

# XDoGE: Multilingual Data Reweighting to Enhance Language Inclusivity in LLMs

Iñaki Lacunza<sup>†</sup>

Language Technologies Lab  
Barcelona Supercomputing Center  
Barcelona, Spain  
inaki.lacunza@bsc.es

José Javier Saiz<sup>†</sup>

Language Technologies Lab  
Barcelona Supercomputing Center  
Barcelona, Spain  
jose.saiz@bsc.es

Alexander Shvets<sup>†</sup>

Language Technologies Lab  
Barcelona Supercomputing Center  
Barcelona, Spain  
aleksandr.shvets@bsc.es

Aitor Gonzalez-Agirre

Language Technologies Lab  
Barcelona Supercomputing Center  
Barcelona, Spain  
aitor.gonzalez@bsc.es

Marta Villegas

Language Technologies Lab  
Barcelona Supercomputing Center  
Barcelona, Spain  
0000-0003-0711-0029

**Abstract**—Current large language models (LLMs) are trained on massive amounts of text data, primarily from a few dominant languages. Studies suggest that this over-reliance on high-resource languages, such as English, hampers LLM performance in mid- and low-resource languages. To mitigate this problem, we propose to (i) optimize the language distribution by training a small proxy model within a domain-reweighting DoGE algorithm that we extend to XDoGE for a multilingual setup, and (ii) rescale the data and train a full-size model with the established language weights either from scratch or within a continual pre-training phase (CPT). We target six languages possessing a variety of geographic and intra- and inter-language-family relations, namely, English and Spanish (high-resource), Portuguese and Catalan (mid-resource), Galician and Basque (low-resource). We experiment with Salamandra-2b, which is a promising model for these languages. We investigate the effects of substantial data repetition on minor languages and under-sampling on dominant languages using the IberoBench framework for quantitative evaluation. Finally, we release a new promising IberianLLM-7B-Instruct model centering on Iberian languages and English that we pretrained from scratch and further improved using CPT with the XDoGE weights.

**Index Terms**—linguistic diversity, multilingual data distribution optimization, large language model pretraining, low-resource languages

## I. INTRODUCTION

The development of large language models (LLMs) has predominantly focused its advances on high-resource languages like English, Spanish or Chinese [1]–[4]. Given that most other languages may not have enough data to train LLMs from scratch, the only option left for them is to be included among the target languages of a multilingual model. Multilingual LLMs have the potential to handle multiple languages comprehensively; however, approaches to dealing with unbalanced pre-training data have traditionally been conservative and have not specifically addressed how to improve performance in underrepresented languages, leaving their contribution to

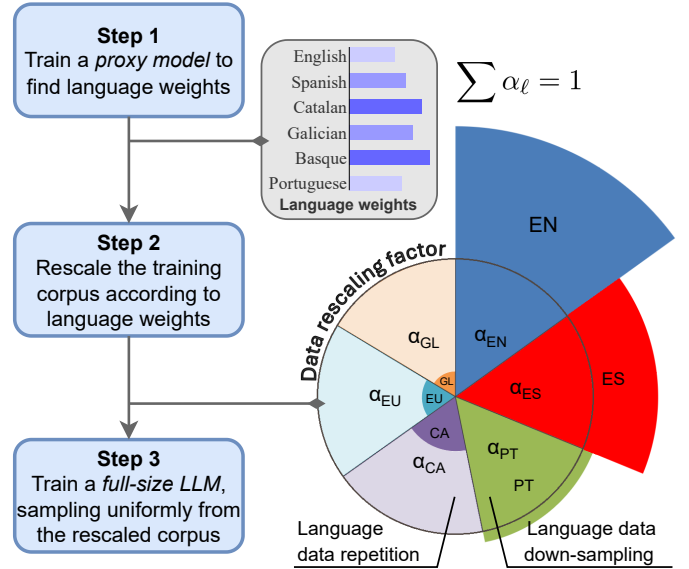


Fig. 1. Summary of our proposed XDoGE method.

the mixture of training data residual compared to dominant languages [5]–[7].

We propose to enhance the linguistic diversity of LLMs by building upon a state-of-the-art method in domain reweighting, DoGE [8], which optimizes the distribution of sources, distinguishing those that help the most to model a language within a given set of domains. In a pilot study, [8] showed the potential of their approach (primarily designed for domains in a single language) to generalize from a mixture of high-resource languages to a certain out-of-mixture target language, when there is a grammatical relatedness. However, the possibility of generalization to several languages at once as well as addressing constraints on data availability for resampling the pre-train corpus remained unexplored.

<sup>†</sup> Equal contribution.

We adapt the DoGE method to a multilingual scenario, resulting in XDoGE. In particular, we modify DoGE to handle multilingual sources that we split beforehand by language and weight individually, thereby avoiding reliance on a prior per-source language distribution, which is often skewed towards high-resource languages. We also introduce a weight-clipping mechanism so that no language becomes excessively over- or under-represented to promote fair modeling of all the target languages. We also rescale the training corpus to satisfy the optimized language distribution. This allows us to down-sample data for dominant languages and compensate for the lack of data in mid- and low-resource languages without compromising the model’s performance for most languages. The summary of our approach is provided in Fig. 1.

We show that XDoGE succeeds in generalizing to languages with distinct grammars, and overcomes the uniform-weighting baseline in many languages when small base models are trained from scratch as well as within a continual pre-training phase (CPT) for larger models.

Our main contributions are:

- We introduce the XDoGE framework, based on DoGE, which re-weights sources in multilingual pre-training data for improved generalization on the target languages. We add mechanisms to handle multiple sources in different languages in an imbalanced data scenario (see Section III).
- We demonstrate that XDoGE is scalable to various proxy and base model sizes and yields improvements for many languages in both from-scratch pre-training and continual pre-training scenarios, as measured by perplexity and a range of downstream metrics in few-shot tasks (see Section V). Note that in studies with relatively small models within a limited token budget, we could aim only at minimal signals of linguistic proficiency improvement.
- We also release IberianLLM-7B-Instruct<sup>12</sup> centered in Iberian languages and English, well-performing in various high-, medium-, and low-resource languages that meet our aim to promote linguistic inclusion in the field.

## II. RELATED WORK

While there have been significant efforts to develop multilingual LLMs, most attempts have sampled pre-training data with heuristic approaches, which lack generalization and still may be biased towards high-resource languages. In most cases, heuristic approaches only consider data size and ignore language similarity, which has been found important for cross-lingual transfer learning [9]–[11].

On the other hand, it is a natural assumption in machine learning that principled techniques yield more generalizable weightings than manually tuned hyperparameters [12]. Among

the latest principled techniques, “offline” data weighting algorithms optimize the domain weights before actually training the language model, providing an optimized distribution of data from the start. DoReMi [13] optimizes domain weights by aiming to achieve low loss on all domains, formalizing this problem as a group distributionally robust optimization [14], [15]. Specifically, DoReMi optimizes the domain weights with a proxy model by reducing the excess loss on each domain compared to the loss on a previous reference model, and then uses the optimized domain weights to train a larger model. More recently, the DoGE algorithm [8] has been proposed and shown to be scalable to different model sizes. Similar to DoReMi, the algorithm involves a two-step process by first training a small **proxy** model that uses the gradient alignment between each domain and all other target domains in order to give more weight to domains with better transfer to other domains. Second, a larger **base** model is trained using data sampled according to the final weight distribution. However, to the best of our knowledge, there has not yet been an attempt to reproduce these principled approaches in a setting that covers many target languages. There are studies on multilingual data mixing, but the weights are often chosen manually and concern data types rather than languages [16].

## III. METHODOLOGY

Building on the offline domain reweighting algorithm DoGE [8], we adapt it to enhance cross-lingual transfer in model pre-training for a set of target languages, resulting in XDoGE. Our adaptation concerns (i) the ability to target a combination of languages and domains within a multilingual corpus aiming to weight language-domain sources based on their contribution to the overall cross-lingual learning; (ii) an introduction of a threshold mechanism that ensures an adequate number of instances for each language used at every distribution update step, maintaining its representation up to the weight convergence; (iii) a corpus rescale assessment to balance an undesired cut in data for high-resource languages and excessive data repetition for low-resource languages. In the following, we outline XDoGE, adapting the original notation.<sup>3</sup>

### A. Proxy Model Training.

Let  $S$  denote the set of sources (domains) and  $\mathcal{L}$  – the set of languages with  $k_\ell$  sources in a language  $\ell \in \mathcal{L}$ :  $S_\ell = \{S_{\ell,1}, S_{\ell,2}, \dots, S_{\ell,k_\ell}\}$  (e.g.,  $S_{ES,1}$  – Spanish Wikipedia,<sup>4</sup>  $S_{ES,2}$  – Spanish part of Community OSCAR<sup>5</sup>). We optimize sampling weights for all sources  $k = \sum_{\ell \in \mathcal{L}} k_\ell$  over the probability simplex  $\alpha \in \Delta_k \subset \mathbb{R}^k$  using an original bi-level optimization framework adjusted to our split of domains into language-based subsets with an addition of a minimum weight constraint:

<sup>1</sup><https://huggingface.co/langtech-language modeling/IberianLLM-7B-Instruct>

<sup>2</sup>We release the instruction-tuned version of the model developed in this work to facilitate its integration into general-purpose applications. See the model card for details on instruction-tuning.

<sup>3</sup>We disclose only core elements of DoGE. Cf. [8], for the individual decisions in the optimization design.

<sup>4</sup><https://es.wikipedia.org/>

<sup>5</sup>Subsets with the `es` code from <https://huggingface.co/datasets/oscar-corpus/community-oscar>

$$\begin{aligned}
\alpha &\in \arg \min_{\alpha \in \Delta_k} \sum_{\ell \in \mathcal{L}} \sum_{i=1}^{k_\ell} l_{\ell,i}(\theta^*(\alpha)), \\
s.t. \quad \theta^*(\alpha) &\in \arg \min_{\theta} \sum_{\ell \in \mathcal{L}} \sum_{i=1}^{k_\ell} \alpha_{\ell,i} l_{\ell,i}(\theta), \\
s.t. \quad \alpha_{min} &\geq \gamma; \gamma \lll |S|^{-1},
\end{aligned}$$

where  $l_{\ell,i}(\theta)$  is the next-token prediction (cross-entropy) loss of the model of parameters  $\theta$  on  $S_{\ell,i}$ , and  $\gamma$  is a constant.

A train batch at time-step  $t$  is sampled from the dynamically updated instance-wise distribution:

$$P_\alpha \triangleq \sum_{\ell \in \mathcal{L}} \sum_{i=1}^{k_\ell} \alpha_{\ell,i}^{(t)} \cdot \text{UNIF}(S_{\ell,i}). \quad (1)$$

where  $\alpha_{\ell,i}$  is the weight for source  $S_{\ell,i}$ , and  $\text{UNIF}(S_{\ell,i})$  is a uniform batch-sampling over  $S_{\ell,i}$ .

### B. Bi-level Optimization.

The framework alternates between two key steps: the inner loop (2), which updates the model parameters  $\theta$  using weighted losses, and the outer loop (3-6) which updates weights  $\alpha$  based on gradient alignments across sources (domain-language pairs that are better aligned with the other pairs receive a larger weight).

The  $\theta$  parameters are updated as (*the inner loop*):

$$\theta^{(t+1)} = \theta^{(t)} - \eta^{(t)} \sum_{\ell \in \mathcal{L}} \sum_{i=1}^{k_\ell} \alpha_{\ell,i}^{(t)} \nabla l_{\ell,i}(\theta^{(t)}), \quad (2)$$

where  $\alpha^{(t)}$  is used to re-weight the loss from sources at time-step  $t$ ,  $\eta(t)$  is the step size, and  $\nabla l_{\ell,i}(\theta(t))$  is a stochastic gradient for  $S_{\ell,i}$  samples.

The weights  $\alpha$  are updated as (*the outer loop*):

$$W_{\ell,i}^{(t)} = \left\langle \nabla l_{\ell,i}(\theta^{(t)}), \sum_{\ell' \in \mathcal{L}} \sum_{j=1}^{k_{\ell'}} \nabla l_{\ell',j}(\theta^{(t)}) \right\rangle, \quad (3)$$

$$\hat{\alpha}_{\ell,i}^{(t)} = \alpha_{\ell,i}^{(t-1)} \odot \exp \left( \frac{\eta^{(t)} W_{\ell,i}^{(t)}}{\mu} \right), \quad (4)$$

$$\tilde{\alpha}_{\ell,i}^{(t)} = \frac{\hat{\alpha}_{\ell,i}^{(t)}}{\sum_{\ell' \in \mathcal{L}} \sum_{j=1}^{k_{\ell'}} \hat{\alpha}_{\ell',j}^{(t)}}, \quad (5)$$

$$\alpha_{\ell,i}^{(t)} = \text{Proj}_{\Delta_{k,\gamma}} \left( \tilde{\alpha}_{\ell,i}^{(t)} \right). \quad (6)$$

where  $W_{\ell,i}^{(t)}$  is the stochastic generalization estimation function (a higher value means learning  $S_{\ell,i}$  will also contribute to learning other languages and domains),  $\mu$  - a hyperparameter for the strength of regularization,  $\text{Proj}_{\Delta_{k,\gamma}}$  - our new projection operator that ensures  $\alpha_{\ell,i} \geq \gamma \forall \ell, i$  while preserving  $\sum_{\ell,i} \alpha_{\ell,i} = 1$ . The projection is applied as:

- 1) Clip weights:  $\alpha'_{\ell,i} = \max(\tilde{\alpha}_{\ell,i}, \gamma)$
- 2) Compute excess weight:  $E = \sum_{\ell,i} \alpha'_{\ell,i} - 1$
- 3) Redistribute  $E$  proportionally from unclipped sources:

$$\alpha_{\ell,i}^{(t)} = \begin{cases} \alpha'_{\ell,i} - E \cdot \frac{\alpha'_{\ell,i}}{\sum_{(\ell',j) \in U} \alpha'_{\ell',j}}, & (\ell, i) \in U \\ \gamma, & \text{otherwise,} \end{cases} \quad (7)$$

where  $U = \{(\ell, j) \mid \alpha'_{\ell,j} > \gamma\}$ . This prevents language under-representation while promoting domain diversity within a language. This threshold mechanism adds negligible overhead compared to standard training, preserving DoGE's computational efficiency.

### C. Final Language Weight Assignment.

As the algorithm avoids relying on the prior domain size and therefore the obtained weights within a language could be excessively disproportional to the number of tokens per domain available in the corpus for language modeling, we assign a single domain-independent language weight for the final sampling. The final language weight  $\bar{\alpha}_\ell$  is computed as the sum of its domain weights:

$$\bar{\alpha}_\ell = \sum_{i=1}^{k_\ell} \alpha_{\ell,i}. \quad (8)$$

Finally, training instances are sampled for each batch based on the language weights:

$$P_{\bar{\alpha}} \triangleq \sum_{\ell \in \mathcal{L}} \bar{\alpha}_\ell \cdot \text{UNIF}(S_\ell). \quad (9)$$

This mitigates the lack of domain-language data to support the allocated weight and makes large-scale training practically feasible.

### D. Corpus Rescaling at Large Model Training.

The data for full-size model training is usually highly unbalanced across languages. Therefore, if a language is assigned a higher probability than its token count in a corpus would suggest, the corresponding part of the corpus will be sampled several times. Conversely, if the assigned weight is lower, the corresponding language instances will be sampled fewer times than the corpus allows (see Fig. 1).

In our approach, we aim to find a trade-off between over- and under-sampling by continuously evaluating intermediate checkpoints, resulting in a data rescaling factor respectful of all languages. This ensures that low-resource languages, gaining more repetitions, are sufficiently represented without causing model degradation, while for high-resource languages, a considerable portion of the available data is utilized. We train the *base* models from scratch, and *pre-trained* models within a continual pre-training phase (CPT). In both cases, we use the standard next-token prediction loss.

## IV. EXPERIMENT SETUP

### A. Data

We collect a training dataset to test the effect of languages on cross-lingual transfer during pre-training and continual pre-training. We choose a set of 6 languages, including English, Spanish, Catalan, Galician, Basque and Portuguese, which represent a group of four typologically related Romance languages (along with English, which is typologically more

distant), with varying degrees of availability of NLP resources [17]. Basque, on the other hand, is a language isolate that provides a contrast to the Romance languages, which can help to assess how structural and lexical similarities, as well as resource availability, influence cross-lingual transfer performance. We note that our approach can be readily extended to other languages beyond those selected in this study.

For the XDoGE base and proxy trainings, we source the data from Wikipedia dumps from May 2024, as it consists of roughly similar encyclopedic register across languages; and from Community OSCAR<sup>6</sup> [18], which is a processed multilingual corpus from web-crawled data, and represents a more heterogeneous source of data for each language. Both sources are widely used for the training of LLMs [5]. For the CPT phase, we replace web data with FineWeb-Edu<sup>7</sup> for English and FineWeb2<sup>8</sup> for other languages [19] to ensure that pre-trained models that used OSCAR as a major source are exposed to new information, allowing them to show improvement with a smaller token budget.

Within each dataset, we perform exact document deduplication for all languages (to control the information repetition across training epochs better) and random shuffling followed by tokenizing with the six-language tokenizer described in Section IV-B. The dataset sizes in tokens are shown in Table I.

TABLE I  
SIZE IN BILLIONS (B) OF TOKENS OF THE TRAINING CORPUS FOR EACH LANGUAGE, WHERE PROXY AND BASE TRAININGS USED OSCAR AND WIKIPEDIA, AND CPT TRAININGS USED FINEWEB-EDU, FINEWEB2, AND WIKIPEDIA. VALIDATION AND TEST SETS REPRESENT EACH 1% OF THE TRAINING DATA.

Language	OSCAR	Wikipedia	FineWeb2/-Edu
English	327,98 B	4,76 B	198,54 B
Spanish	140,61 B	1,19 B	159,36 B
Portuguese	62,12 B	0,59 B	117,00 B
Catalan	3,48 B	0,45 B	9,37 B
Basque	0,33 B	0,12 B	3,29 B
Galician	0,11 B	0,10 B	1,50 B
Total	534,62 B	7,21 B	489,07 B

## B. Technical Setup

a) *Proxy models*: While [8] run the experiments at three different scales up to 125M parameters and demonstrate effective weight estimation at 10k training steps within a monolingual setup, we allocate more resources to account for potentially more demanding domain-language interactions and consider proxy models with 70M, 125M, 250M, and 500M parameters (see Table II).<sup>9</sup> We perform longer runs of 10k-50k steps to validate the robustness of our approach across different-sized models.

To train the proxy models, we use the DoGE PyTorch-based framework<sup>10</sup> and choose LLaMA architecture [20] to

be compatible with Salamandra [21] used for the *base* and *pretrained models* in this work.

We use a maximum context length of 128 with a global batch size of 16,384 tokens, a learning rate of  $5e^{-4}$ , a linear warmup cosine weight decay of  $1e^{-2}$ , a warmup ratio of 0.05, and 500 warmup steps. We set  $\gamma = 0.02$  empirically: with a batch size of 128 document instances, this threshold ensures that at least several instances from each source would be seen in each iteration.

b) *Base models*: We train 250M, 500M and 900M base models with several reweighted data configurations (the resulting optimized weights are provided in Section V-A) and compare them in their respective sizes against *baselines* that we train with uniform source weights, which is the preferred heuristic for cross-lingual generalization without prior knowledge on language relatedness. Each model has a context length of 8,192 tokens and was trained with a global batch size of  $\sim 4M$  tokens, equivalent to 512 instances per batch.

We train the models for about 150k steps to observe up to 30% of the English corpus, although making an extreme number of repetitions for downstream languages (e.g., up to 450 repetitions of Galician; cf. Fig. 1 for better illustration).

The base models were trained using the NeMo framework<sup>11</sup> for its efficient parallelization and scalability in large-scale training. We use the architecture of Salamandra to ease the comparison with the *pre-trained models* from the same model family described in the next subsection.

c) *Pretrained models*: To mitigate excessive data repetition as happens with the base models, we tune already pre-trained models in a continual pretraining stage (CPT) that potentially requires fewer data to adapt to target languages (we limit to 80k steps).

We selected the 2b-parameter variant from the Salamandra family of models [21]. This model was originally pre-trained on a corpus spanning 35 European languages and programming code, including all the languages targeted in our work.

We also pretrain IberianLLM-7b of the Salamandra architecture from scratch, focusing on six target languages (using the Salamandra pretraining data for corresponding languages).

For both models, we use Salamandra’s tokenizer designed for multilingual support, based on byte-pair encoding [22] with a vocabulary of 256,000 tokens and an equal number of documents per language to ensure fair representation across languages [21], [23]. We train the tokenizer on 6 languages and use it for proxy, base, and IberianLLM models. For Salamandra-2b we use its pre-trained tokenizer.

## V. EXPERIMENTS

In this section, we discuss the findings from training the XDoGE proxies and its version without the thresholding mechanism (unthresholded proxy), as well as the choice of final language weights averaged from several runs and the results from applying XDoGE weights to from-scratch and continuous pre-training scenarios. We use a comprehensive set

<sup>6</sup>We used 35 monthly dumps of the Community OSCAR corpus.

<sup>7</sup>datasets/HuggingFaceFW/fineweb-edu

<sup>8</sup>datasets/HuggingFaceFW/fineweb-2

<sup>9</sup>Rounded sizes correspond to actual of 68M, 122M, 273M, and 509M.

<sup>10</sup><https://github.com/Olivia-fsm/doge>

<sup>11</sup><https://docs.nvidia.com/nemo-framework/>

TABLE II  
DIMENSIONS OF THE PROXY MODELS. THE LLAMA-2 ARCHITECTURE HAS BEEN USED TO BUILD THESE MODELS.

Model ID	Proxy-70M	Proxy-125M	Proxy-250M	Proxy-500M
Parameters	67,936,768	122,358,528	273,122,304	509,199,360
Layers	8	12	16	24
Hidden Size	512	768	1,024	1,024
FFN Size	512	512	2024	4096
Attention Heads	8	12	16	16
K/V Heads	8	12	16	16
Context Length	128	128	128	128
Vocabulary Size	52,000	52,000	52,000	52,000

of few-shot tasks from a recently published benchmark for the Iberian languages, IberoBench [24], to evaluate the models.

IberoBench comprises 60+ tasks that evaluate various language capabilities of LLMs. The tasks are of two main types: multiple choice (*accuracy* metric) and open-ended generation (BLEU [25] and ROUGE [26]). It extends LM-Evaluation Harness,<sup>12</sup> reusing the Harness execution pipeline but supplies additional task YAMLS, dataset splits and prompt templates for Iberian languages. We run experiments with the same IberoBench and Harness configurations; therefore, the prompt text, the five examples (chosen with a fixed default seed for reproducibility), their order and formatting, and generation parameters are identical across models for all languages.

#### A. Proxy Model Training

Training unthreshold proxy models often results in some domain language weights dropping and stabilizing at near-zero values soon after the initial training steps (see Fig. 2). Since the proxy model is trained with the updated weights at each step, once a given language reaches near-zero weight values, the proxy model no longer receives a meaningful amount of data from that language and cannot recover its weight. On the other hand, languages that are more distinct, and therefore more “difficult” to learn, may receive increasingly larger weights if the threshold is not used (e.g., as seen with Basque). This may negatively affect the entire distribution.

Table III shows the benefits of using a minimum threshold for domain-language weights: a) XDoGE converges stably to more similar distributions,<sup>13</sup> avoiding drastic skews towards certain languages; b) the unthresholded variants often reverse the decision about which domain within a language to take for most of the tokens,<sup>14</sup> (both or only one, with strong but not exclusive preferences for Wikipedia of higher data quality), while the thresholded setup ensures that both domains contribute rather equally to the modeling of the language (no domain gets the minimum possible weight). This leads to the favorable conclusion that *learned weights can be assigned to languages independent of the domain*, providing flexibility in domain selection and token assignment without having to retrain proxy models when new domains are added.

a) *Weight dynamics visualization*: To visually validate the weight stability achieved through our thresholded approach, Fig. 2 compares the training dynamics of language-domain weights between our implementation and the original DoGE configuration. The side-by-side comparison demonstrates how thresholding prevents extreme divergences while maintaining the competition between sources and stable weight evolution. Notably, even if the optimization results in the largest weight for a high-resource language, the remaining are represented more fairly than with a natural ad hoc distribution. For example, the final weights of 13% for Galician and 20% for English in Fig. 2, in fact, correspond to an increase from a natural 0.19%, and a decrease from a natural 44.65%.

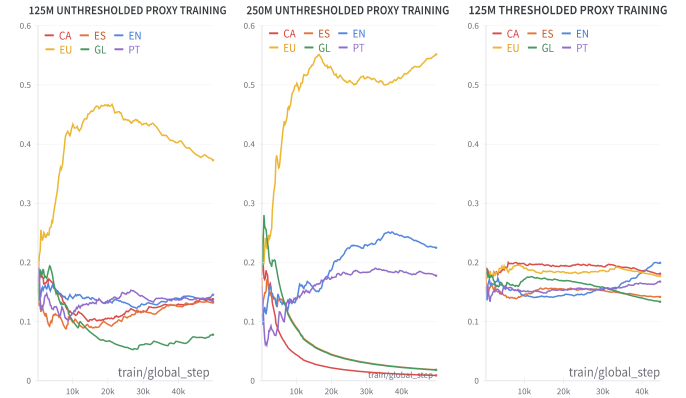


Fig. 2. Language weight updates during the proxy model training without (125M and 250M on the left) and with the threshold mechanism (125M on the right).

b) *Training stability*: A critical requirement for any weighting adaptation is maintaining stable training dynamics. Fig. 3 provides this validation by showing the loss trajectories across all language sources in our thresholded setup. The coordinated downward trend confirms that our modifications preserve the fundamental training behavior: no language is destabilized. Fig. 4 quantifies how thresholding improves overall optimization.

c) *Model scaling and weight selection*: We also observe in Table III that *the larger the proxy model the more similar the distributions are and especially stable for the thresholded setup*. Although intuitively they could be considered the best found, we opt for averaging across various proxy model sizes and different seed runs following [8]. Since XDoGE features

<sup>12</sup><https://github.com/EleutherAI/lm-evaluation-harness>

<sup>13</sup>Compare  $D_{KL}$  values for language weights.

<sup>14</sup>Compare  $D_{KL}$  values for domain-language weights.

TABLE III

DOMAIN-LANGUAGE (D-L) AND LANGUAGE (L) WEIGHTS WITH DIFFERENT SIZES OF PROXY MODELS SMOOTHED OVER 10K TRAINING STEPS.  $D_{KL}$  – KL DIVERGENCE  $\cdot 100\%$ ; AGAINST THE DISTRIBUTION OBTAINED BY THE 500M MODELS.

		Unthresholded								Thresholded							
Domain-Lang	Lang	70M		125M		250M		500M		70M		125M		250M		500M	
		D-L	L	D-L	L	D-L	L	D-L	L	D-L	L	D-L	L	D-L	L	D-L	L
OSCAR-EN	EN	6.37	13.98	5.66	13.96	1.41	13.94	1.82	12.74	7.54	14.81	6.84	14.39	6.29	13.94	6.41	14.05
WIKI-EN		7.61		8.30		12.53		10.92		7.27		7.55		7.65		7.64	
OSCAR-ES	ES	7.71	13.51	1.82	9.95	0.42	9.00	0.39	13.53	8.31	14.74	7.96	15.17	7.37	17.28	7.95	17.83
WIKI-ES		5.80		8.13		8.58		13.14		6.43		7.21		9.91		9.88	
OSCAR-CA	CA	13.01	19.85	5.30	11.46	2.59	4.39	1.84	6.69	11.46	19.08	10.62	18.17	8.13	16.27	9.57	18.65
WIKI-CA		6.84		6.16		1.80		4.85		7.62		7.55		8.14		9.08	
OSCAR-GL	GL	0.99	14.24	0.72	10.22	0.78	8.81	0.62	7.69	6.37	14.50	6.65	15.91	7.30	16.86	7.56	17.44
WIKI-GL		13.25		9.50		8.03		7.07		8.13		9.26		9.56		9.88	
OSCAR-EU	EU	9.50	24.54	7.85	43.08	7.30	50.00	0.28	44.52	11.25	21.85	9.75	21.54	7.22	19.40	6.81	16.36
WIKI-EU		15.04		35.23		42.70		44.24		10.60		11.79		12.18		9.55	
OSCAR-PT	PT	9.61	13.88	2.85	11.33	0.68	13.86	1.14	14.84	9.21	15.04	7.72	14.86	7.30	16.25	7.28	15.67
WIKI-PT		4.27		8.48		13.18		13.70		5.83		7.14		8.95		8.39	
Total		100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0	100.0
$D_{KL}$ (D-L, base 500M)		92.85	-	29.74	-	19.48	-	0.00	-	3.30	-	1.60	-	0.56	-	0.00	-
$D_{KL}$ (L, base 500M)		-	16.11	-	2.83	-	1.80	-	0.00	-	1.42	-	1.05	-	0.45	-	0.00

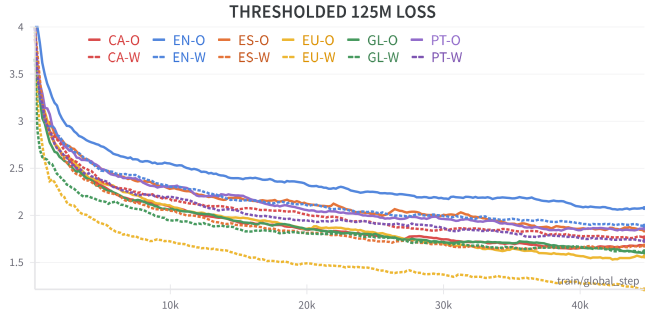


Fig. 3. Thresholded proxy train loss per source. The dotted lines – Wikipedia, continuous – OSCAR.

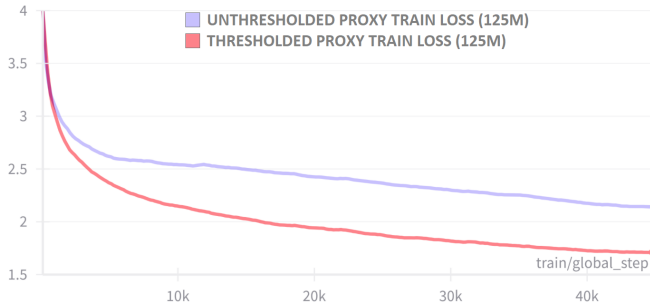


Fig. 4. 125M proxy train loss for DoGE and XDoGE.

less stable convergence than DoGE across model sizes, we select two averaging schemes to test how sensitive the training is to the deviations in the language balancing: averaging over our four largest models (two 125M models of different seeds, 250M, and 500M models) and three middle-size models (500M model is excluded) using the last checkpoints.<sup>15</sup> We do not include 70M models that give the largest weight dispersion with a distribution significantly diverging from those of larger models. We also consider a uniform distribution for baselines

<sup>15</sup>Due to computing limitations we trained the 500M model only to up to 11k steps and smaller models to up to 50K steps. The number of models is also subject to available computing.

(i.e., sampling the same number of documents per language). The final sets of weights are provided in Fig. 5.

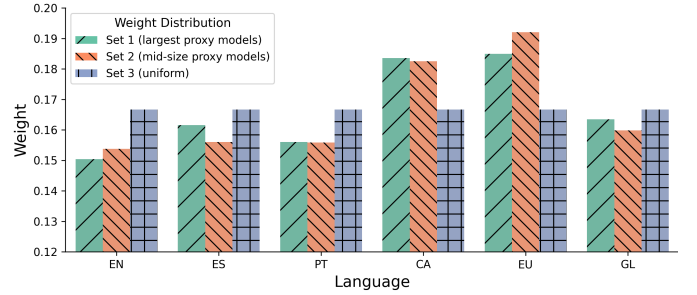


Fig. 5. Selected weight distributions.

TABLE IV  
PERPLEXITY SCORES (THE LOWER THE BETTER) AND IBEROBENCH AVERAGE ACCURACY SCORES (THE HIGHER THE BETTER) FOR THE MEDIUM 500M BASE MODELS.

Language	Perplexity			IberoBench (acc)		
	Bsl.	Set1	Set2	Bsl.	Set1	Set2
CA	14.43	<b>14.42</b>	<b>14.42</b>	41.01	<b>41.21</b>	40.29
PT	<b>31.33</b>	31.52	32.30	51.73	<b>53.19</b>	43.66
EN	31.56	<b>31.40</b>	32.13	38.15	38.03	<b>38.87</b>
GL	<b>29.38</b>	32.07	29.71	32.52	<b>33.00</b>	32.60
ES	35.15	35.30	<b>35.13</b>	43.43	44.17	<b>44.73</b>
EU	<b>23.11</b>	23.61	23.55	37.06	36.12	<b>37.24</b>
Total Avg	<b>27.49</b>	28.05	27.87	<b>40.65</b>	<b>40.95</b>	39.57

### B. Base Models

We trained base models of varying sizes using three weight setups from Fig. 5. Tables IV and V compare Set1 and Set2 against the uniform-set baseline.<sup>16</sup>

500M XDoGE exhibits the highest accuracy across all languages, but shows lower perplexity for some (whereas 250M surpasses the baseline across all tasks). We do not find

<sup>16</sup>The evaluations are performed on all accuracy tasks on IberoBench.



TABLE V  
IBEROBENCH AVERAGE ACCURACY SCORES (THE HIGHER THE BETTER)  
FOR SMALL 250M AND LARGE 900M BASE MODELS.

Language	250M			900M		
	Bsl.	Set1	Set2	Bsl.	Set1	Set2
CA	<b>39.88</b>	39.32	39.40	41.85	40.88	<b>42.07</b>
PT	46.00	<b>52.32</b>	48.85	48.64	<b>48.83</b>	47.24
EN	37.03	<b>37.18</b>	36.25	38.17	<b>38.65</b>	38.22
GL	32.64	33.65	<b>33.68</b>	31.05	<b>32.10</b>	31.22
ES	<b>45.02</b>	43.63	43.89	<b>43.15</b>	40.64	41.87
EU	36.75	<b>36.87</b>	36.24	<b>37.07</b>	36.47	36.76
<b>Total Avg</b>	39.55	<b>40.50</b>	39.72	<b>39.99</b>	39.60	39.56

stable language-set associations across model sizes, but, in general, more languages benefit from Set1.

### C. IberianLLM-7b Pre-training

This subsection outlines the architecture, technical setup, and training methodology employed to develop IberianLLM. It also discusses the rationale behind our design choices, particularly in handling the unique challenges posed by low-resource Iberian languages.

*a) Challenges in Training LLMs with Low-Resource Languages:* Training large language models for low-resource languages is inherently challenging due to the limited availability of high-quality data.

Maximizing data utilization is crucial, yet even with careful curation, the overall volume for low-resource languages remains much smaller compared to high-resource ones like Spanish or English. A common strategy to mitigate this issue is to pre-train models on diverse linguistic distributions while filtering out low-quality sources [27].

Table I shows that the token counts for each language clearly indicate the impracticality of training a very large model from scratch with our available data. In particular, when aiming for a balanced representation across languages, one must consider that effective large-scale training typically requires around 30 tokens per parameter [28]. Given that Galician, the lowest-resourced language in our set, comprises approximately 213 million tokens, a balanced approach would cap each language to a similar token count. This results in a combined dataset of roughly 1.28 billion tokens across six languages, for a total token consumption of about 3.84 billion tokens over three epochs. According to the heuristic, this corresponds to a model with approximately 128 million parameters, ensuring that no language is oversampled at the expense of another.

#### *b) Optimizing Training Data for Target Languages:*

The Salamandra family of models was pre-trained on data spanning 35 European languages and code, which includes all of our target languages. However, instead of adopting a highly multilingual model, we opted to focus on a more specialized set: Iberian languages, English (selected for its abundant data), and code. In addition to our primary targets – English, Spanish, Catalan, Basque, Galician, and Portuguese – we also included Occitan, Aragonese, and Balearic. These additional languages,

due to their close linguistic ties to our main Iberian languages, help enhance the model’s ability to capture related language structures.

Notably, while the Salamandra models upscaled Iberian languages by doubling their data share, our approach inherently centers on these languages, eliminating the need for further upscaling. Moreover, to avoid an English-centric bias, we adjusted the composition of the English data. Instead of merely subsampling the original English dataset used in Salamandra, we replaced the Oscar portion with the higher-quality FineWeb-Edu dataset [19]. This change reduced the overall quantity of English data – bringing it closer in line with Spanish, the next most-resourced target language – while significantly improving data quality.

Overall, the IberianLLM model was pre-trained on approximately  $0.53T$  unique tokens, in contrast to the  $2.6T$  unique tokens used for the Salamandra models. Fig. 6 compares the language distribution in our pre-training data with that of the Salamandra family, underscoring our targeted focus.

*c) Architecture & Training Strategy:* Table VI shows the architectural configuration of our model, which mirrors the 7b version of the Salamandra model.

We trained IberianLLM using the NeMo framework, following a multi-epoch strategy similar to that of the Salamandra models. Pre-training was conducted over three epochs to ensure comprehensive exposure to our curated dataset. Table VII outlines the key training parameters.

Further details on the tokenizer implementation and training can be found in Section IV-B, and additional insights are provided in the Salamandra technical report [21].

### D. Continual Pretraining Phase

Within CPT, we used only Set1 weights of XDoGE (since they showed superiority in most cases across various base model sizes). We focus on comparing with uniform weights and ad hoc weights (i.e., defined by the original data distribution).

Table VIII summarizes the evaluation, which shows that Salamandra-2b with XDoGE achieves a lower or equal perplexity in all languages, suggesting improved overall generalization on the target languages. On IberoBench tasks, our CPTs achieve higher average scores over the original Salamandra-2b and IberianLLM-7b and the baselines (see also the relative improvements in dynamic over the uniform baseline in Fig. 7), with Catalan, Spanish and Basque benefiting the most. English encounters under-sampling, while Portuguese and Galician received lower weights, preventing them from sufficiently reinforcing each other.

As we used the updated web data for all CPT setups, including baselines, we conclude that the new information is not solely responsible for the gains over the original models (*No-CPT*). This justifies the key role of our weights. Over-repetition of the data also contributes to the gains. Although uniform weights also imply this, we observe that our weights benefit the model in a longer training run, delaying model degradation: a slight decrease after 60k-80k steps for XDoGE

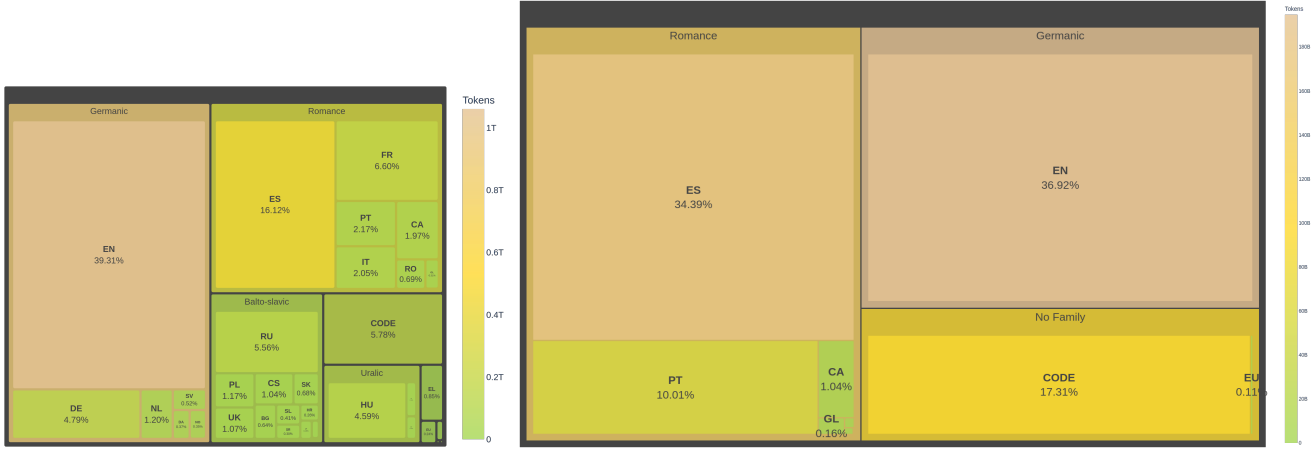


Fig. 6. Data distribution weights for each language during Salamandra pre-training (left) and 7b model pre-training (right). The total number of tokens for each model are  $2.68T$  and  $542B$ , respectively.

TABLE VI

DIMENSIONS OF THE BASE MODELS AND THE LARGE ARCHITECTURES WE CREATED AND USED FOR CPT EXPERIMENTS. THESE MODELS HAVE BEEN CREATED USING THE SALAMANDRA ARCHITECTURE.

Model ID	Base-250M	Base-500M	Base-900M	Salamandra-2b	IberianLLM-7b
Parameters	277,881,344	497,050,368	901,825,536	2,253,490,176	7,768,117,224
Layers	8	16	24	24	32
Hidden Size	512	768	1,024	2,048	4,096
FFN Size	768	2,048	4,096	5,440	11,008
Attention Heads	8	16	16	16	32
K/V Heads	4	8	8	16	8
Context Length	8,192	8,192	8,192	8,192	8,192
Vocabulary Size	256,000	256,000	256,000	256,000	256,000

TABLE VII

TRAINING HYPERPARAMETERS AND OPTIMIZER SETTINGS FOR PRE-TRAINING THE IBERIANLLM.

Training Configuration	Value
Epochs and Steps	3 training epochs, 377215 steps
Context Length	8192
Global Batch Size	512
Optimizer Name	Distributed Fused Adam
Learning Rate (lr)	$3.0e^{-4}$
Weight Decay	0.1
Betas	[0.9, 0.95]
Scheduler	CosineAnnealing
Warmup Steps	2000
Minimum Learning Rate	$3.0e^{-5}$

TABLE VIII

PERPLEXITY (LOWER IS BETTER) AND IBEROBENCH SCORES (HIGHER IS BETTER) FOR EACH LANGUAGE BETWEEN THE ORIGINAL CHECKPOINT FROM WHICH THE CPT WAS PERFORMED (No-CPT) AND CPT WITH UNIFORM (BSL-1), AD HOC (BSL-2), AND XDOGE WEIGHTS.

Language	Salamandra-2B					IberianLLM-7B		
	Perplexity		IberoBench			IberoBench (acc)		
	Bsl-1	XDoGE	No-CPT	Bsl-1	XDoGE	No-CPT	Bsl-2	XDoGE
CA	6.42	<b>6.37</b>	51.66	51.78	<b>52.69</b>	53.76	53.40	<b>54.77</b>
PT	11.27	<b>11.16</b>	<b>53.34</b>	52.87	52.20	48.87	<b>54.01</b>	52.94
EN	<b>10.80</b>	<b>10.80</b>	<b>46.49</b>	45.51	45.71	52.24	<b>53.08</b>	52.10
GL	7.79	<b>7.66</b>	29.54	29.66	<b>29.74</b>	34.31	<b>34.64</b>	34.28
ES	12.19	<b>12.15</b>	50.19	50.14	<b>50.85</b>	46.75	48.70	<b>50.46</b>
EU	5.21	<b>5.20</b>	38.47	39.35	<b>39.97</b>	39.76	41.03	<b>41.44</b>
Total Avg	8.95	<b>8.89</b>	44.95	44.89	<b>45.19<sup>a</sup></b>	45.95	47.48	<b>47.67<sup>b</sup></b>

<sup>a</sup>The differences with No-CPT and Bsl-1 are statistically significant at the 5% level.

<sup>b</sup>The difference with No-CPT is statistically significant at the 5% level.

and a significant decline after 40k-60k steps for uniform is seen for most languages.

### E. Evaluation and Analysis of IberianLLM

To position the IberianLLM-7b introduced in Section IV-B, we provide a summary of the IberoBench evaluations in Table IX, comparing it with the Salamandra-7b.

Our key findings are the following: (i) IberianLLM notably outperforms Salamandra in Belebele multilingual reading comprehension; (ii) trained with only  $0.53T$  unique tokens (vs. Salamandra’s  $2.6T$ ), IberianLLM achieves 95.58% of Salamandra’s average score on Iberian tasks. This suggests that our targeted training yields roughly  $4.9\times$  greater token

efficiency; (iii) IberianLLM shows competence in few-shot translation involving never-seen languages like German and French, achieving 76% of Salamandra’s performance on average despite the presumable absence of these languages and parallel data in the training corpus. This highlights the model’s strong cross-lingual capabilities.

## VI. DISCUSSION

The results presented in Table VIII and illustrated in Fig. 7 highlight the success of our measures to enhance linguistic inclusivity in LLMs. With two differently sized models capable of supporting various numbers of languages, we verified that



TABLE IX  
MODEL COMPARISON WITH RELATIVE IMPROVEMENTS ON IBEROBENCH. IBER. = IBERIANLLM-7B, SALA. = SALAMANDRA-7B. POSITIVE PERCENTAGES INDICATE WHEN IBERIANLLM-7B OVERCOMES SALAMANDRA-7B. ALL VALUES AVERAGED PER GROUP.

	Flores Iberian		Flores Non-Iber		Belebele		English		Catalan		Spanish		Basque		Galician		Portuguese	
	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.	Iber.	Sala.
Scores	23.63	25.42	18.42	24.23	28.04	25.91	43.84	43.46	50.99	52.70	35.37	36.03	37.36	39.87	26.76	27.12	58.80	63.78
$\Delta\%$	<b>-7.04%</b>		<b>-23.98%</b>		<b>+8.22%</b>		<b>+0.87%</b>		<b>-3.24%</b>		<b>-3.44%</b>		<b>-6.29%</b>		<b>-1.33%</b>		<b>-7.81%</b>	

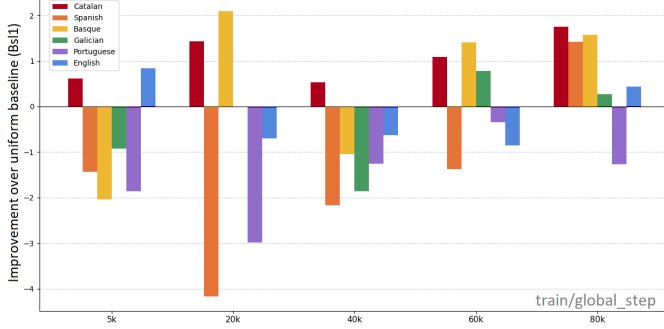


Fig. 7. Gains (in %) of the Salamandra-2b CPT with XDoGE-Set1 over Bsl1.

they can be improved in more than half of the target languages. This aligns with the original DoGE’s ambition to excel in more than half of the target domains [8]. Continual pretraining with XDoGE weights outperforms a uniform baseline (Bsl-1) that neglects language transferability, across five of six target languages (with a statistical significance at the 5% level). It is also not inferior to training with ad hoc (natural) weights (Bsl-2), on average across languages, while the advantage is the ability to steer the model towards highly weighted languages (Catalan and Basque in our case) and towards languages sharing commonalities (Spanish). Thus, if a specific combination of target languages makes proxy models converge to multiple options, we recommend choosing one that best satisfies needs in specific downstream use cases. We presume it is unlikely to excel in all languages at once by adjusting language sampling weights alone.

The results also offer several insights validating our design choices.

Firstly, our thresholded DoGE implementation reduces weight variance, which stabilizes proxy training and substantially decreases the loss compared to the original approach (Fig. 4). On the contrary, unconstrained proxy training leads to considerable discrepancies in weights, making them far unsuitable for real-world data volumes.

Secondly, applying the XDoGE framework to *pre-trained* models yields lower perplexity and higher scores on downstream tasks more consistently than with *base* models (cf. Tables IV, V and VIII), suggesting its better suitability for continual training, especially by delaying model degradation caused by data over-repetition.

Finally, the full-size training is not very sensitive to slight changes in weights overall (cf. scores for Set1 vs Set2 in Table IV). The exclusion in our case was for Portuguese, with a 10-point drop on downstream tasks. We assume this happened

due to a simultaneous decrease in the weights of the supporting Spanish and Galician languages.

## VII. CONCLUSION AND FUTURE WORK

We introduced XDoGE – an extension of the DoGE algorithm for multilingual pre-training that finds the optimized set of language weights that enhances cross-lingual transfer capability. In line with previous work, applying learnt weights for full-size training from scratch yields gains for only about half of the targets. Our proposal to use the weights in continual pre-training appears more beneficial, resulting in lower validation perplexity and improved downstream performance compared to the baselines for the majority of targets. We result in a series of LLMs of different sizes and origins, centred on Iberian languages and English.

For future work, we plan to explore adapting highly multilingual models to target languages while compromising less on non-target languages. Specifically, we will compare holistic language inclusion with a staged approach to changing target languages. We will also consider explicitly using language-family correspondences to improve the capture of interdependencies among languages.

## VIII. LIMITATIONS

While our approach demonstrates the potential of optimizing multilingual data distributions in a principled manner, there are several limitations that we consider when interpreting the findings and designing future research.

a) *Data availability as a contributing factor for weight estimation:* The XDoGE framework is originally meant to be efficient when estimating the weights with the proxy models, making the scale of the proxy models too small to see all available data. This is particularly relevant for low-resource languages, where data repetition is often necessary for effective pre-training. Since proxy models do not necessarily replicate the availability of data across languages, they may not fully capture the impact of data scarcity and repetition. While excluding data availability as a factor in estimating the language weights allows for clearer language interaction modeling, it does not address practical constraints in large-scale multilingual pre-training. Although our experiments have shown that controlled repetition can enhance language model learning in low-resource languages, future research should explore the impact of integrating data availability into the weight computation process. In particular, when computing weights for continual pre-training (CPT) experiments, it may be beneficial to train proxy models starting with the same weights used during pre-training instead of initializing them

with uniform weights, as in our current setup. This could enable more accurate estimation of the optimal data mixture for CPT experiments.

*b) Better efficiency in highly multilingual scenarios:* The computational cost of estimating  $W$  increases with both model size and the number of data sources. While the original DoGE framework was applied to only 7 sources for domain-weight estimation in a monolingual setting, we scaled this to 12 domains across 6 languages. However, extending this approach to highly multilingual language models remains a challenge. Future research with the XDoGE or similar frameworks should take into account the scalability of data size and sources while maintaining their benefits for generalization.

*c) Algorithm overhead clarity:* The optimization procedure that our method implies may not feel fully justified compared to the baselines that are straightforward to apply. Still, we must consider that we work with models with only up to 7 billion parameters: it is challenging to develop advanced language capabilities, such as reasoning, needed for downstream tasks; therefore, the improvement we observe is limited. We expect that scaling the conclusions to larger models would considerably save computational resources. In this case, the XDoGE proxy training time becomes truly negligible, clarifying the contribution of experimenting with smaller models.

## IX. ACKNOWLEDGEMENTS

This work was funded by the Ministerio para la Transformación Digital y de la Función Pública and Plan de Recuperación, Transformación y Resiliencia - Funded by EU – NextGenerationEU within the framework of the project Modelos del Lenguaje and the ILENIA Project with reference 2022/TL22/00215337, 2022/TL22/00215336, 2022/TL22/00215335, 2022/TL22/00215334. It was also funded by the Project Desarrollo de Modelos ALIA with the framework of the Plan Nacional de Tecnologías de Lenguaje -ENIA 2024 and PRTR, NextGeneration EU, Resol. SEDIA 19.08.2024. We also acknowledge the EuroHPC Joint Undertaking and the Spanish Supercomputing Network for awarding us access to MareNostrum5 and the technical support provided by the Barcelona Supercomputing Center (EHPC-EXT-2024E01-009, EHPC-AI-2024A05-048, RES-IM-2024-2-0031, RES-IM-2024-3-0021).

## REFERENCES

- [1] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian *et al.*, “The LLaMA 3 herd of models,” *arXiv e-prints*, arXiv:2407, 2024.
- [2] B. Adler, N. Agarwal, A. Aithal, D. H. Anh, P. Bhattacharya *et al.*, “Nemotron-4 340B technical report,” *CoRR*, abs/2406.11704, 2024.
- [3] A. Yang, B. Yang, B. Hui, B. Zheng, B. Yu *et al.*, “Qwen2 technical report,” *CoRR*, abs/2407.10671, 2024.
- [4] A. Young, B. Chen, C. Li, C. Huang, G. Zhang *et al.*, “Yi: Open foundation models by 01.AI,” *CoRR*, abs/2403.04652, 2024.
- [5] T. Le Scao, A. Fan, C. Akiki, E. Pavlick, S. Ilić *et al.*, “BLOOM: A 176B-parameter open-access multilingual language model,” *arXiv preprint* arXiv:2211.05100, 2023.
- [6] A. Q. Jiang, A. Sablayrolles, A. Roux, A. Mensch, B. Savary *et al.*, “Mixtral of experts,” *arXiv preprint* arXiv:2401.04088, 2024.
- [7] P. H. Martins, P. Fernandes, J. Alves, N. M. Guerreiro, R. Rei *et al.*, “EuroLLM: Multilingual language models for Europe,” *Procedia Computer Science*, vol. 255, pp. 53–62, 2025.
- [8] S. Fan, M. Pagliardini, and M. Jaggi, “DOGE: Domain reweighting with generalization estimation,” in *Proc. 41st Int. Conf. Machine Learning (ICML’24)*, pp. 12895–12915, 2024.
- [9] H. W. Chung, N. Constant, X. Garcia, A. Roberts, Y. Tay *et al.*, “Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining,” *arXiv preprint* arXiv:2304.09151, 2023.
- [10] H. Nigatu, A. Tonja, and J. Kalita, “The less the merrier? Investigating language representation in multilingual models,” in *Findings of the Assoc. Comput. Linguistics: EMNLP 2023*, pp. 12572–12589, 2023.
- [11] E. Gogoulou, T. Lesort, M. Boman, and J. Nivre, “Continual learning under language shift,” in *Text, speech, and dialogue: 27th Int. Conf. (TSD 2024)*, Brno, Czech Republic, 2024, pp. 71–84. Springer-Verlag.
- [12] H.-Y. Lee, S.-W. Li, and T. Vu, “Meta learning for natural language processing: A survey,” in *Proc. 2022 Conf. North American Chapter Assoc. Comput. Linguistics: Human Language Technologies*, pp. 666–684, 2022.
- [13] S. M. Xie, H. Pham, X. Dong, N. Du, H. Liu *et al.*, “DoReMi: Optimizing data mixtures speeds up language model pretraining,” in *Adv. Neural Information Processing Systems*, vol. 36, pp. 69798–69818, Curran Associates Inc., 2023.
- [14] Y. Oren, S. Sagawa, T. B. Hashimoto, and P. Liang, “Distributionally robust language modeling,” in *Proc. EMNLP-IJCNLP 2019*, pp. 4227–4237, 2019.
- [15] S. Sagawa, P. W. Koh, T. B. Hashimoto, and P. Liang, “Distributionally robust neural networks,” in *Int. Conf. Learning Representations*, 2020.
- [16] A. Üstün, V. Aryabumi, Z.-X. Yong, W. Ko, D. D’souza *et al.*, “Aya model: An instruction finetuned open-access multilingual language model,” *arXiv preprint* arXiv:2402.07827, 2024.
- [17] P. Joshi, S. Santy, A. Budhiraja, K. Bali, and M. Choudhury, “The state and fate of linguistic diversity and inclusion in the NLP world,” in *Proc. 58th Annu. Meeting Assoc. Comput. Linguistics*, pp. 6282–6293, 2020.
- [18] M. Brack, M. Ostendorff, P. O. Suarez, J. Saiz, I. Castilla, J. Palomar-Giner, A. Shvets, P. Schramowski, G. Rehm, M. Villegas and K. Kersting, “Community OSCAR: A community effort for multilingual web data,” in *Proc. 4th Workshop Multilingual Representation Learning (MRL 2024)*, pp. 232–235, 2024.
- [19] G. Penedo, H. Kydliček, A. Lozhkov, M. Mitchell, C. Raffel *et al.*, “The FineWeb datasets: Decanting the web for the finest text data at scale,” in *NeurIPS Datasets and Benchmarks Track*, 2024.
- [20] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi *et al.*, “LLaMA 2: Open foundation and fine-tuned chat models,” *arXiv preprint* arXiv:2307.09288, 2023.
- [21] A. Gonzalez-Agirre, M. Pàmies, J. Llop, I. Baucells, S. Da Dalt *et al.*, “Salamandra technical report,” *arXiv preprint* arXiv:2502.08489, 2025.
- [22] R. Sennrich, B. Haddow, and A. Birch, “Neural machine translation of rare words with subword units,” in *Proc. 54th Annu. Meeting Assoc. Comput. Linguistics (Vol. 1: Long Papers)*, pp. 1715–1725, 2016.
- [23] S. Da Dalt, J. Llop, I. Baucells, M. Pàmies, Y. Xu, A. Gonzalez-Agirre, and M. Villegas, “FLOR: On the effectiveness of language adaptation,” in *Proc. Joint Int. Conf. Comput. Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pp. 7377–7388, 2024.
- [24] I. Baucells, J. Aula-Blasco, I. de Dios-Flores, S. P. Suárez, N. Pérez, A. Salles, S. S. Docio, J. Falcão, J. J. Saiz, R. Sepúlveda-Torres and J. Barnes, “Iberobench: A benchmark for LLM evaluation in Iberian languages,” in *Proc. 31st Int. Conf. Comput. Linguistics*, pp. 10491–10519, 2025.
- [25] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, “Bleu: a method for automatic evaluation of machine translation,” in *Proc. 40th Annual Meeting of the ACL*, pp. 311–318, 2002.
- [26] C.-Y. Lin, “ROUGE: A package for automatic evaluation of summaries,” in *Text Summarization Branches Out*, pp. 74–81, 2004.
- [27] K. Micallef, A. Gatt, M. Tanti, L. van der Plas, and C. Borg, “Pre-training data quality and quantity for a low-resource language: New corpus and BERT models for Maltese,” in *Proc. 3rd Workshop Deep Learning for Low-Resource NLP*, pp. 90–101, 2022.
- [28] J. Hoffmann, S. Borgeaud, A. Mensch, E. Buchatskaya, T. Cai *et al.*, “Training compute-optimal large language models,” in *Proc. 36th Int. Conf. Neural Information Processing Systems (NeurIPS ’22)*, Curran Associates Inc., 2022.