

Continually Adding New Languages to Multilingual Language Models

Anonymous authors

Paper under double-blind review

Abstract

Multilingual language models are trained on a fixed set of languages. To support new languages, the models need to be re-trained from scratch. This is an expensive endeavor and is often infeasible, as model developers tend not to release their pre-training data. Naive approaches like continued pretraining suffer from catastrophic forgetting but mitigation strategies like experience replay cannot be applied due to lack of original pretraining data. In this work, we investigate the problem of continually adding new languages to a multilingual model assuming access to pretraining data in only the target languages. We propose Layer-selective LoRA (LAYRA), which adds Low-Rank Adapters (LoRA) to selected initial and final layers while keeping the rest of the model frozen. LAYRA builds on two insights: (1) LoRA reduces forgetting, and (2) multilingual models encode inputs in the source language in the initial layers, reason in English in intermediate layers, and translate back to the source language in final layers. Our experiments with adding multiple combinations of Galician, Swahili and Urdu to Llama 3.1 reveal the effectiveness of our approach across diverse multilingual tasks. We also demonstrate that using model arithmetic, the adapted models can be equipped with instruction following abilities without access to any instruction tuning data in the target languages.

1 Introduction

Although several recently released language models (LMs) are advertised as multilingual (Grattafiori et al., 2024; Faysse et al., 2024; Gemma, 2025), they only support a handful of predetermined high-resource languages. As resources for new languages become available, continually supporting them in such models is not trivial. Retraining them from scratch is often prohibitively expensive, so practitioners typically adopt an incremental continued pretraining (CPT) strategy to incorporate new languages (Csaki et al., 2023). However, it often results in catastrophic forgetting of previously supported languages (Cahyawijaya et al., 2023; Chalkidis et al., 2021; Vu et al., 2022).

The most common solution to avoid forgetting is experience replay—reintroducing data in previously supported languages during CPT (Winata et al., 2023; Wang et al., 2024b). Unfortunately, most recent model releases are not accompanied by their pretraining data (Touvron et al., 2023; Jiang et al., 2023; Gemma, 2025; Bai et al., 2023). Even if the data were available or approximated using public sources, as the number of supported languages in an LM grows, replaying data in all of them can also become computationally infeasible. Recent works proposing alternative approaches to mitigate forgetting have also been shown to work well only in conjunction with replay (Winata et al., 2023; Alexandrov et al., 2024; Chen et al., 2023; Aggarwal et al., 2024).

We propose a lightweight replay-free continued pretraining method called LAYRA (Layer-selective LoRA; see Figure 1 and §2). We add Low-Rank Adapter modules (Biderman et al., 2024a) to selected transformer layers in an LM during training while keeping other layers frozen. Our method takes inspiration from two recent works. Biderman et al. (2024a) show that LoRA based training can reduce forgetting but also reduce learning. To improve learning, we apply LoRA only to subset of the model layers inspired by Zhao

et al. (2024); Wendler et al. (2024) which showed evidence that multilingual LMs process an input sequence in three stages. Using logit lens based analysis (Nostalgebraist, 2020), they demonstrated that the earliest layers of LMs process the sequence in the language in which it is written, the middle layers process the sequence in English (the most dominant language in the pretraining corpora) and the final layers translate back and generate a response in the input language. By performing targeted updates to only the layers responsible for handling non-English text, we show that we improve learning while further reducing forgetting. We further show that by combining LAYRA with model merging methods, we can sequentially continue to add new languages to an LM. Finally, beyond adapting to new languages, we show that we can enable instruction following in the adapted models using instruction residuals extracted from already instruction-tuned models.

We validate our method by adding different combinations of three typologically different languages with limited pretraining resources (Galician, Urdu, Swahili) to Llama 3.1 (§3). We choose these languages to understand the impact of writing script and relatedness of target languages with the original model on both learning and forgetting. Our results (§4) show that LAYRA surpasses baselines in new language acquisition while outperforming or matching them in preventing forgetting. Our analysis reveals that a target language irrespective of its relatedness to the originally supported languages can be adapted successfully as long as its writing script is represented by the model.

2 The Continual Learning Problem

2.1 Problem Setup

Suppose we have a pretrained autoregressive LM θ_N that supports N languages (where $N > 1$). Given pretraining data in n new languages, $\{L_1, L_2, \dots, L_n\}$, our goal is to create a model θ_{N+n} that supports all $N + n$ languages. Crucially, θ_{N+n} should retain its performance in the original N languages (stability) while acquiring competence in the new n languages (plasticity). We also consider a generalized continual learning setup where given θ_{N+n} , we update it to include n' more languages, thus creating $\theta_{N+n+n'}$. In principle, this process can go on indefinitely as resources for new languages emerge reflecting common practice in language modeling and machine learning, where new training data arrive incrementally (BLOOM (Leong et al., 2022), Wura (Oladipo et al., 2023), or FineWeb 2 (Penedo et al., 2024)).

Furthermore, we assume no access to the pretraining data in the original N languages that led to the creation of θ_N . This also reflects a new reality in open-weights release of language models where the pretraining data or its constitution is often not publicly shared by the organizations building them (Dubey et al., 2024; Bai et al., 2023; Abdin et al., 2024).

Supporting instruction following in new languages To create models that can respond to user queries, pretrained LMs typically go through an instruction tuning phase using curated labeled datasets (Ouyang et al., 2022; Chung et al., 2024; Lambert et al., 2024). However, instruction-tuning datasets remain scarce or unavailable for most non-English languages. Hence, we explore data-free methods to add instruction following abilities to updated models θ_{N+n} , assuming access to an instruction tuned model that supports the original N languages, θ_N^{it} .

2.2 Method

We summarize our methodology in Figure 1. Our goal is to add new languages to θ_N without causing the model to forget the previously learned languages. The simplest and most naive approach to do this is to continue pretraining (CPT) θ_N with new data using a language modeling objective (such as next token prediction). However, it has been widely observed (and confirmed in our experiments) that CPT leads to severe catastrophic forgetting, causing large drops in performance on previously supported languages (Csaki et al., 2023). Experience replay, reintroducing past data during CPT, has been proposed as a viable approach to address this issue but it requires access to pretraining data of θ_N which is not available to us. While multilingual pretraining datasets in many languages

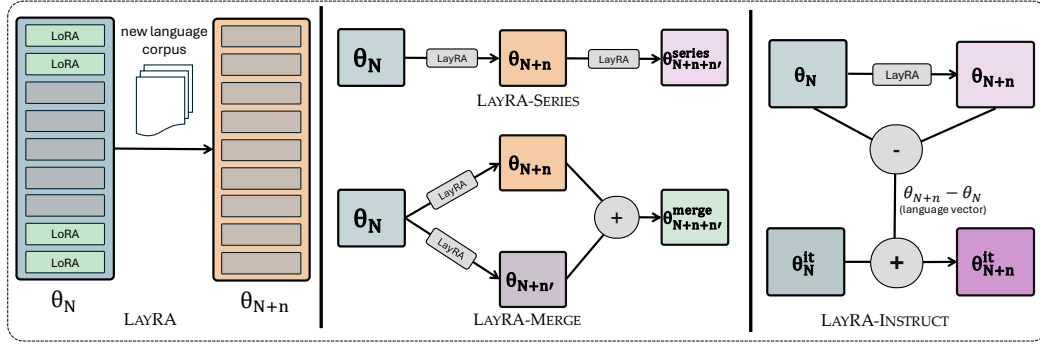


Figure 1: Problem setup for continually adding new languages and instruction tuning to a model. Left: Layer-selective LoRA which performs LoRA updates to only selected initial and final transformer layers. Middle: Sequential continual learning techniques to enable multiple stages of adaptation. Right: Enable instruction following in the adapted model without instruction data using model arithmetic.

95 have recently been open-sourced, recent state-of-the-art multilingual LLMs are trained with
 96 such a large number of languages that it can become extremely computationally expensive
 97 to perform replay. For example, BloomZ (Muennighoff et al., 2022), Aya 23 (Aryabumi
 98 et al., 2024), and Gemma 3 (Gemma, 2025) were trained on 23, 46 and over 140 languages
 99 respectively. This number is likely to grow over time.

100 Instead of experience replay, in this work, we explore continued pretraining with parameter
 101 efficient approaches, such as LoRA (Hu et al., 2022), which have been shown in prior work
 102 to “learn less and forget less” (Biderman et al., 2024a), thus finding a better balance between
 103 learning new tasks and retaining past knowledge. Our experiments on adding languages
 104 also reveal similar findings. To further minimize forgetting and improve learning, we
 105 propose the following improvements to the training procedure updating only a subset of
 106 the layers during CPT.

107 **Layer-selective Continual Learning** Recent work has indicated that multilingual LMs
 108 with layers $\{\mathcal{T}_l\}_{l=1}^L$ process sequences in three main stages (Zhao et al., 2024; Wendler et al.,
 109 2024).

- 110 1. The **earliest layers** ($\mathcal{T}_{\sim 1}$) encode the input in its source language.
- 111 2. The **middle layers** ($\mathcal{T}_{\sim L/2}$) handle the model’s internal “reasoning language” (often
 112 English in models such as Llama series of models).
- 113 3. The **final layers** ($\mathcal{T}_{\sim L}$) convert the representation back into the target language
 114 during generation.

115 Based on this observation, we hypothesize that training only the layers responsible for
 116 handling non-English text and freezing the English-specific layers should be sufficient to
 117 support new languages to the model, preserving its core reasoning abilities. We combine
 118 this layer-selective training method with parameter-efficient updates to propose our final
 119 training approach, which we call **Layer-selective LoRA (LAYRA)** which finds the best
 120 balance between learning and forgetting in our extensive experiments.¹

121 Given n new languages, we first continue pretraining θ_N to obtain an adapted model ϕ_{N+n} .
 122 We then obtain θ_{N+n} as,

$$\theta_{N+n} = \theta_N + \lambda(\phi_{N+n} - \theta_N) \quad (1)$$

¹A variation of this method as evaluated by Remy et al. (2024) showed its validity in faster adaptation to low-resource languages. However, their method still suffered from catastrophic forgetting which we aim to address in this work.

Here λ is a hyperparameter that controls the weights of added language vector that is learned during training. Similar to its usage in previous works (Morrison et al., 2024; Wang et al., 2024a), the addition here is a vector operation and this helps to obtain the right balance of language vectors to add when working with multiple language vectors. The language vector in this equation could also be swapped for a task vector as later seen in Equation 4

Sequentially adding new languages over multiple stages To support multiple stages of continual learning, where n languages are added in the first stage, and n' languages are added in the second stage (and so on), we explore the following two methods.

In the first setup, which we call LAYRA-SERIES, we iteratively apply LAYRA by first training θ_N to create θ_{N+n} . We then continue training θ_{N+n} on n' newer languages again following LAYRA creating

$$\theta_{N+n+n'}^{\text{series}} = \theta_{N+n} + \lambda'(\phi_{N+n+n'} - \theta_{N+n'}). \quad (2)$$

Here, λ' is another hyperparameter tuned separately from λ . This process can be continued indefinitely to add more languages.

In the second setup, which we call LAYRA-MERGE, we first create $\theta_{N+n'}$ separately without relying on θ_{N+n} by applying LAYRA with the n' languages on θ_N . We then merge the weights of these specialized models, θ_{N+n} and $\theta_{N+n'}$, to yield a single model $\theta_{N+n+n'}$. Concretely,

$$\theta_{N+n+n'}^{\text{merge}} = \mu\theta_{N+n} + (1 - \mu)\theta_{N+n'} \quad (3)$$

This approach aims to combine gains for each set of training without introducing additional forgetting issues from a previous LAYRA stage as in series. In practice, a practitioner may use different combinations of SERIES and MERGE depending on the languages being added to the underlying model.

Adding Instruction Following Capabilities So far, we describe an approach to add new languages to a pretrained model using raw text available in target languages. Recent open-weights models (Dubey et al., 2024; Bai et al., 2023; Abdin et al., 2024; OLMo et al., 2024) all follow a pattern of releasing both a pretrained (base) and an instruction-tuned model (θ_N^{it}). To add instruction following abilities in the adapted model θ_{N+n} without any labeled data in the n new languages, we compute a *language vector* as the difference between θ_{N+n} and θ_N and apply it to θ_N^{it} as,

$$\theta_{N+n}^{\text{it}} = \gamma(\theta_{N+n} - \theta_N) + \theta_N^{\text{it}} \quad (4)$$

By doing so, we inherit the instruction-following capabilities learned by θ_N^{it} while supporting newly added languages to create LAYRA-INSTRUCT.² The scaling factor γ can be tuned to balance instruction performance and new-language retention. Given instruction-tuning datasets in the target language, this model can further be improved but we do not assume any such access in this work.

3 Experimental Setup

3.1 Languages, Datasets, and Model

We use Llama 3.1 8B (Grattafiori et al., 2024) for our experiments which supports $N = 8$ languages. Six of them use Latin script (English, German, French, Italian, Portuguese, and Spanish), while two use non-Latin scripts (Hindi in Devanagari and Thai in Thai script).

²While more sophisticated methods of model merging have recently been developed such as TIES (Yadav et al., 2023) and DARE (Yu et al., 2024), our initial experiments did not show improvements with them over task vectors.

161 **New languages** To test the impact of writing script and their relatedness to existing
 162 languages in the model, we experimented with adding the following languages:

- 163 • **Galician:** a mid-resource Romance language (Latin script) spoken in northwestern
 164 Spain. Given its similarities to Portuguese and Spanish (both of which exist in the
 165 original model), Galician is well suited for leveraging prior multilingual knowledge.
- 166 • **Swahili:** a low-resource Bantu language predominantly spoken in East Africa by
 167 roughly 100 million speakers. It is not related to any of the languages in Llama 3.1
 168 but is written in Latin script, which is well represented in the model.
- 169 • **Urdu:** a low-resource Indo-Aryan language (Perso-Arabic script). Although Urdu
 170 shares substantial linguistic commonalities with Hindi (which is supported by the
 171 base model), its script differs from Hindi.

172 This diverse selection of languages provides a robust test of how effectively the model can
 173 learn distinct scripts and linguistic structures. For each of the three languages, we create
 174 our pretraining datasets using FineWeb 2 (Penedo et al., 2024). We use all available data for
 175 Swahili which was $\sim 1.2\text{B}$ tokens. For the other two languages, we subsampled the corpus
 176 to contain the same number of tokens to control for the impact of dataset size.

177 We conduct two sets of experiments: (1) a single-stage continual learning setup with $n = 1$
 178 where we add only one of the three languages at a time to the base model, and, (2) a two-
 179 stage setup with $n = 1$ and $n' = 1$, where we first add one of the languages to the pretrained
 180 model, and incorporate a second language later on.

181 3.2 LAYRA Hyperparameters

182 For our single-stage experiments (Equation 1), we set $\lambda = 1$ which adds the entire language
 183 vector that is obtained after CPT. This is analogous to a straightforward LoRA CPT (with
 184 selected layers). In LAYRA-SERIES where we iteratively add more languages following
 185 Equation 2, we empirically determine $\lambda' = 0.5$ to perform the best.³ In our second setup for
 186 adding multiple languages via merging, LAYRA-MERGE, we set $\mu = 0.5$ which is analogous
 187 to averaging all the adapted models from CPT. To add instruction following abilities to the
 188 adapted model as in Equation 4, we use a value of $\gamma = 0.7$ which adds part of the language
 189 vector to Llama 3.1 Instruct. We determine the value of γ with a small scale experiment
 190 with varying values. For the the adapters, we use a rank (r) of 8 and α as 16. We use LoRA
 191 dropout of 0.05. We only use this setup as results from Biderman et al. (2024a) shows that
 192 this rank and alpha results to the least forgetting during continual pretraining of an LLM.
 193 For all LAYRA experiments, we apply LoRA to the earliest 10 and the final 2 layers. In
 194 addition, we also finetune the embedding layer and the LM head. We provide analysis and
 195 a blation studies that provide the reasoning for choosing these hyperparameters in §5. All
 196 other training hyperparameters can be found in the Appendix A.4 Table 19.

197 3.3 Baselines

198 We compare LAYRA with the following methods.

- 199 • **Full CPT** In this baseline, we continue pretraining all the base model parameters.
- 200 • **LoRA CPT** In this baseline, we continue pretraining the base model using low-rank
 201 adapters (LoRA) applied at all layers following (Biderman et al., 2024a).
- 202 • **Layer-Selective Full CPT** In this baseline, we fully train the first and the last
 203 transformer layer of the base model along with the embedding layer and the LM
 204 head (these layers were empirically determined to give the best performance). This
 205 baseline also serves as an ablation of LAYRA with the adapters removed.

³While we perform experiments with only two stages, future stages may require an even smaller multiplier

3.4 Evaluation

We evaluate the adapted pretrained models on XNLI (Natural Language Inference; [Conneau et al., 2018](#)), PAWS-X (Paraphrasing; [Yang et al., 2019](#)), XCOPA (Commonsense Reasoning; [Ponti et al., 2020](#)), and XStoryCloze (Commonsense Reasoning; [Lin et al., 2021](#)). We evaluate the instruction adapted models on XNLI, MGSM (Math; [Shi et al., 2022](#)), and MMLU-Lite (MCQs; [Singh et al., 2024](#)). For MGSM and MMLU, we generate the answer by greedily decoding from the model with a temperature of 0.4 which we determine with earlier experiments we do not report. We evaluate the rest as classification tasks by choosing the labels with the the highest probability. We use the LM harness evaluation framework ([Biderman et al., 2024b](#)) our evaluations. We use 0-shot evaluation for all tasks except for MGSM for which we use a 3-shot setup following prior work ([Group et al., 2024](#)). Not all languages with which we experiment have datasets available for all tasks. For languages for which datasets are not available, we translate the English subset of the task to the missing language using Google Machine Translate (see Table 20 in the Appendix A.4 for language that required translations). We perform qualitative analysis to ensure that the translations are accurate. For all tasks, we report accuracy. For each task, we track *retention* measured by a minimal drop in performance in the originally supported languages and *gain* which is measured by improvement in performance in the newly added language(s). An ideal solution leads to the highest gains while maximizing retention.

4 Results

4.1 Adding One Language to the Pretrained Model

We provide results for one-stage continual learning by adding one language at time in Table 1. As expected, full CPT consistently exhibits the highest level of catastrophic forgetting across all languages and tasks—regardless of script, resource availability, or linguistic similarity of the target language to previously supported languages. Layer-Selective full CPT, by freezing most of the model layers and finetuning only the top and bottom layers improves the learning-forgetting tradeoff. However, significant forgetting still persists. LoRA CPT with its lightweight parameter updates closes the gap even further. Our proposed approach LayRA, considerably outperforms LoRA in terms of forgetting while being competitive and oftentimes exceeding LoRA in terms of acquisition across all tasks and target languages. We provide detailed results across multiple tasks and additional languages in Tables 3 4 5 and 6, which are included in the Appendix (A.1). In these results, we observe a 2 point margin over the baselines with our Urdu-trained model on over 5 languages for XNLI. We also see a similar trend with 2 languages with our Galician-trained model. With LAYRA, we see the least amount of forgetting in English across all tasks, which is the anchor language Llama 3.1 8B. This aligns with our hypothesis that freezing the model’s middle transformer layers preserves the core capabilities encoded in the anchor language.

Impact of language relatedness and scripts Using LAYRA, the Galician-trained model exhibits the strongest retention-gain tradeoff. We hypothesize that this result is due to positive transfer from Spanish and Portuguese (which are both supported by the original model). In fact, this model improves English performance across multiple tasks, highlighting that cross-lingual transfer can happen in both directions with our proposed approach. Swahili, while unrelated to any of the languages in the original model, also responds well to our training strategy with a substantial performance gain with a good amount of retention (albeit slightly worse than Galician). We speculate that this result is due to Latin script being well supported in the original model as well as the heavy usage of English loanwords in Swahili ([Martin et al., 2021](#)). On the other hand, models trained with Urdu, which is very closely related to Hindi, achieve the poorest overall performance, both in terms of gain as well as retention. This is due to Urdu’s writing script not being well represented in the base model’s tokenizer, leading to overfragmentation and poor adaptation. Although prior work has sought to address such issues with vocabulary expansion techniques ([Kim et al., 2024](#)), the resulting change in the number of model parameters hinders the use of model merging or parameter-efficient techniques to reduce catastrophic forgetting. Tokenizer-free models

XNLI evaluation for models trained on Swa/Urd/Glg							
Model	eng	spa	hin	swa	urd	glg	Avg
Pretrained	54.90	51.33	48.96	39.24	36.43	47.57	47.55
Full	52.93/49.88/50.88	44.42/34.86/47.31	34.18/35.82/33.90	45.46/34.26/32.61	33.45/39.68/37.43	37.19/34.26/50.93	40.95/37.57/41.21
layer-sel.	55.06/54.70/53.01	45.50/36.27/47.11	41.93/39.08/37.83	46.71/35.06/36.14	33.45/42.49/34.50	41.94/37.09/53.82	43.91/39.34/43.40
LoRA	55.90/55.78/54.34	48.84/42.89/49.32	45.30/39.12/45.70	47.71/36.10/35.90	36.55/42.65/36.75	44.08/48.55/54.02	46.82/44.82/46.49
LAYRA	54.22/57.11/55.81	47.43/43.82/50.64	46.83/42.01/48.07	45.34/36.83/34.66	34.98/40.96/37.55	45.48/47.81/54.02	46.64/45.58/47.14

PAWS-X evaluation for models trained on Swa/Urd/Glg							
Model	eng	spa	hin	swa	urd	glg	Avg
Pretrained	67.45	65.30	64.45	61.00	54.25	63.55	63.17
Full	65.25/66.60/61.70	61.50/54.90/62.65	61.50/52.95/54.15	63.00/52.95/53.35	55.05/46.85/47.50	50.80/47.15/62.85	59.36/54.11/57.37
layer-sel.	64.65/66.95/58.90	60.60/58.40/63.95	59.00/54.35/56.05	60.60/49.65/51.00	55.05/52.20/49.10	55.70/53.60/67.80	59.82/56.10/58.57
LoRA	68.75/70.45/67.15	64.00/61.10/64.60	63.35/64.30/62.85	61.70/48.00/46.15	59.90/49.35/50.45	56.70/59.10/66.60	62.97/59.58/60.89
LAYRA	68.95/68.45/67.15	63.85/60.35/64.30	64.20/63.70/63.30	61.20/53.50/50.35	59.75/54.85/49.45	59.65/60.80/65.25	63.25/60.79/61.03

XCOPA evaluation for models trained on Swa/Urd/Glg							
Model	eng	spa	tha	swa	urd	glg	Avg
Pretrained	87.00	81.40	57.60	55.00	58.80	57.60	67.14
Full	72.00/73.00/77.00	57.40/50.20/70.20	55.60/52.40/55.00	66.80/53.60/53.60	53.40/59.60/54.80	53.40/50.80/59.00	58.97/56.09/60.09
layer-sel.	83.00/77.00/80.00	60.20/56.40/76.20	57.60/54.20/56.60	66.00/54.00/53.40	53.20/57.60/57.20	54.00/53.40/58.00	61.00/58.37/62.34
LoRA	88.00/86.00/85.00	70.60/74.80/77.00	56.60/58.00/55.40	66.20/53.60/53.00	53.80/59.40/59.80	56.00/58.00/62.20	64.74/65.57/65.14
LAYRA	87.00/86.00/86.00	69.60/76.80/76.80	57.80/60.60/58.40	64.60/53.40/54.40	56.00/61.00/57.40	52.40/56.00/63.20	64.09/66.26/65.20

XStoryCloze evaluation for models trained on Swa/Urd/Glg							
Model	eng	spa	hin	swa	-	glg	Avg
Pretrained	78.16	70.75	64.46	55.86	-	64.46	66.74
Full	70.22/69.09	59.03/65.39	47.12/48.91	47.12/48.38	-	64.33/68.56	57.56/60.07
layer-sel.	76.57/75.84	66.05/69.49	55.46/52.61	64.99/49.90	-	54.47/70.42	63.51/63.65
LoRA	76.17/76.77	66.71/69.82	63.40/63.67	57.91/51.42	-	65.12/70.81	65.86/66.50
LAYRA	76.51/76.11	66.51/69.69	63.27/63.20	63.73/50.69	-	57.91/69.89	65.59/65.92

Table 1: Performance of different CPT methods across languages for XNLI, PAWS-X, XStoryCloze and XCOPA. See Tables 3, 4, 5 & 6 in the Appendix A.1 for full results with more languages which we use to compute the average.

may be a viable future direction in addressing this issue (Ahia et al., 2024). Furthermore, we observe that non-Latin-scripted languages such as Thai and Hindi also disproportionately suffer at retention across all methods highlighting a broader trend of negative transfer between languages that do not share scripts. We leave further examination of this trend to future work. Due to these observations, we exclude Urdu from further experiments and only report results with Galician and Swahili for two-stage continual learning.

4.2 Sequentially Adding Multiple Languages to the Pretrained Model

We provide the results for sequential addition of Galician and Swahili, LAYRA-SERIES and -MERGE in Table 2. Both of these methods assume that the resources of the languages arrived in order and not at the same time. For reference, we also include results for CPT assuming datasets for both languages were indeed available at the same time (referred to as PARALLEL). Unsurprisingly, PARALLEL produced the highest overall gain-retention trade-off, indicating the effectiveness of single-stage adaptation with multiple languages ($n = 2$). This method serves as the upper bound for the multi-stage learning approaches. With LAYRA-SERIES, the gains tend to shift toward the most recently added language with a slight forgetting of previously acquired languages. Adding related languages at the second stage (Swahili then Galician) leads to better retention. In comparison, LAYRA-MERGE performs much better matching or even surpassing PARALLEL, yielding the highest retention of knowledge for languages employing Latin scripts.

4.3 Adding Instruction Tuning to the Adapted Model

In Figure 2, we observe that adding an instruction residual (as described in Equation 4) can enable our adapted models to follow user instructions. XNLI shows clear trends of improvement of the base instruct model across all three languages in our experiments. With MGSM and MMLU the trends are not consistent. For MGSM, we observe an increase in accuracy for Swahili and Urdu but a decline in performance for Galician. For MMLU both Urdu and Galician show declines. While we do not identify clear reasons for this performance drop, we attribute it to tokenization issues with Urdu and previously identified issues with simple

Tasks	XNLI / PAWS-X						
Model	deu	eng	spa	fra	swa	glg	Avg
Pretrained	52.05/66.20	54.90/67.45	51.33/65.30	50.12/64.45	39.24/61.00	47.57/63.55	48.95/64.66
Parallel	50.00/66.25	53.53/66.05	50.32/64.65	45.94/62.25	43.73/58.05	56.08/67.35	49.11/64.10
Series (Glg→Swa)	48.88/65.10	54.74/66.80	50.96/62.90	48.63/63.50	43.86/53.75	52.01/55.65	48.98/61.28
Series (Swa→Glg)	51.81/64.10	56.06/67.35	51.24/64.55	48.11/65.45	42.97/47.50	54.82/64.20	49.51/62.19
Merging	51.00/67.30	54.50/67.45	51.81/65.30	49.92/64.45	43.57/61.00	53.98/63.55	49.00/64.66

Tasks	XStoryCloze / XCOPA						
Model	eng	spa	hin/tha	-	swa	glg	Avg
Pretrained	78.16/87.00	70.75/81.40	64.46/57.60	-	55.86/55.00	64.46/57.60	66.74/68.53
Parallel	76.64/86.00	69.16/76.80	63.47/57.00	-	62.14/63.40	69.95/60.80	68.27/67.87
Series (Glg→Swa)	75.91/81.00	68.23/74.40	64.13/57.40	-	58.44/59.40	67.31/61.80	66.80/66.37
Series (Swa→Glg)	75.45/83.00	68.56/76.40	63.53/57.60	-	57.45/60.00	67.90/59.20	66.58/66.37
Merging	76.64/87.00	69.03/78.80	64.26/58.00	-	56.45/58.20	66.51/60.20	66.58/68.10

Table 2: Performance of different LAYRA setups for adding two languages (Galician + Swahili) across XNLI, PAWS-X, XStoryCloze and XCOPA. Top for XNLI & PAWS-X. Bottom for XStoryCloze & XCOPA. See Table 7, 9, 8 & 10 for full results with more languages.

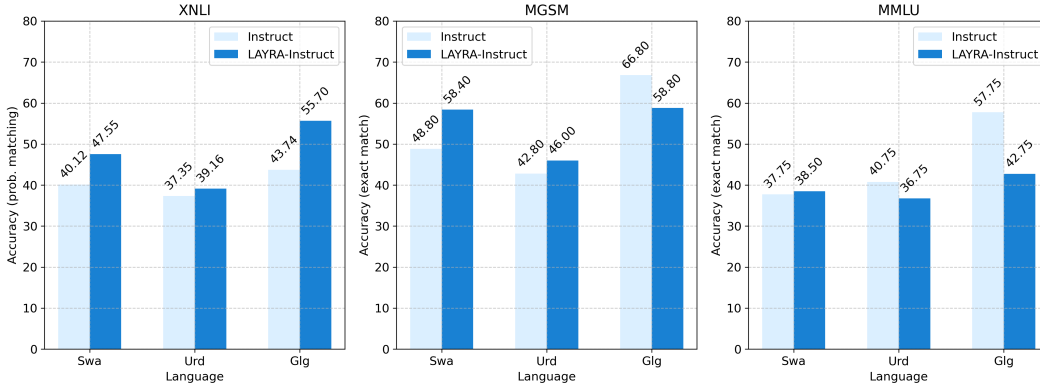


Figure 2: Accuracy of the instruction base model vs the adapted model on XNLI, MGSM and MMLU

model arithmetic techniques [Yadav et al. \(2023\)](#); [Tao et al. \(2024\)](#). Given a small amount of instruction tuning data in the target languages, this gap may be filled. We also measure how much our LAYRA-INSTRUCT models forget its previous knowledge by evaluating them on the English version of our mentioned tasks in [Figure 3](#) in Appendix A.4. We find that our adapted models still achieve high accuracy in English and at times outperforms Llama 3.1 Instruct.

5 Ablations

Varying the number of layers for LAYRA. To choose the optimal number of layers that balances the learning-forgetting tradeoff during CPT, we ranged the number of early and final transformer layers to be finetuned from 1 to 10 each. We conducted this evaluation with Swahili and find that the combination with first 10 and last 2 layers gives us the best balance. See Table 11, 12, 13, and 14 in Appendix A.2 for details.

Language vector Scaling in LAYRA-SERIES. We investigate the impact of changing the language vector added during the sequential addition of multiple languages. We continually increase λ' (from [Equation 1](#)) for LAYRA-SERIES from 0 to 1 in our Galician and Swahili series experiment (Glg→Swa). For all the task we evaluated on, (see Tables 15, 16, 17, 18 in Appendix A.3), as λ' tends to 1, we observe more retention of the Swahili and more

303 forgetting of Galician while there is a general drop in accuracy for all the previously learned
 304 languages.

305 6 Related Work

306 **Language Adaptation** There exists extensive prior research to adapt LMs to new languages
 307 (Ogueji et al., 2021; Alabi et al., 2022; Lu et al., 2024). Most studies have focused on continued
 308 pretraining of all parameters of the models (Csaki et al., 2023; Alabi et al., 2022), adding
 309 new parameters such as adapters (Yong et al., 2022), or training a small subset of the model
 310 parameters (Pfeiffer et al., 2020; Houlsby et al., 2019; Remy et al., 2024). Similarly, our work
 311 uses adapters (LoRA) and applies them on a subset of model layers. These approaches are
 312 motivated by training LMs in the target languages(s), not preserving the performance in
 313 the original ones. They benefit from cross-lingual transfer of encoded knowledge in the
 314 pretrained models. If the script of the target knowledge is not supported by the pretrained
 315 models’ tokenizer, Han et al. (2024) show that adapting can be challenging. We demonstrate
 316 a similar issue with adapting Llama 3.1 to Urdu. A commonly proposed solution to address
 317 this issue is the expanding of vocabulary before continuing pretraining (Liu et al., 2023;
 318 Dobler & De Melo, 2023; Mundra et al., 2024). However, these techniques are known to
 319 exacerbate the forgetting issue (Mundra et al., 2024); model merging techniques to mitigate
 320 the issue cannot be applied due to different model sizes.

321 **Mitigating Catastrophic Forgetting** Catastrophic forgetting is a well known issue in neural
 322 models and remains a challenge for modern LMs even for other cases beyond language
 323 adaptation. Reintroducing original data during adaptation (known as experience replay) is
 324 a commonly adopted remedy (Rolnick et al., 2019; Csaki et al., 2023; Winata et al., 2023). We
 325 explore strategies that do not assume access to the original data which is a new reality in
 326 the case of modern LMs.

327 Specifically, we modify LoRA by restricting it to select layers to improve this tradeoff. We
 328 leave exploration of learning rate schedules along with LAYRA for future work.

329 **Model Arithmetic** As a way to combine multiple models without training, model merging
 330 has been widely explored in the context of modern LMs (Hammoud et al., 2024; Dziadzio
 331 et al., 2024; Yang et al., 2024). Since its inception, many advanced merging techniques have
 332 been explored in recent works (Yadav et al., 2023; Yu et al., 2024; Kim et al., 2023). In our
 333 early exploration, they did not outperform the simplest arithmetic technique for creating
 334 task vectors proposed in Ilharco et al. (2022). Hence, we adopt it for sequential adaption
 335 and for creating our instruction adapted model. Multiple works have also explored model
 336 arithmetic during continued pretraining or finetuning showing it can match or improve
 337 the performance of training from scratch. Most related to our work is BAM (Alexandrov
 338 et al., 2024) which perform full finetuning and merge after every few iterations but they
 339 used experience replay which is not application to our setup.

340 7 Conclusion

341 We introduced LAYRA, a layer-selective adapter-based method to continuously add new
 342 languages to a multilingual LLM. By strategically updating only the first and the last
 343 few transformer layers, LAYRA effectively preserves knowledge of previously supported
 344 languages while learning new ones. Our experiments demonstrate that this targeted, low-
 345 rank adaptation approach not only mitigates catastrophic forgetting but also benefits from
 346 cross-lingual transfer, and can improve performance on existing languages. In addition, our
 347 merging strategies enable sequential continual learning, maintaining a favorable balance
 348 between stability and plasticity. Finally, we showed the potential of LAYRA to integrate
 349 instruction-following capabilities, even in scenarios where instruction-tuning data for newly
 350 added languages is not available. We tested our approach only with an 8B model and
 351 low-resource languages. We leave the exploration of model size and data for future work.

References

- Marah Abdin, Jyoti Aneja, Hany Awadalla, Ahmed Awadallah, Ammar Ahmad Awan, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Jianmin Bao, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Divyanshu Aggarwal, Ashutosh Sathe, and Sunayana Sitaram. Exploring pretraining via active forgetting for improving cross lingual transfer for decoder language models. *arXiv preprint arXiv:2410.16168*, 2024.
- Orevaoghene Ahia, Sachin Kumar, Hila Gonen, Valentin Hofmann, Tomasz Limisiewicz, Yulia Tsvetkov, and Noah A. Smith. MAGNET: Improving the multilingual fairness of language models with adaptive gradient-based tokenization. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, 2024. URL <https://openreview.net/forum?id=1e3MOWHSIX>.
- Jesujoba O. Alabi, David Ifeoluwa Adelani, Marius Mosbach, and Dietrich Klakow. Adapting pre-trained language models to African languages via multilingual adaptive fine-tuning. In *Proceedings of the 29th International Conference on Computational Linguistics*, pp. 4336–4349, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics. URL <https://aclanthology.org/2022.coling-1.382>.
- Anton Alexandrov, Veselin Raychev, Mark Niklas Mueller, Ce Zhang, Martin Vechev, and Kristina Toutanova. Mitigating catastrophic forgetting in language transfer via model merging. In Yaser Al-Onaizan, Mohit Bansal, and Yun-Nung Chen (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 17167–17186, Miami, Florida, USA, November 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-emnlp.1000. URL <https://aclanthology.org/2024.findings-emnlp.1000/>.
- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, et al. Aya 23: Open weight releases to further multilingual progress. *arXiv preprint arXiv:2405.15032*, 2024.
- Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, et al. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- Dan Biderman, Jacob Portes, Jose Javier Gonzalez Ortiz, Mansheej Paul, Philip Greengard, Connor Jennings, Daniel King, Sam Havens, Vitaliy Chiley, Jonathan Frankle, et al. Lora learns less and forgets less. *arXiv preprint arXiv:2405.09673*, 2024a.
- Stella Biderman, Hailey Schoelkopf, Lintang Sutawika, Leo Gao, Jonathan Tow, Baber Abbasi, Alham Fikri Aji, Pawan Sasanka Ammanamanchi, Sidney Black, Jordan Clive, et al. Lessons from the trenches on reproducible evaluation of language models. *arXiv preprint arXiv:2405.14782*, 2024b.
- Samuel Cahyawijaya, Holy Lovenia, Tiezheng Yu, Willy Chung, and Pascale Fung. Instruct-align: teaching novel languages with to llms through alignment-based cross-lingual instruction. *arXiv preprint arXiv:2305.13627*, 2023.
- Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. Multieurlex—a multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer. *arXiv preprint arXiv:2109.00904*, 2021.
- Yihong Chen, Kelly Marchisio, Roberta Raileanu, David Adelani, Pontus Lars Erik Saito Stenetorp, Sebastian Riedel, and Mikel Artetxe. Improving language plasticity via pretraining with active forgetting. *Advances in Neural Information Processing Systems*, 36:31543–31557, 2023.
- Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Yunxuan Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *Journal of Machine Learning Research*, 25(70):1–53, 2024.

- 402 Alexis Conneau, Guillaume Lample, Ruty Rinott, Adina Williams, Samuel R Bowman,
403 Holger Schwenk, and Veselin Stoyanov. Xnli: Evaluating cross-lingual sentence represen-
404 tations. *arXiv preprint arXiv:1809.05053*, 2018.
- 405 Zoltan Csaki, Pian Pawakapan, Urmish Thakker, and Qiantong Xu. Efficiently adapting
406 pretrained language models to new languages. *arXiv preprint arXiv:2311.05741*, 2023.
- 407 Konstantin Dobler and Gerard De Melo. Focus: Effective embedding initialization for
408 monolingual specialization of multilingual models. *arXiv preprint arXiv:2305.14481*, 2023.
- 409 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle,
410 Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3
411 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 412 Sebastian Dziadzio, Vishaal Udandaraao, Karsten Roth, Ameya Prabhu, Zeynep Akata,
413 Samuel Albanie, and Matthias Bethge. How to merge your multimodal models over time?
414 *arXiv preprint arXiv:2412.06712*, 2024.
- 415 Manuel Faysse, Patrick Fernandes, Nuno Guerreiro, António Loison, Duarte Alves, Caio
416 Corro, Nicolas Boizard, Jaee Alves, Ricardo Rei, Pedro Raphaël Martins, et al. Crois-
417 santllm: A truly bilingual french-english language model. 2024.
- 418 Team Gemma. Gemma 3 technical report. Google, 2025. URL <https://storage.googleapis.com/deepmind-media/gemma/Gemma3Report.pdf>.
- 420 Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
421 Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The
422 llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 423 Rakuten Group, Aaron Levine, Connie Huang, Chenguang Wang, Eduardo Batista, Ewa
424 Szymanska, Hongyi Ding, Hou Wei Chou, Jean-François Pessiot, Johannes Effendi,
425 et al. Rakutenai-7b: Extending large language models for japanese. *arXiv preprint*
426 *arXiv:2403.15484*, 2024.
- 427 Hasan Abed Al Kader Hammoud, Umberto Michieli, Fabio Pizzati, Philip Torr, Adel Bibi,
428 Bernard Ghanem, and Mete Ozay. Model merging and safety alignment: One bad model
429 spoils the bunch. *arXiv preprint arXiv:2406.14563*, 2024.
- 430 Hyojung Han, Akiko Eriguchi, Haoran Xu, Hieu Hoang, Marine Carpuat, and Huda
431 Khayrallah. Adapters for altering llm vocabularies: What languages benefit the most?
432 *arXiv preprint arXiv:2410.09644*, 2024.
- 433 Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Larous-
434 silhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer
435 learning for nlp. In *International conference on machine learning*, pp. 2790–2799. PMLR, 2019.
- 436 Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang,
437 Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*,
438 1(2):3, 2022.
- 439 Gabriel Ilharco, Marco Tulio Ribeiro, Mitchell Wortsman, Suchin Gururangan, Ludwig
440 Schmidt, Hannaneh Hajishirzi, and Ali Farhadi. Editing models with task arithmetic.
441 *arXiv preprint arXiv:2212.04089*, 2022.
- 442 Albert Q Jiang, A Sablayrolles, A Mensch, C Bamford, D Singh Chaplot, Ddl Casas,
443 F Bressand, G Lengyel, G Lample, L Saulnier, et al. Mistral 7b. arxiv. *arXiv preprint*
444 *arXiv:2310.06825*, 10, 2023.
- 445 Dahyun Kim, Chanjun Park, Sanghoon Kim, Wonsung Lee, Wonho Song, Yunsu Kim,
446 Hyeonwoo Kim, Yungi Kim, Hyeonju Lee, Jihoo Kim, et al. Solar 10.7 b: Scaling large lan-
447 guage models with simple yet effective depth up-scaling. *arXiv preprint arXiv:2312.15166*,
448 2023.

- Seungduk Kim, Seungtaek Choi, and Myeongho Jeong. Efficient and effective vocabulary expansion towards multilingual large language models. *arXiv preprint arXiv:2402.14714*, 2024.
- Nathan Lambert, Jacob Morrison, Valentina Pyatkin, Shengyi Huang, Hamish Ivison, Faeze Brahman, Lester James V Miranda, Alisa Liu, Nouha Dziri, Shane Lyu, et al. Tulu 3: Pushing frontiers in open language model post-training. *arXiv preprint arXiv:2411.15124*, 2024.
- Colin Leong, Joshua Nemecek, Jacob Mansdorfer, Anna Filighera, Abraham Owodunni, and Daniel Whitenack. Bloom library: Multimodal datasets in 300+ languages for a variety of downstream tasks. *arXiv preprint arXiv:2210.14712*, 2022.
- Xi Victoria Lin, Todor Mihaylov, Mikel Artetxe, Tianlu Wang, Shuohui Chen, Daniel Simig, Myle Ott, Naman Goyal, Shruti Bhosale, Jingfei Du, et al. Few-shot learning with multilingual language models. *arXiv preprint arXiv:2112.10668*, 2021.
- Yihong Liu, Peiqin Lin, Mingyang Wang, and Hinrich Schütze. Ofa: A framework of initializing unseen subword embeddings for efficient large-scale multilingual continued pretraining. *arXiv preprint arXiv:2311.08849*, 2023.
- Wei Lu, Rachel K Luu, and Markus J Buehler. Fine-tuning large language models for domain adaptation: Exploration of training strategies, scaling, model merging and synergistic capabilities. *arXiv preprint arXiv:2409.03444*, 2024.
- Gati L Martin, Medard E Mswahili, and Young-Seob Jeong. Sentiment classification in swahili language using multilingual bert. *arXiv preprint arXiv:2104.09006*, 2021.
- Jacob Morrison, Noah A Smith, Hannaneh Hajishirzi, Pang Wei Koh, Jesse Dodge, and Pradeep Dasigi. Merge to learn: Efficiently adding skills to language models with model merging. *arXiv preprint arXiv:2410.12937*, 2024.
- Niklas Muennighoff, Thomas Wang, Lintang Sutawika, Adam Roberts, Stella Biderman, Teven Le Scao, M Saiful Bari, Sheng Shen, Zheng-Xin Yong, Hailey Schoelkopf, et al. Crosslingual generalization through multitask finetuning. *arXiv preprint arXiv:2211.01786*, 2022.
- Nandini Mundra, Aditya Nanda Kishore, Raj Dabre, Ratish Puduppully, Anoop Kunchukuttan, and Mitesh M Khapra. An empirical comparison of vocabulary expansion and initialization approaches for language models. *arXiv preprint arXiv:2407.05841*, 2024.
- Nostalgebraist. Interpreting gpt: the logit lens. less-wrong. *arXiv preprint arXiv:2104.09006*, 2020. URL <https://www.lesswrong.com/posts/AcKRB8wDpdaN6v6ru/interpreting-gpt-the-logit-lens>.
- Kelechi Ogueji, Yuxin Zhu, and Jimmy Lin. Small data? no problem! exploring the viability of pretrained multilingual language models for low-resourced languages. In Duygu Ataman, Alexandra Birch, Alexis Conneau, Orhan Firat, Sebastian Ruder, and Gozde Gul Sahin (eds.), *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pp. 116–126, Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.mrl-1.11. URL <https://aclanthology.org/2021.mrl-1.11/>.
- Akintunde Oladipo, Mofetoluwa Adeyemi, Orevaoghene Ahia, Abraham Owodunni, Odunayo Ogundepo, David Adelani, and Jimmy Lin. Better quality pre-training data and t5 models for african languages. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 158–168, 2023.
- Team OLMo, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, et al. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*, 2024.

- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744, 2022.
- Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Martin Jaggi, Leandro von Werra, and Thomas Wolf. Fineweb2: A sparkling update with 1000s of languages, December 2024. URL <https://huggingface.co/datasets/HuggingFaceFW/fineweb-2>.
- Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052*, 2020.
- Edoardo Maria Ponti, Goran Glavaš, Olga Majewska, Qianchu Liu, Ivan Vulić, and Anna Korhonen. Xcopa: A multilingual dataset for causal commonsense reasoning. *arXiv preprint arXiv:2005.00333*, 2020.
- François Remy, Pieter Delobelle, Hayastan Avetisyan, Alfiya Khabibullina, Miryam de Lhoneux, and Thomas Demeester. Trans-tokenization and cross-lingual vocabulary transfers: Language adaptation of llms for low-resource nlp. *arXiv preprint arXiv:2408.04303*, 2024.
- David Rolnick, Arun Ahuja, Jonathan Schwarz, Timothy Lillicrap, and Gregory Wayne. Experience replay for continual learning. *Advances in neural information processing systems*, 32, 2019.
- Freda Shi, Mirac Suzgun, Markus Freitag, Xuezhi Wang, Suraj Srivats, Soroush Vosoughi, Hyung Won Chung, Yi Tay, Sebastian Ruder, Denny Zhou, Dipanjan Das, and Jason Wei. Language models are multilingual chain-of-thought reasoners, 2022.
- Shivalika Singh, Angelika Romanou, Clémentine Fourrier, David I Adelani, Jian Gang Ngui, Daniel Vila-Suero, Peerat Limkonchotiawat, Kelly Marchisio, Wei Qi Leong, Yosephine Susanto, et al. Global mmlu: Understanding and addressing cultural and linguistic biases in multilingual evaluation. *arXiv preprint arXiv:2412.03304*, 2024.
- Mingxu Tao, Chen Zhang, Quzhe Huang, Tianyao Ma, Songfang Huang, Dongyan Zhao, and Yansong Feng. Unlocking the potential of model merging for low-resource languages. *arXiv preprint arXiv:2407.03994*, 2024.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Tu Vu, Aditya Barua, Brian Lester, Daniel Cer, Mohit Iyyer, and Noah Constant. Overcoming catastrophic forgetting in zero-shot cross-lingual generation. *arXiv preprint arXiv:2205.12647*, 2022.
- Ke Wang, Nikolaos Dimitriadis, Alessandro Favero, Guillermo Ortiz-Jimenez, Francois Fleuret, and Pascal Frossard. Lines: Post-training layer scaling prevents forgetting and enhances model merging. *arXiv preprint arXiv:2410.17146*, 2024a.
- Yifan Wang, Yafei Liu, Chufan Shi, Haoling Li, Chen Chen, Haonan Lu, and Yujiu Yang. Insl: A data-efficient continual learning paradigm for fine-tuning large language models with instructions. *arXiv preprint arXiv:2403.11435*, 2024b.
- Chris Wendler, Veniamin Veselovsky, Giovanni Monea, and Robert West. Do llamas work in english? on the latent language of multilingual transformers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15366–15394, 2024.
- Genta Indra Winata, Lingjue Xie, Karthik Radhakrishnan, Shijie Wu, Xisen Jin, Pengxiang Cheng, Mayank Kulkarni, and Daniel Preotiuc-Pietro. Overcoming catastrophic forgetting in massively multilingual continual learning. *arXiv preprint arXiv:2305.16252*, 2023.

- 546 Prateek Yadav, Derek Tam, Leshem Choshen, Colin A Raffel, and Mohit Bansal. Ties-
 547 merging: Resolving interference when merging models. *Advances in Neural Information*
 548 *Processing Systems*, 36:7093–7115, 2023.
- 549 Enneng Yang, Li Shen, Guibing Guo, Xingwei Wang, Xiaochun Cao, Jie Zhang, and Dacheng
 550 Tao. Model merging in llms, mllms, and beyond: Methods, theories, applications and
 551 opportunities. *arXiv preprint arXiv:2408.07666*, 2024.
- 552 Yinfei Yang, Yuan Zhang, Chris Tar, and Jason Baldridge. Paws-x: A cross-lingual adversarial
 553 dataset for paraphrase identification. *arXiv preprint arXiv:1908.11828*, 2019.
- 554 Zheng-Xin Yong, Hailey Schoelkopf, Niklas Muennighoff, Alham Fikri Aji, David Ifeoluwa
 555 Adelani, Khalid Almubarak, M Saiful Bari, Lintang Sutawika, Jungo Kasai, Ahmed
 556 Baruwa, et al. Bloom+ 1: Adding language support to bloom for zero-shot prompting.
 557 *arXiv preprint arXiv:2212.09535*, 2022.
- 558 Le Yu, Bowen Yu, Haiyang Yu, Fei Huang, and Yongbin Li. Language models are su-
 559 per mario: Absorbing abilities from homologous models as a free lunch. In *Forty-first*
 560 *International Conference on Machine Learning*, 2024.
- 561 Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. How do
 562 large language models handle multilingualism? *arXiv preprint arXiv:2402.18815*, 2024.

563 A Appendix

564 A.1 Full results for LAYRA

565 This section provides comprehensive tables detailing the complete set of experimental
 566 results for LAYRA and baseline methods across all languages and evaluation tasks (XNLI: 3,
 567 PAWS-X: 4, XCOPA: 5, XStoryCloze: 6). These results offer further evidence supporting the
 568 conclusions drawn in the main text §4, allowing for deeper comparison and validation of
 performance metrics.

Model	deu	eng	spa	fra	hin	tha	swa	urd	glg	Avg
Pretrained	52.05	54.90	51.33	50.12	48.96	47.39	39.24	36.43	47.57	47.55
Swa Full	41.33	52.93	44.42	44.54	34.18	35.02	45.46	33.45	37.19	40.95
Swa layer-sel.	44.78	55.06	45.50	48.27	41.93	37.59	46.71	33.45	41.94	43.91
Swa LoRA	48.63	55.90	48.84	48.31	45.30	46.02	47.71	36.55	44.08	46.82
Swa LAYRA	49.92	54.22	47.43	49.24	46.83	46.35	45.34	34.98	45.48	46.64
Urd Full	37.11	49.88	34.86	35.86	35.82	36.39	34.26	39.68	34.26	37.57
Urd layer-sel.	37.91	54.70	36.27	36.87	39.08	34.58	35.06	42.49	37.09	39.34
Urd LoRA	46.18	55.78	42.89	49.32	39.12	42.77	36.10	42.65	48.55	44.82
Urd LAYRA	48.39	57.11	43.82	50.00	42.01	43.29	36.83	40.96	47.81	45.58
Glg Full	43.17	50.88	47.31	40.84	33.90	33.82	32.61	37.43	50.93	41.21
Glg layer-sel.	47.67	53.01	47.11	43.57	37.83	36.99	36.14	34.50	53.82	43.40
Glg LoRA	51.77	54.34	49.32	45.70	45.70	44.90	35.90	36.75	54.02	46.49
Glg LAYRA	50.84	55.81	50.64	46.22	48.07	46.43	34.66	37.55	54.02	47.14

Table 3: Performance of different CPT methods across languages for XNLI (%)

570 A.2 LAYRA layer selection ablation results

571 Here we present results for ablation studies done to find the best layer combinations for
 572 LAYRA. The tables in this section contains all the languages and task we evaluate our
 573 models on. The tables are also split up and expanded for easier comprehension.

Model	deu	eng	spa	fra	swa	urd	glg	Avg
Pretrained	66.20	67.45	65.30	64.45	61.00	54.25	63.55	63.17
Swa Full	58.45	65.25	61.50	61.50	63.00	55.05	50.80	59.36
Swa layer-sel.	63.15	64.65	60.60	59.00	60.60	55.05	55.70	59.82
Swa LoRA	66.40	68.75	64.00	63.35	61.70	59.90	56.70	62.97
Swa LAYRA	65.15	68.95	63.85	64.20	61.20	59.75	59.65	63.25
Urd Full	57.35	66.60	54.90	52.95	52.95	46.85	47.15	54.11
Urd layer-sel.	57.55	66.95	58.40	54.35	49.65	52.20	53.60	56.10
Urd LoRA	64.75	70.45	61.10	64.30	48.00	49.35	59.10	59.58
Urd LAYRA	63.90	68.45	60.35	63.70	53.50	54.85	60.80	60.79
Glg Full	59.40	61.70	62.65	54.15	53.35	47.50	62.85	57.37
Glg layer-sel.	63.20	58.90	63.95	56.05	51.00	49.10	67.80	58.57
Glg LoRA	68.40	67.15	64.60	62.85	46.15	50.45	66.60	60.89
Glg LAYRA	67.40	67.15	64.30	63.30	50.35	49.45	65.25	61.03

Table 4: Performance of different CPT methods across languages for PAWS-X

Model	eng	spa	hin	swa	glg	Avg
Pretrained	78.16	70.75	64.46	55.86	64.46	66.74
Swa Full	70.22	59.03	47.12	47.12	64.33	57.56
Swa layer-sel.	76.57	66.05	55.46	64.99	54.47	63.51
Swa LoRA	76.17	66.71	63.40	57.91	65.12	65.86
Swa LAYRA	76.51	66.51	63.27	63.73	57.91	65.59
Glg Full	69.09	65.39	48.91	48.38	68.56	60.07
Glg layer-sel.	75.84	69.49	52.61	49.90	70.42	63.65
Glg LoRA	76.77	69.82	63.67	51.42	70.81	66.50
Glg LAYRA	76.11	69.69	63.20	50.69	69.89	65.92

Table 5: Performance of different CPT methods across languages for XStoryClose (%)

574 A.3 Changing the language vector in LAYRA-SERIES.

575 This appendix shows the sensitivity of the model to the scaling factor λ' used when sequen-
576 tially adding new languages with LAYRA-SERIES. We include detailed tables and analyses
577 to illustrate how changing this hyperparameter influences retention of previously acquired
578 languages versus performance gains on newly added ones.

579 A.4 Other Tables and Figures

580 We put all other tables and figure in this section such as hyperparameter 19 table containing
581 exhaustive details regarding experimental setups, including training hyperparameters such
582 as learning rates, batch sizes, etc. We also add a table to show the languages we translate in
583 Table 20 and Figure to show the forgetting rate of LAYRA-INSTRUCT on English (Figure 3).
584

Model	eng	spa	ita	tha	swa	urd	glg	Avg
Pretrained	87.00	81.40	72.60	57.60	55.00	58.80	57.60	67.14
Swa Full	72.00	57.40	54.20	55.60	66.80	53.40	53.40	58.97
Swa layer-sel.	83.00	60.20	53.00	57.60	66.00	53.20	54.00	61.00
Swa LoRA	88.00	70.60	62.00	56.60	66.20	53.80	56.00	64.74
Swa LAYRA	87.00	69.60	61.20	57.80	64.60	56.00	52.40	64.09
Urd Full	73.00	50.20	53.00	52.40	53.60	59.60	50.80	56.09
Urd layer-sel.	77.00	56.40	56.00	54.20	54.00	57.60	53.40	58.37
Urd LoRA	86.00	74.80	69.20	58.00	53.60	59.40	58.00	65.57
Urd LAYRA	86.00	76.80	70.00	60.60	53.40	61.00	56.00	66.26
Glg Full	77.00	70.20	51.00	55.00	53.60	54.80	59.00	60.09
Glg layer-sel.	80.00	76.20	55.00	56.60	53.40	57.20	58.00	62.34
Glg LoRA	85.00	77.00	63.60	55.40	53.00	59.80	62.20	65.14
Glg LAYRA	86.00	76.80	60.20	58.40	54.40	57.40	63.20	65.20

Table 6: Performance of different CPT methods across languages for XCOPA

Model	deu	eng	spa	fra	hin	tha	swa	glg	Avg
Pretrained	52.05	54.90	51.33	50.12	48.96	47.39	39.24	47.57	48.95
Series (Glg→Swa)	48.88	54.74	50.96	48.63	47.95	44.82	43.86	52.01	48.98
Series (Swa→Glg)	51.81	56.06	51.24	48.11	46.87	44.18	42.97	54.82	49.51
Parallel	50.00	53.53	50.32	45.94	48.63	44.66	43.73	56.08	49.11
Merging	51.00	54.50	51.81	49.92	44.86	42.33	43.57	53.98	49.00

Table 7: Performance of different LAYRA setups for adding two languages (Galician + Swahili) on XNLI

Model	eng	spa	hin	swa	glg	Avg
Pretrained	78.16	70.75	64.46	55.86	64.46	66.74
Series (Glg→Swa)	75.91	68.23	64.13	58.44	67.31	66.80
Series (Swa→Glg)	75.45	68.56	63.53	57.45	67.90	66.58
Parallel	76.64	69.16	63.47	62.14	69.95	68.27
Merging	76.64	69.03	64.26	56.45	66.51	66.58

Table 8: Performance of different LAYRA setups for adding two languages (Galician + Swahili) on XStoryCloze

Model	deu	eng	spa	fra	swa	glg	Avg
Pretrained	66.20	67.45	65.30	64.45	61.00	63.55	64.66
Series (Glg→Swa)	65.10	66.80	62.90	63.50	53.75	55.65	61.28
Series (Swa→Glg)	64.10	67.35	64.55	65.45	47.50	64.20	62.19
Parallel	66.25	66.05	64.65	62.25	58.05	67.35	64.10
Merging	64.05	67.30	64.55	63.50	56.50	65.15	63.51

Table 9: Performance of different LAYRA setups for adding two languages (Galician + Swahili) on PAWS-X

Model	eng	spa	ita	tha	swa	glg	Avg
Pretrained	87.00	81.40	72.60	57.60	55.00	57.60	68.53
Series (Glg→Swa)	81.00	74.40	64.20	57.40	59.40	61.80	66.37
Series (Swa→Glg)	83.00	76.40	62.00	57.60	60.00	59.20	66.37
Parallel	86.00	76.80	63.20	57.00	63.40	60.80	67.87
Merging	87.00	78.80	66.40	58.00	58.20	60.20	68.10

Table 10: Performance of different LAYRA setups for adding two languages (Galician + Swahili) on XCOPA

Model	deu	eng	spa	fra	hin	tha	swa	Avg
Pretrained	52.05	54.90	51.33	50.12	48.96	47.39	39.24	49.14
Swa LoRA	44.78	55.06	45.50	48.27	41.93	37.59	45.46	45.51
Swa LAYRA (1,10)	49.24	53.98	46.22	48.47	44.58	43.78	47.39	47.67
Swa LAYRA (2,10)	49.92	55.42	46.63	50.64	44.98	43.94	46.39	48.27
Swa LAYRA (6,10)	50.04	55.70	47.51	48.31	45.42	45.82	46.47	48.47
Swa LAYRA (10,10)	49.68	54.98	46.67	48.84	44.02	46.39	47.15	48.25
Swa LAYRA (10,6)	50.36	54.74	47.15	48.47	46.91	46.83	46.83	48.76
Swa LAYRA (10,2)	49.92	54.22	47.43	49.24	46.83	46.35	45.34	48.48
Swa LAYRA (10,1)	49.04	54.82	48.92	48.55	45.90	45.06	46.99	48.47

Table 11: Ablation on LAYRA Configurations for Swahili XNLI

Model	deu	eng	spa	fra	swa	Avg
Pretrained	66.20	67.45	65.30	64.45	61.00	64.88
Swa LoRA	66.40	68.75	64.00	63.35	61.70	64.80
Swa LAYRA (1,10)	65.40	69.40	61.45	64.40	64.40	65.01
Swa LAYRA (2,10)	65.00	66.70	62.60	62.75	63.00	64.01
Swa LAYRA (6,10)	65.25	69.70	64.50	62.30	62.40	64.83
Swa LAYRA (10,10)	66.40	68.65	63.45	63.00	62.50	64.80
Swa LAYRA (10,6)	65.10	69.70	65.35	63.40	60.45	64.80
Swa LAYRA (10,2)	65.15	68.95	63.85	64.20	61.20	64.67
Swa LAYRA (10,1)	66.45	69.45	64.15	64.20	63.30	65.51

Table 12: Ablation on LAYRA Configurations for Swahili PAWS-X

Model	eng	spa	ita	tha	swa	Avg
Pretrained	87.00	81.40	72.60	57.60	55.00	70.72
Swa LoRA	88.00	70.60	62.00	56.60	66.20	68.68
Swa LAYRA (1,10)	86.00	68.00	56.40	57.80	63.60	66.36
Swa LAYRA (2,10)	85.00	69.60	59.60	57.20	64.60	67.20
Swa LAYRA (6,10)	88.00	70.00	60.80	57.20	65.00	68.20
Swa LAYRA (10,10)	86.00	70.40	58.00	56.60	65.40	67.28
Swa LAYRA (10,6)	88.00	69.80	59.60	58.60	66.00	68.40
Swa LAYRA (10,2)	87.00	69.60	61.20	57.80	64.60	68.04
Swa LAYRA (10,1)	85.00	69.40	59.00	56.40	62.80	66.52

Table 13: Ablation on LAYRA Configurations for Swahili XCOPA

Model	eng	spa	hin	swa	Avg
Pretrained	78.16	70.75	64.46	55.86	67.31
Swa LoRA	76.17	66.71	63.40	65.12	67.85
Swa LAYRA (1,10)	76.17	66.51	63.34	63.07	67.27
Swa LAYRA (2,10)	76.70	66.64	56.59	63.53	65.87
Swa LAYRA (4,10)	77.04	66.51	62.41	64.46	67.61
Swa LAYRA (10,10)	77.04	67.97	62.41	65.06	68.12
Swa LAYRA (10,6)	76.57	67.64	63.53	64.53	68.07
Swa LAYRA (10,2)	76.51	66.51	63.27	63.73	67.51
Swa LAYRA (10,1)	75.58	67.57	63.20	62.94	67.32

Table 14: Ablation on LAYRA Configurations for Swahili XStoryCloze

Model	deu	eng	spa	fra	hin	tha	swa	glg	Avg
Swa 0.0	50.84	55.81	50.64	46.22	48.07	46.43	34.66	54.02	48.34
Swa 0.1	49.76	54.50	51.16	47.79	47.91	44.62	37.23	55.54	48.56
Swa 0.2	50.32	54.18	51.45	48.92	48.63	44.82	39.52	54.74	49.07
Swa 0.3	50.16	54.58	51.37	49.60	48.31	44.54	40.92	54.32	49.23
Swa 0.4	49.24	54.78	51.24	49.64	48.55	44.30	42.21	53.24	49.15
Swa 0.5	48.88	54.74	50.96	48.63	47.95	44.82	43.86	52.01	48.98
Swa 0.6	47.55	54.30	49.92	47.59	46.67	45.14	45.06	49.69	48.24
Swa 0.7	46.55	54.06	48.80	46.72	45.18	45.18	46.34	47.55	47.55
Swa 0.8	46.22	53.09	47.15	44.66	44.02	44.86	46.71	44.66	46.42
Swa 0.9	45.34	50.88	44.18	41.08	40.68	44.74	46.59	41.92	44.43
Swa 1.0	44.90	48.15	42.01	40.28	40.28	44.18	45.30	37.23	42.79

Table 15: LAYRA-SERIES Ablation: Accuracy of varying λ' with Galician adapted model on XNLI

Model	deu	eng	spa	fra	swa	glg	Avg
Swa 0.0	67.40	67.15	64.30	63.30	50.35	65.25	62.96
Swa 0.1	66.60	67.15	64.85	64.20	52.85	63.65	63.22
Swa 0.2	65.95	67.10	64.50	63.45	53.80	63.10	62.98
Swa 0.3	65.55	67.20	63.55	63.25	54.80	61.50	62.64
Swa 0.4	65.15	67.40	63.60	62.95	54.80	59.15	62.18
Swa 0.5	65.10	66.80	62.90	63.50	53.75	55.65	61.28
Swa 0.6	65.10	65.95	62.40	61.60	52.45	51.85	59.89
Swa 0.7	63.75	65.00	60.75	60.30	51.30	48.50	58.27
Swa 0.8	62.40	64.25	60.70	58.80	50.80	47.10	57.34
Swa 0.9	61.55	63.15	58.70	56.65	50.80	46.55	56.23
Swa 1.0	61.55	62.10	59.00	55.40	52.75	46.15	56.16

Table 16: LAYRA-SERIES Ablation: Accuracy of varying λ' with Galician adapted model on PAWS-X

Model	eng	spa	ita	tha	swa	glg	Avg
Swa 0.0	86.00	76.80	60.20	58.40	54.40	63.20	66.50
Swa 0.1	83.00	78.80	61.40	58.80	56.40	62.00	66.73
Swa 0.2	83.00	77.40	63.20	57.60	57.20	61.80	66.70
Swa 0.3	86.00	77.00	65.20	57.60	57.00	61.20	67.33
Swa 0.4	83.00	76.00	64.60	57.60	58.20	61.40	66.80
Swa 0.5	81.00	74.40	64.20	57.40	59.40	61.80	66.37
Swa 0.6	79.00	71.60	64.00	55.80	61.80	60.00	65.37
Swa 0.7	78.00	69.60	62.60	55.60	62.20	59.00	64.50
Swa 0.8	79.00	67.20	62.20	56.60	63.00	56.80	64.13
Swa 0.9	80.00	63.80	59.20	56.80	62.80	55.40	63.00
Swa 1.0	78.00	62.60	56.00	57.00	62.40	55.80	61.97

Table 17: LAYRA-SERIES Ablation: Accuracy of varying λ' with Galician adapted model on XCOPA

Model	eng	spa	hin	swa	glg	Avg
Swa 0.0	76.11	69.69	63.20	50.69	69.89	65.92
Swa 0.1	76.37	69.82	63.27	52.22	70.28	66.39
Swa 0.2	76.44	69.69	64.00	54.20	70.68	67.00
Swa 0.3	76.37	69.16	64.00	55.72	69.95	67.04
Swa 0.4	76.17	68.83	64.46	56.92	68.70	67.02
Swa 0.5	75.91	68.23	64.13	58.44	67.31	66.80
Swa 0.6	75.58	67.90	63.60	59.43	65.25	66.35
Swa 0.7	74.98	66.98	63.80	60.75	62.94	65.89
Swa 0.8	74.32	65.85	64.00	61.02	61.15	65.27
Swa 0.9	72.67	63.40	63.67	61.48	58.44	63.93
Swa 1.0	70.62	61.68	63.53	61.48	56.19	62.70

Table 18: LAYRA-SERIES Ablation: Accuracy of varying λ' with Galician adapted model on XStoryCloze

Hyperparameter	Description	Value
Epochs	Training epochs	1
Batch Size	1 language: 32	32
	2 languages: 64	64
	3 languages: 128	128
Sequence Length	Maximum sequence length	2048
Warm-up Steps	Proportion of total optimization steps	5%
Learning Rate (α)	Initial learning rate	3e-4
Learning Rate Schedule		Linear
Weight Decay (λ)	Regularization parameter	0.1
Optimizer		AdamW
Epsilon (ϵ)	Optimizer stability parameter	1.0e-5
β_1	First moment decay rate	0.9
β_2	Second moment decay rate	0.95
GPUS	Hardware	H100 X2

Table 19: Hyperparameters for Experiments

Task	Swa	Urd	Glg
XNLI	✓	✓	✓
XStoryCloze	✓	-	✓
PAWS	✗	✗	✓
XCOPA	✓	✗	✓
MGSM	✓	✓	✗
MMLU-Lite	✓	✗	✗

Table 20: New languages Translated Using Google Machine Translate (GMT). All other languages used for our evaluation but not listed here were obtained from the original dataset release.

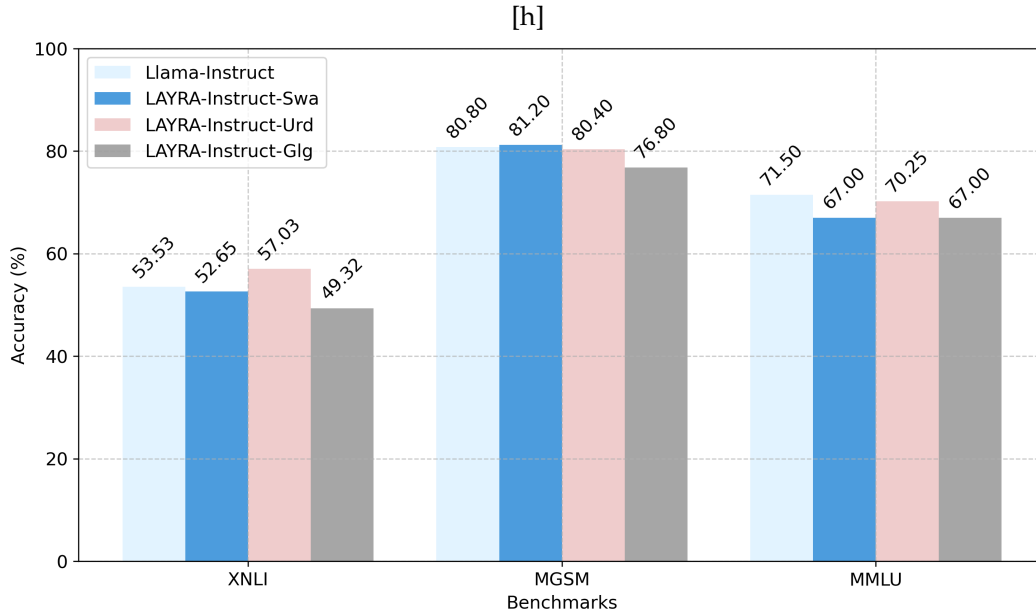


Figure 3: Accuracy of the instruction base model vs the LAYRA Instruct on XNLI, MGSM and MMLU for English