afrique : Une femme motivée et concentrée	Not to be missed on BBC Africa: A motivated and focused woman	111
8	mBART, freeze dec & emb, with IT	La professeure africaine de nanotechnologie a été lauréate du prix Kwame Nkrumah de l'Union africaine.	The African nanotechnology female professor was the recipient of the African Union Kwame Nkrumah Prize.	115
9	mBART, train data translation	Un scientifique africain a été lauréat du prix Kwame Nkrumah de l'Union africaine.	An African scientist has been awarded the African Union Kwame Nkrumah Prize.	108

Figure 4: Example predictions for a selection of models. Avg. len. over evaluation corpora in French, in characters. Red highlights errors or extra tokens.

for mT5 or <sep> for mBART, see rows 2 and 4 in Figure 4. In most cases this does not affect meaningfulness of predictions, but in rare cases leads to mT5 producing incomplete sentences, which may look unreasonable in summarization, e.g. “<extra_id_0> Guinea-Bissau President Alberto Dabo said.” (translated from French). The reason is that in mT5 pretraining tokens <extra_id_{N}> were followed by fragments of input sentences. The described effect is eliminated by intermediate tuning (row 3 in Fig. 4).

In the same fashion, mBART average lengths are closer to groundtruth average lengths than mT5 in summarization, and the reverse effect takes place in QA. The reason is that in mT5 pretraining, the targets are only fragments masked in the input, which are shorter than targets in mBART pretraining represented by full sequences (they need to be reconstructed from the masked inputs).

Notably, data translation can produce translation-related errors, e.g. in rows 5 and 9 models generate a wrong male article "Un", probably because this

was a dominating article in the translated data.

6 Conclusion

In this work, we conducted a deep systematic study of how to achieve high-performing zero-shot cross-lingual generation. Our study highlights the high importance of careful learning rate tuning and the usefulness of the intermediate tuning. We show that with these two ingredients, mT5 and mBART achieves strong results with simple full finetuning, i.e. closely approach the performance of the translate-train approach in summarization and reaching it in question answering. The performance gap in summarization is closed by using several source languages in mT5 and freezing decoder and embeddings in mBART. Translation-pretrained NLLB-200 shows surprisingly good performance in summarization but lags behind in question answering. We suggest that future works report more rigorously their experimental setup and details on hyperparameter search, and consider wider spectrum of models and baselines in the experiments.

7 Limitations and broader impact

We aim at conducting a deep, thoughtful study of various design choices in zero-shot cross-lingual generation, but acknowledge the impossibility of considering all possible options, given the resource constraints. In particular, we could not perform full fine-grained grid search of LR for each task-model-adaptation method combination. Instead, we use a well-designed simplified strategy described in Section 3, which already gave strong results. In the same fashion, we had to limit our study to three models (we picked most commonly used models) and adaptation methods which do not require model pretraining, e.g. we do not consider mmT5 model. Nonetheless, we hope our study provides helpful insights on zero-shot cross-lingual transfer in generative tasks and shows that it can achieve the performance of the data translation method, which is usually considered as an unreachable upper bound.

We do not anticipate any negative impact of our work and on the reverse hope that it will help to make higher-quality language technologies accessible to a broader set of languages.

References

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020a. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.

Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020b. [On the cross-lingual transferability of monolingual representations](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4623–4637, Online. Association for Computational Linguistics.

Ankur Bapna and Orhan Firat. 2019. [Simple, scalable adaptation for neural machine translation](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1538–1548, Hong Kong, China. Association for Computational Linguistics.

Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [mbert blog post](#). <https://github.com/google-research/bert/blob/master/multilingual.md>.

Tahmid Hasan, Abhik Bhattacharjee, Md. Saiful Islam, Kazi Mubasshir, Yuan-Fang Li, Yong-Bin Kang, M. Sohel Rahman, and Rifat Shahriyar. 2021. [XL-sum: Large-scale multilingual abstractive summarization for 44 languages](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4693–4703, Online. Association for Computational Linguistics.

Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. [Parameter-efficient transfer learning for NLP](#). In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 2790–2799. PMLR.

Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, H erve J egou, and Tomas Mikolov. 2016. [Fasttext.zip: Compressing text classification models](#). *arXiv preprint arXiv:1612.03651*.

Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. [Bag of tricks for efficient text classification](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.

Anne Lauscher, Vinit Ravishankar, Ivan Vuli c, and Goran Glava s. 2020. [From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.

Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. [The power of scale for parameter-efficient prompt tuning](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Tianjian Li and Kenton Murray. 2023. [Why does zero-shot cross-lingual generation fail? an explanation and a solution](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 12461–12476, Toronto, Canada. Association for Computational Linguistics.

Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Yinhan Liu, Jiatao Gu, Naman Goyal, Xian Li, Sergey Edunov, Marjan Ghazvininejad, Mike Lewis, and

741	Luke Zettlemoyer. 2020. Multilingual denoising pre-training for neural machine translation . <i>Transactions of the Association for Computational Linguistics</i> , 8:726–742.	pages 3450–3466, Online. Association for Computational Linguistics.	798
742			799
743			
744			
745	Kaushal Kumar Maurya, Maunendra Sankar Desarkar, Yoshinobu Kano, and Kumari Deepshikha. 2021. Zm-BART: An unsupervised cross-lingual transfer framework for language generation . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> , pages 2804–2818, Online. Association for Computational Linguistics.	Tu Vu, Aditya Barua, Brian Lester, Daniel Cer, Mohit Iyyer, and Noah Constant. 2022. Overcoming catastrophic forgetting in zero-shot cross-lingual generation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9279–9300, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	800
746			801
747			802
748			803
749			804
750			805
751			806
752	Jonas Pfeiffer, Francesco Piccinno, Massimo Nicosia, Xinyi Wang, Machel Reid, and Sebastian Ruder. 2023. mmt5: Modular multilingual pre-training solves source language hallucinations .	Shijie Wu and Mark Dredze. 2019. Beto, bentz, becas: The surprising cross-lingual effectiveness of BERT . 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 833–844, Hong Kong, China. Association for Computational Linguistics.	807
753			808
754			809
755			810
756	Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. MAD-X: An Adapter-Based Framework for Multi-Task Cross-Lingual Transfer . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 7654–7673, Online. Association for Computational Linguistics.		811
757			812
758			813
759			814
760			
761			
762			
763	Telmo Pires, Eva Schlinger, and Dan Garrette. 2019. How multilingual is multilingual BERT? In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 4996–5001, Florence, Italy. Association for Computational Linguistics.	Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. mT5: A massively multilingual pre-trained text-to-text transformer . In <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.	815
764			816
765			817
766			818
767			819
768			820
769	Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text . In <i>Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing</i> , pages 2383–2392, Austin, Texas. Association for Computational Linguistics.	Mengjie Zhao, Yi Zhu, Ehsan Shareghi, Ivan Vulić, Roi Reichart, Anna Korhonen, and Hinrich Schütze. 2021. A closer look at few-shot crosslingual transfer: The choice of shots matters . 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 5751–5767, Online. Association for Computational Linguistics.	821
770			822
771			823
772			824
773			825
774			826
775	Machel Reid and Mikel Artetxe. 2023. On the role of parallel data in cross-lingual transfer learning . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5999–6006, Toronto, Canada. Association for Computational Linguistics.		827
776			828
777			829
778			830
779			831
780	Teven Le Scao, Angela Fan, Christopher Akiki, Ellie Pavlick, Suzana Ilić, Daniel Hesslow, Roman Castagné, Alexandra Sasha Luccioni, François Yvon, Matthias Gallé, et al. 2022. Bloom: A 176b-parameter open-access multilingual language model .		
781			
782			
783			
784			
785	Fabian David Schmidt, Ivan Vulić, and Goran Glavaš. 2022. Don't stop fine-tuning: On training regimes for few-shot cross-lingual transfer with multilingual language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 10725–10742, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		
786			
787			
788			
789			
790			
791			
792			
793	Yuqing Tang, Chau Tran, Xian Li, Peng-Jen Chen, Naman Goyal, Vishrav Chaudhary, Jiatao Gu, and Angela Fan. 2021. Multilingual translation from denoising pre-training . In <i>Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021</i> ,		
794			
795			
796			
797			

A Experimental setup

Data. We experiment with news summarization on the XL-Sum dataset (Hasan et al., 2021) (released under the CC BY-NC-SA 4.0 license) and question answering on the XQuAD dataset (Artetxe et al., 2020b) (released under the CC BY-SA 4.0 license). Both datasets were released for research purposes. The XL-Sum dataset was obtained by crawling BBC news in 44 languages, with corpus size per language varying from 1K (Scottish Gaelic) to 300K (English) article-summary pairs. Inputs are composed of 1–2 paragraphs and targets are usually 2–3 sentences. We evaluate on test sets and crop test sets larger than 2K samples, to 2K. The XQuAD dataset was obtained by translating SQuAD validation set (Rajpurkar et al., 2016) into 11 languages, thus all language corpora are parallel. We use this dataset for evaluation and train on the training set of SQuAD (80K training instances). Each input is composed of a paragraph and a question about this paragraph appended in the end of the paragraph. Each output is an answer to a question, a short segment copied from the paragraph.

Preprocessing and postprocessing. We tokenize data using each model’s tokenizer. We crop model inputs and outputs to the maximum lengths supported by models, which equal to 1024 tokens for mBART and 512 tokens for mT5-base and NLLB-600M. Due to the design of pretraining, models may generate extra tokens such as `<extra_id_{N}>` for or `<sep>` for mBART. We remove such extra tokens from predictions before computing metrics.

Models and training. We consider three models: mT5 (base and large, released under the Apache License 2.0 license), mBART (MIT license), and NLLB-200 (cc-by-nc-4.0 license). All models allow use for research purposes. We train models on English data for 20k steps with batch size of 4000 tokens on a single A100 GPU, and conduct validation on considered target languages each 2k steps. We use Adam optimizer with standard inverse square root LR schedule and warm up of 4k steps, and update model parameters after each mini-batch. We estimated the total computational budget of our experiments to be 4K GPU hours.

Hyperparameter search. For full finetuning, adapters and prompt tuning, we run a search over a range of LR. For each task and model (without intermediate tuning), we search the LR best for non

English languages on average, looking at ROUGE-2 for summarization and F-measure for QA. We start with the set of three LRs: $\{10^{-k}, k = 3, 4, 5\}$. If the optimal $k^* \neq 4$ then we extend search correspondingly to $k = 2, 1$ or $k = 6, 7$ until performance stops improving. For full finetuning, after we find optimal k^* we also consider $3 \cdot 10^{-k^*}$. The motivation is that the optimal k^* usually corresponds to the maximal k that still allows generation in the correct language, and considering $3 \cdot 10^{-k^*}$ enables more accurate search for this maximum. We report chosen LRs for full finetuning and adapters in Table 4. For prompt tuning we chose LR of 0.01 for both tasks.

Evaluation. For summarization, we report ROUGE metrics (Lin, 2004), and for QA, we report F-measure. In QA, a lot of answers contain numbers or English words which could inflate metrics even if the model does not generate in the correct language. Moreover, the accuracy of language identification decreases on short answers, resulting in false indication of generation in wrong language. To avoid these issues, we compute metrics in QA only over questions for which groundtruth answers do not contain numbers and are correctly identified to be written in the target language ($\sim 50\%$ of 1190 questions satisfy this criteria).

For ROUGE metric, we use the `gem-metrics` package. For F1 metric in XQuAD, we use the script provided by the dataset authors. To identify language, we use `fasttext` library (Joulin et al., 2017, 2016) and its `lid.176.bin` model⁶.

B Preliminary experiments with intermediate tuning

Figure 5 reports comparison of two self-supervised objectives for intermediate tuning: Prefix-LM and ZmBART-like objective. PrefixLM objective implies predicting the continuation of the chunk of text based on its beginning, while ZmBART-like objective implies citing random sentences from the input chunk of text. We compare two objectives using the freezing of the decoder and embeddings as an adaptation method, applied after intermediate tuning with the chosen objective, because we found intermediate tuning to be essential for this adaptation method in the preliminary experiments. Finetuning LR equals to the PLR defined in Section 4, intermediate tuning LR was chosen to optimize

⁶<https://fasttext.cc/docs/en/language-identification.html>

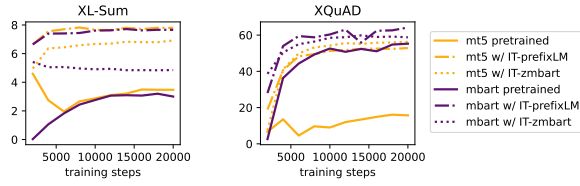


Figure 5: Comparison of self-supervised objectives for intermediate tuning, with freezing decoder and embeddings as an adaptation method. Task metric: Rouge-2 for XL-Sum, F1 for XQuAD. Correct language rate is close to 100% in all cases except pretrained mBART on XL-Sum.

930 fluency of model generations, inspected manually.
 931 Intermediate tuning is performed on the Common-
 932 Crawl dataset.

933 We observe that for XL-Sum, the Prefix-LM ob-
 934 jective leads to substantially higher Rouge-2 values,
 935 while for XQuAD both objectives lead to close re-
 936 sults. Based on these results, we decided to use the
 937 Prefix-LM objective in all experiments.

938 C Preliminary experiments with 939 mixing-in target languages

940 Figure 6 reports results of preliminary experiments
 941 with mixing-in a self-supervised task in target lan-
 942 guages. For each base model, namely mT5-base
 943 and mBART, we consider its pretraining task and a
 944 Prefix-LM task used for intermediate tuning. We
 945 consider several options for the probability of sam-
 946 pling target language examples when forming mini-
 947 batches. CommonCrawl data is used for the self-
 948 supervised task. The experiment is conducted for
 949 the XL-Sum task, with LR being equal to the PLR
 950 defined in Section 4, without intermediate tuning.

951 For mt5, we observe that using the span corrup-
 952 tion pretraining task leads to empty outputs with
 953 any mixing-in probability (with smaller probabili-
 954 ties this effect happens later in the training). This
 955 is because task examples do not contain any mask
 956 tokens, and empty generation is a default response
 957 of the pretrained mT5 to such inputs. Mixing-in
 958 PrefixLM task examples performs similarly to the
 959 standard finetuning of mT5, with mixing-in proba-
 960 bility of 1% performing best, same as in (Vu et al.,
 961 2022). Qualitatively, mixing-in self-supervised
 962 task increases the length of generated outputs in
 963 the tasks of interest.

964 For mBART, all mixing-in strategies lead to
 965 modest improvements in performance, with Pre-

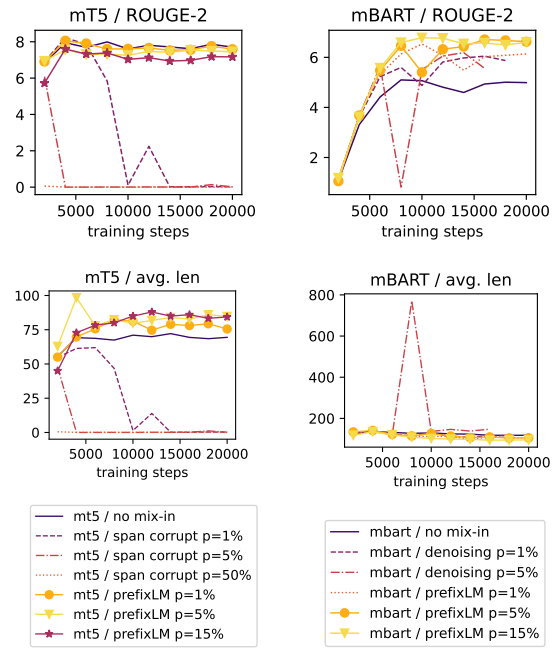


Figure 6: Preliminary experiments with mixing-in a self-supervised task for target languages. The probability in the legend denotes the probability of sampling target language examples when forming mini-batches. Two self-supervised tasks considered: Prefix-LM and the pretraining task of the model. Correct language rate is close to 100% in all cases

966 fixLM task performing slightly better. All consid-
 967 ered mixing-in probabilities lead to similar results.
 968 Based on these observations, we decided to use the
 969 PrefixLM task with mixing-in probability of 1% in
 970 our experiments.

971 D Additional experiment with translating 972 English outputs into target languages

973 When reducing the LR for preserving generation
 974 in correct language, a reasonable question could be
 975 whether predictions of higher LR models are higher
 976 quality answers, but just in the wrong language, or
 977 simply hallucinations caused by data distribution
 978 shift. The premise for the former scenario is that
 979 on English data, performance with our chosen LR
 980 is usually slightly lower than with a larger LR.

981 We find that actually the later scenario takes
 982 place, by comparing performance of our chosen
 983 LR (best for non-English) and of the best LR for
 984 English with model predictions being translated
 985 into target languages using NLLB-3.3B⁷, for last

⁷NLLB-3.3B handles well inputs containing code switch-
 ing which are frequent in predictions we are translating, and

		Best-En LR + Tr.		Best-non-En LR	
		LR	Score	LR	Score
Sum	mT5	1e-3	4.02	1e-4	7.7
	mBART	1e-5	4.06	1e-6	5.34
	NLLB-200	1e-4	2.86	1e-5	4.62
QA	mT5	1e-4	46.2	1e-4	58.6
	mBART	1e-5	41.1	1e-5	46.6
	NLLB-200	1e-4	17.4	3e-5	18.2

Table 3: Comparison of best LR for non-English languages and best LR for English with model outputs being translated into target languages. Performance averaged over non-English languages, after 20k of full finetuning. Reported metric: Rouge-2 for summarization, F-measure for QA. mBART — pretrained version, no intermediate tuning is used in this experiment.

Model	Method	XL-Sum		XQuAD	
		LR	IT?	LR	IT?
mT5 (base)	Ft w/o IT	1e-4		1e-4	
	Ft	1e-4		1e-4	✓
	+ Mix tgt langs	1e-4		1e-4	✓
	+ >1 src langs	1e-4		n/a	
	Freeze	1e-4	✓	1e-4	✓
	Adapters	1e-3		1e-3	
	Prompt tuning	1e-2	✓	1e-2	✓
mBART	Ft w/o IT	1e-6		1e-5	
	Ft	1e-6	✓	1e-5	✓
	+ Mix tgt langs	1e-6	✓	1e-5	
	+ >1 src langs	1e-6	✓	n/a	
	Freeze	1e-5	✓	1e-5	✓
	Adapters	1e-5	✓	1e-3	✓
	Prompt tuning	1e-2	✓	1e-3	✓
NLLB	Ft w/o IT	1e-5		3e-5	
mBART (tr)	Ft w/o IT	1e-6		1e-3	

Table 4: Best learning rates for non-English languages. n/a: not applicable.

986 checkpoints of full models finetuning. Accord-
987 ing to Table 3, translated predictions of higher LR
988 model score lower than the (non-translated) pre-
989 dictions of lower LR model. This result further
990 advocates for the importance of careful LR tun-
991 ing for full finetuning in zero-shot cross-lingual
992 generation.

simply copies inputs which are already in the target language.

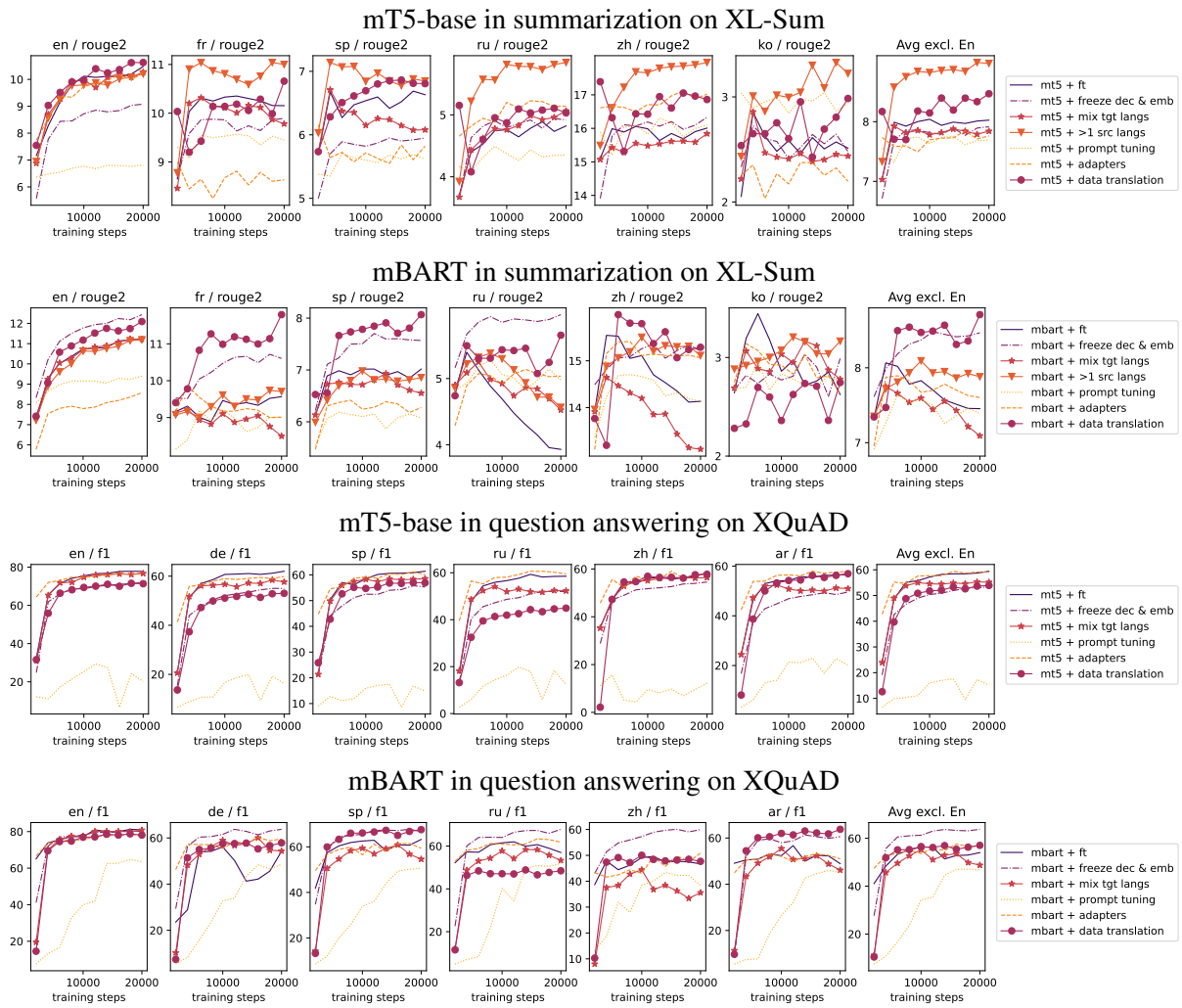


Figure 7: Per-language results on the comparison of adaptation methods. Each plot averaged over 2 runs. Correct language rate is close to 100% in all cases, due to the hyperparameter tuning, except prompt tuning of mT5 in the XQuAD task.

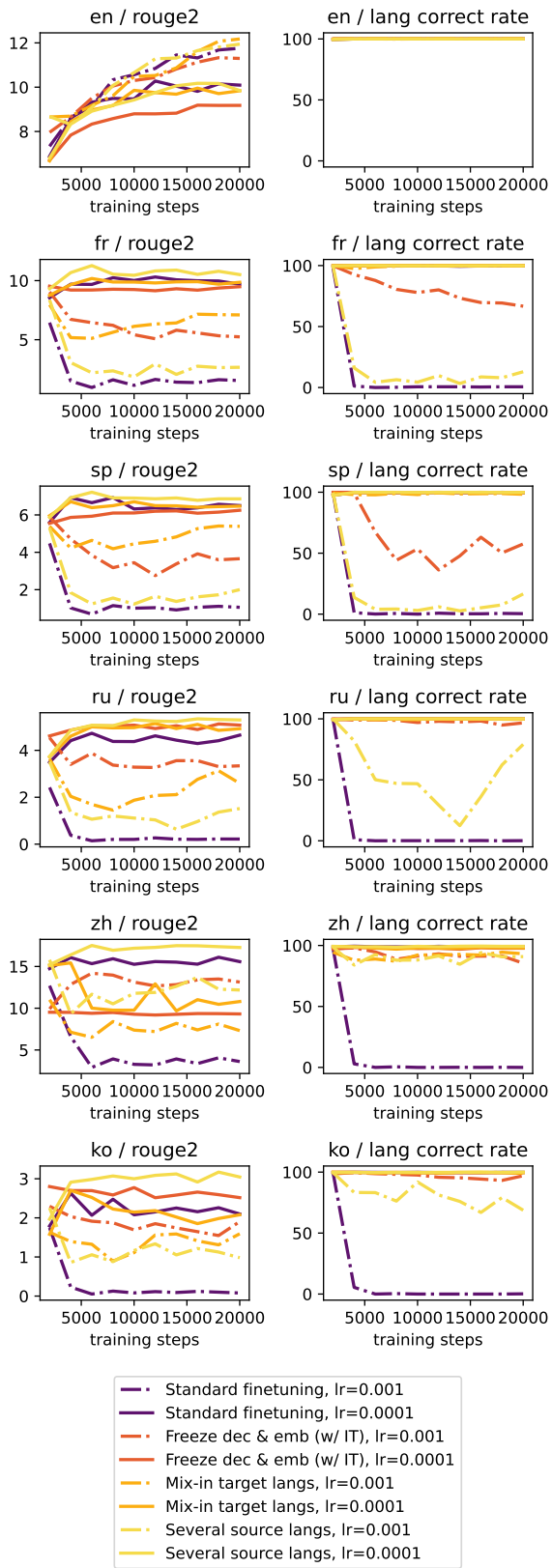


Figure 8: Per-language results on the effect of learning rate, for mT5 on XL-Sum.

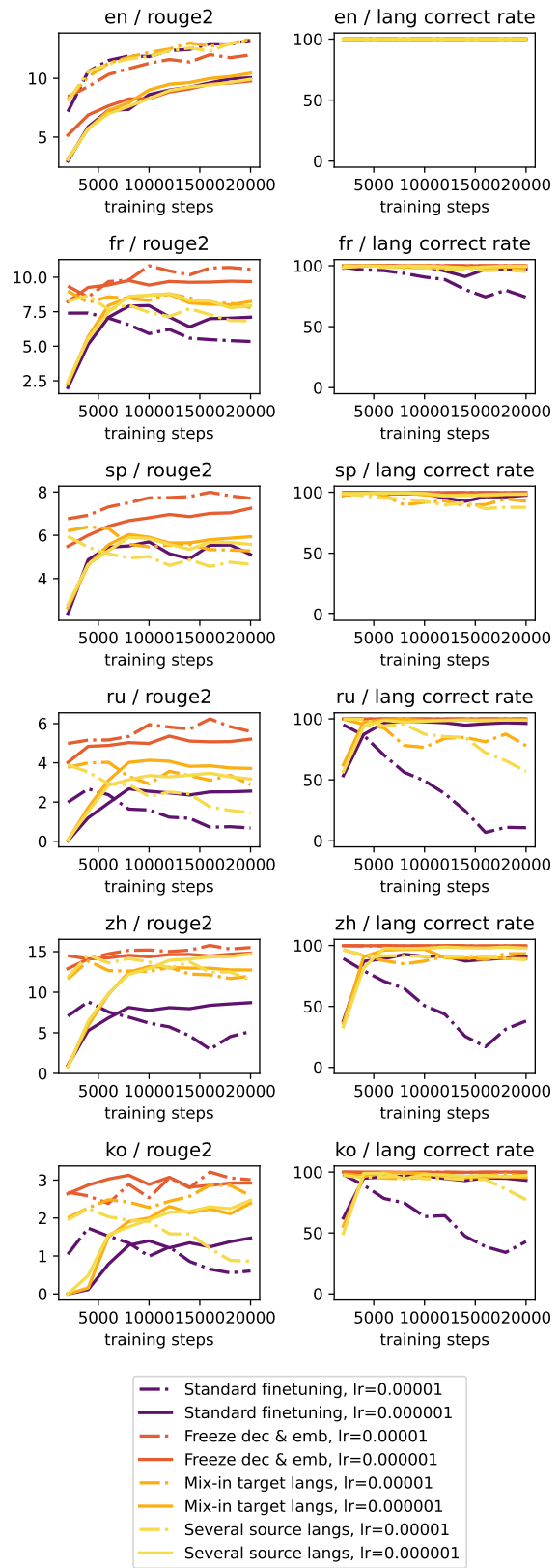


Figure 9: Per-language results on the effect of learning rate, for mBART on XL-Sum.

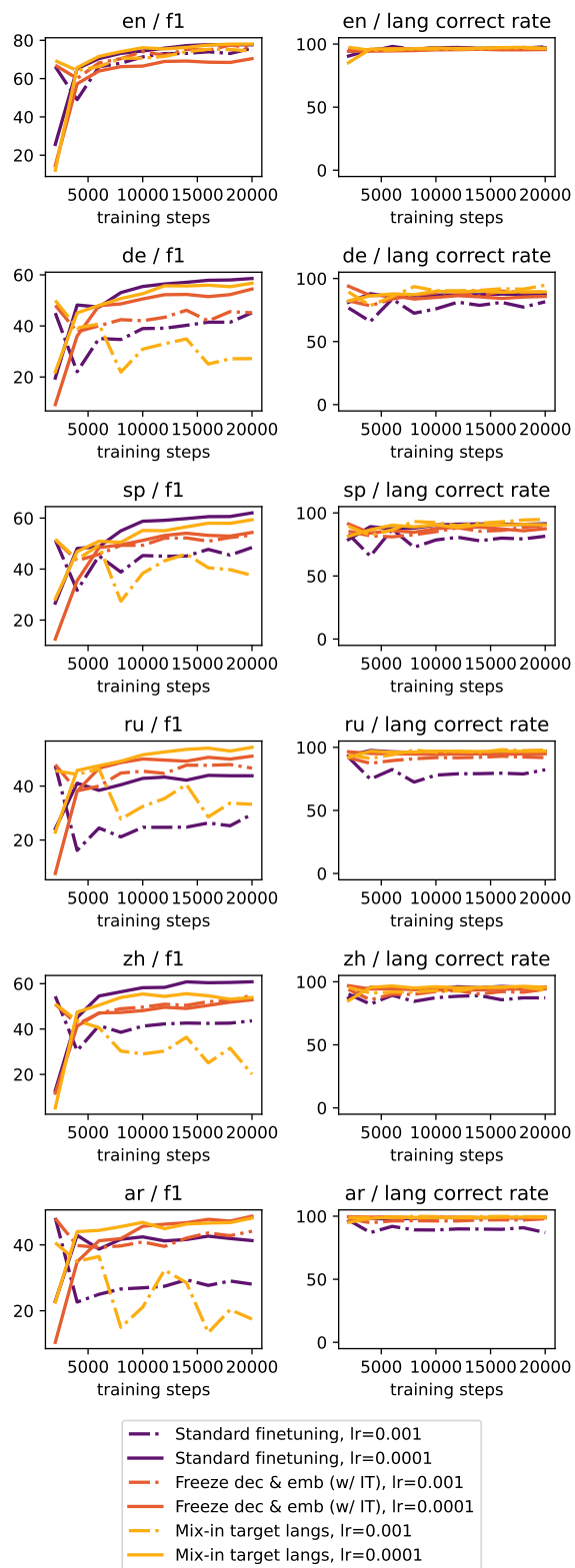


Figure 10: Per-language results on the effect of learning rate, for mT5 on XQuAD.

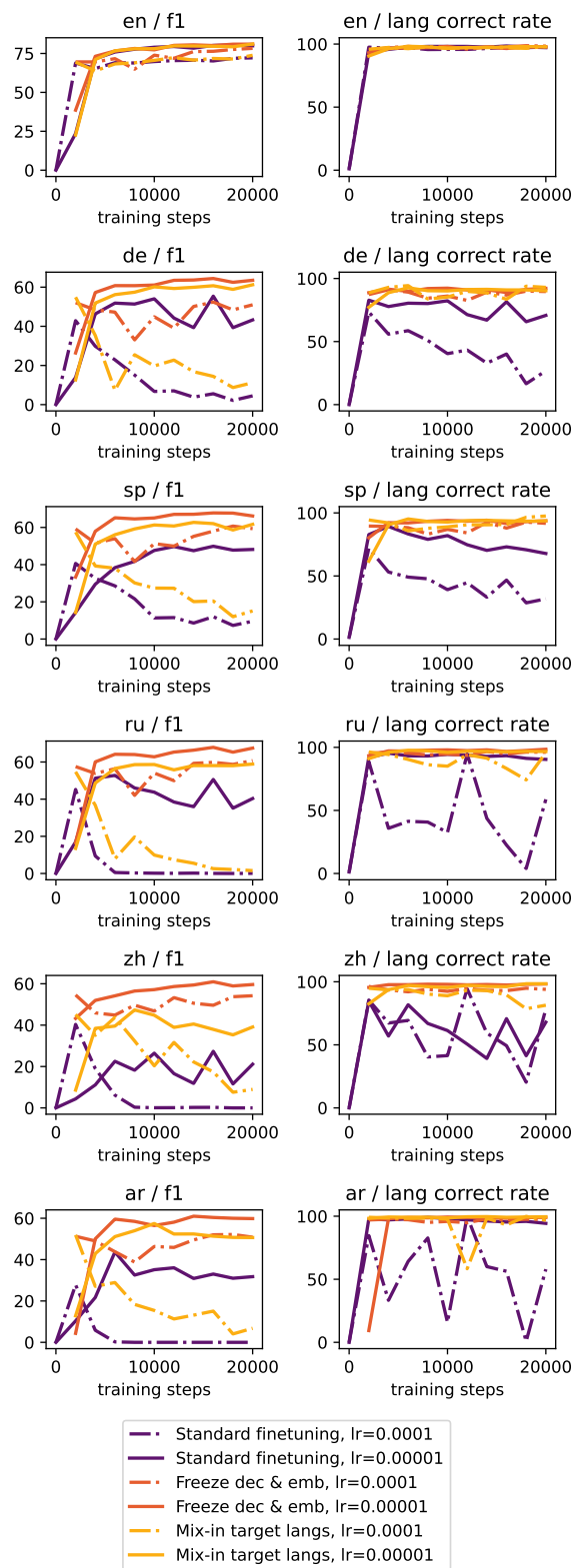


Figure 11: Per-language results on the effect of learning rate, for mBART on XQuAD.