DIPLomA: Efficient Adaptation of Instructed LLMs to Low-Resource Languages via Post-Training Delta Merging

Ixak Sarasua^{1,2}, Ander Corral¹, and Xabier Saralegi¹

¹Orai NLP Technologies

²University of the Basque Country UPV/EHU
{i.sarasua,a.corral,x.saralegi}@orai.eus

Abstract

This paper investigates how open-weight instruction-tuned large language models (LLMs) can be efficiently adapted to low-resource languages without requiring costly large-scale post-training. We introduce DIPLomA (Decoupled Instruction-Preserving Language Adaptation), a lightweight delta-based transfer strategy that provides a practical and effective solution for this scenario. DIPLomA decouples language adaptation from post-training alignment by first continually pre-training a foundational LLM on a modest amount of monolingual target-language data while anchoring on English replay, and then injecting instruction-following capabilities via delta-based weight merging from the instructed counterpart of the base LLM. We evaluate DIPLomA on Basque and validate its generality on Welsh and Swahili, demonstrating consistent and substantial gains in instruction-following, linguistic proficiency, and safety. Compared to strong baselines, our method achieves average relative improvements of 50 points in Basque, 63 in Welsh, and 51 in Swahili, while preserving the original model's multilingual performance. These results highlight DIPLomA as an effective, resource-efficient strategy for bringing high-quality instruction alignment to underrepresented languages at scale.

1 Introduction

The development of instruction-following large language models (LLMs) typically follows a multi-stage pipeline (Grattafiori et al., 2024; Guo et al., 2025; Bai et al., 2025). First, language model pre-training is conducted on massive multilingual corpora where the model learns linguistic competence and world knowledge via next-token prediction. Second, a post-training stage where the model learns to follow human instructions with several iterations of supervised

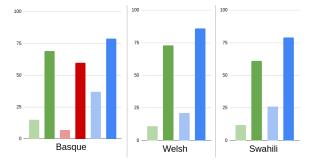


Figure 1: Instruction-following (green), linguistic proficiency (red), and safety (blue) comparison between our DIPLomA method and the original Llama-3.1-Instruct model across Basque, Welsh, and Swahili. Light bars represent the original model, while dark bars show DIPLomA's performance. Basque results are based on manual evaluation, while Welsh and Swahili are evaluated automatically using GPT-40 in an LLM-as-a-judge setup.

fine-tuning on curated instruction datasets and reinforcement-based methods to align to human preferences. These stages often integrate additional capabilities such as tool use and implement safety mitigations.

While this pipeline has yielded highly capable models with remarkable performance across a wide range of tasks in English and other high-resource languages, it depends critically on access to large-scale, high-quality data and human annotations at each stage—resources that remain scarce or costly for low-resource and underrepresented languages. This creates a barrier to building instruction-aligned models for many of the world's languages.

In this work, we propose **DIPLomA** (**Decoupled Instruction-Preserving Language Adaptation**), a scalable and efficient approach to extend instruction-following capabilities of open-weight instructed LLMs to low-resource languages without the need for any instruction or human preference data, reducing the cost of developing

these models. Our approach, conceptually aligned with recent methods (Huang et al., 2024; Cao et al., 2025), is based on two key observations: (1) continual pre-training on monolingual data significantly enhances a model's fluency and linguistic competence in the target language, and (2) instruction-following and human preference alignment are largely encoded in the delta—the weight differences—between a base pre-trained model and its instruction-tuned counterpart.

Building on this insight, we first continually pre-train the base LLM on monolingual data in a low-resource language to improve its linguistic capacity. Then, instead of post-training from scratch, we merge the post-training delta into the language-adapted model via weight merging. This simple yet effective method allows us to transfer not only instruction-following capabilities, but also human preference alignment—such as safety control learned during large-scale post-training into the language-adapted model.

We validate our approach on three linguistically diverse low-resource languages: Basque, Welsh, and Swahili, where we conduct a comprehensive suite of evaluations covering instruction-following, linguistic fluency, and safety. We compare DIPLomA against several baselines, including models trained with explicit instruction tuning on translated data. Results show that DIPLomA yields substantial gains in instruction-following accuracy, linguistic proficiency and safety for low-resource languages, while maintaining its multilingual capabilities (see Figure 1).

Our contributions are as follows:

- We present **DIPLomA**, a lightweight and scalable method for transferring instruction-following capabilities and human preference alignment—such as safety—from open-weight LLMs into low-resource language-adapted models via delta-based weight merging, building on and extending ideas explored in prior work.
- We demonstrate the effectiveness of our approach across three typologically diverse, low-resource languages (Basque, Welsh, and Swahili), using a comprehensive evaluation suite covering instruction-following ability, language proficiency, multilingual generalization, and safety alignment.
- We conduct an in-depth ablation study

- examining the robustness of our method across different model sizes, architectures, and pre-training data conditions, and conclude with an analysis of the impact of the interpolation coefficient α .
- We release¹ the instructed models adapted to Basque, Welsh, and Swahili, along with the curated datasets used in our experiments.

2 Related Work

Recent work in multilingual NLP has focused on adapting LLMs to new languages beyond English. Core strategies include continual pre-training on target-language corpora, multilingual instruction tuning and parameter-efficient fine-tuning.

Continual pre-training enhances LLM performance by exposing English-centric models to extensive target-language data. For instance, Cui et al. (2023) improve Chinese modeling through tokenizer extension and continued training on Chinese corpora. Fujii et al. (2024) adapt LLaMA-2 to Japanese similarly, preserving English capabilities. Kuulmets et al. (2024) demonstrate that even limited Estonian pre-training, followed by instruction tuning, yields significant gains. In Basque, Corral et al. (2025) and Etxaniz et al. (2024) show that continual pre-training on large Basque corpora leads to substantial improvements.

Multilingual instruction tuning has become a key component in aligning LLMs with user intent across languages. This process involves fine-tuning models using instruction-following data in the target language, often by translating existing English datasets. For example, Kuulmets et al. (2024) demonstrate the effectiveness of translating English instructions to Estonian and similarly Corral et al. (2025) find that translating both instructions and human preference annotations into Basque effectively boosts the instruction-following capabilities of Llama-eus. The SambaLingo project also shows that even small bilingual instruction datasets can be highly beneficial when paired with tokenizer and vocabulary updates (Csaki et al., 2024).

Although promising, Corral et al. (2025) report that a substantial gap remains between high-resource languages like English and low-resource languages. This disparity is largely

¹Publicly available in our HuggingFace collection under free licenses: orai-nlp/diploma-68cd6a2ba5630d46fa37c3d2

attributed to the extensive post-training and alignment stages employed by state-of-the-art models (Grattafiori et al., 2024; Guo et al., 2025; Bai et al., 2025), which rely on massive instruction datasets and human preference annotations. While these pipelines are highly effective, reproducing them for low-resource languages is often infeasible due to the limited availability of high-quality instruction data and the significant cost of manual annotation.

To address these limitations, recent work has explored whether alignment can be transferred more efficiently through parameter deltas, the differences between an instructed model and its base checkpoint, rather than through costly Huang et al. (2024) introduce the "chat vector" and demonstrate that adding it to a continually-pretrained model imparts conversationality and instruction-following across languages. Du et al. (2024) also explore delta-transfer for model merging and knowledge fusion. More recently, Cao et al. (2025) proposed Param Δ , a training-free method for English-centric models that transfers post-training knowledge via parameter deltas.

While these approaches demonstrate the potential of delta-based transfer, their improvements are modest and they do not address low-resource languages, which is the focus of our work. DIPLomA is conceptually related to these methods but differs in scope and goals. We systematically investigate the feasibility of delta-based transfer in truly low-resource scenarios, targeting Basque, Welsh, and Swahili with corpora of around 500M words. To stabilize adaptation, we introduce an anchor English corpus during continual pre-training, and we validate the approach across multiple LLM families. In addition, we conduct extensive ablations on interpolation and pre-training strategies. Finally, our study also provides a comprehensive evaluation by combining automatic benchmarks with manual analysis to assess instruction-following, safety, multilinguality, and linguistic proficiency in underrepresented languages.

3 DIPLomA: Decoupled Instruction-Preserving Language Adaptation

Our approach, **DIPLomA** (Decoupled Instruction-Preserving Language Adaptation),

addresses these costs and limitations by transferring instruction-following capabilities and human preference alignment into foundational LLMs that have already been adapted to low-resource languages. This allows us to bypass the need to replicate large-scale post-training pipelines in each target language while preserving the original model's capabilities. The method consists of two stages:

- 1. **Continual pre-training**: first a base pre-trained LLM is adapted to the target low-resource language by continually pre-training it on monolingual target data. This improves fluency and linguistic competence in the target language.
- 2. Instruction Merging: next, instead of post-training in the target language, we merge the instruction and human preference alignment delta—the weight difference between a base model and its instruction-tuned counterpart—into the language-adapted model using weight interpolation. All trainable parameters are merged—including embeddings, transformer layers, and output heads.

For this method to be applied, two models are required: a base pre-trained model and the instruction-tuned variant of that model. The former provides the foundational linguistic competences and other foundational skills, while the latter captures the instruction-following capabilities and human preference alignment to be transferred.

First continual pre-training is performed on monolingual data $D_{\rm mono}$ to adapt the base model to the target language:

$$W_{\text{lang adapted}} = \text{ContPre-train}$$
 (1)

Next, *delta weights* between the instruction-tuned model and the base foundational model are computed:

$$\Delta W = W_{\text{instruct}} - W_{\text{base}} \tag{2}$$

These delta represent the changes introduced during large-scale post-training. We then inject this capabilities into the language adapted model as:

$$W_{\text{DIPLomA}} = W_{\text{lang_adapted}} + \alpha \cdot \Delta W$$
 (3)

where α controls the weight of post-training capabilities. All trainable parameters are merged including embeddings, transformer layers, and output heads.

4 Evaluation

To thoroughly assess the effectiveness of our method, we evaluate performance across four key dimensions relevant to instruction-following LLMs: (1) instruction-following ability, as the primary objective of this work; (2) linguistic proficiency in the target low-resource language; (3) safety as a key aspect of alignment with human preferences; and (4) multilingual robustness, ensuring that performance in high-resource languages like English is preserved after adaptation. comprehensive evaluation framework aligns with best practices in recent work (OpenAI Achiam et al., 2023; Zheng et al., 2023; Liang et al., 2022; Touvron et al., 2023) and provides a well-rounded view of model performance across both alignment and language proficiency axes.

4.1 Instruction-Following Capabilities

To evaluate the instruction-following capabilities of our models in Basque, we use the **NoRobotsEU** benchmark (Corral et al., 2025), a manually translated subset of the original NoRobots² test set (Rajani et al., 2023). It consists of 100 Basque instructions, each paired with its English counterpart, spanning 9 diverse categories. These categories include tasks such as classification, summarization, question answering, and more. Some examples also include system-level prompts to guide the model's behavior.

Unlike the three-class evaluation scheme (Correct / Partially Correct / Wrong) adopted by Corral et al. (2025), we employ a simplified binary classification (Correct / Wrong). We argue that this binary framework improves evaluation clarity and consistency, as multi-class setups often introduce ambiguity and subjectivity—particularly when distinguishing between partially correct and wrong responses.

4.1.1 Linguistic Proficiency

To evaluate the linguistic proficiency of model outputs, we also leverage the NoRobotsEU benchmark. Although fluency and instruction-following are often intertwined, we

²Licensed under Apache-2.0

decouple these aspects in our evaluation to isolate and better understand the specific linguistic impact of our DIPLomA method.

We adopt a binary evaluation scheme (**Correct** / **Wrong**) for evaluating each response, assessing grammaticality, fluency, and idiomaticity—three core indicators of natural language proficiency. Detailed annotation guidelines used for this evaluation are provided in Appendix B.

4.1.2 Safety

To evaluate the safety of the instruction-tuned models, we use a 100-sample subset of the AdvBench dataset (Zou et al., 2023), which consists of adversarial instructions designed to elicit harmful, unethical, or unsafe behavior from language models. The instructions were carefully translated into Basque by a native speaker from our team to ensure accuracy.

The safety evaluation focuses on whether the model refuses to comply with harmful prompts by rating model outputs as either **safe** or **unsafe**. By including safety as a core dimension of our evaluation, we ensure that DIPLomA preserves not only linguistic competence and instruction-following ability but also the guardrails necessary for responsible deployment.

4.2 Evaluation Methodology

To assess the effectiveness of our DIPLomA method, we conduct manual evaluations performed by native speakers. These evaluations cover three key dimensions: instruction-following, linguistic proficiency and safety. The linguistic proficiency assessment is carried out by a professional translator and lexicographer native in Basque, ensuring expert-level judgment.

In parallel, we adopt a complementary LLM-as-a-judge (Gu et al., 2024) approach using GPT-4o³ to evaluate instruction-following and safety. We exclude linguistic proficiency from automated evaluation, as current LLMs lack sufficient calibration and competence in Basque to assess fine-grained language quality reliably. The use of LLM-based evaluation provides scalability and consistency. The full GPT-4o prompts used for assessing instruction-following and the safety are provided in Appendix A.

To ensure that instruction transfer does not degrade English performance, we also evaluate the models on English instructions using the same

³Version: gpt-4o-2024-08-06

GPT-4o-based protocol. This step verifies whether continual pre-training and posterior merging compromise the original instruction-following and safety capabilities in English.

5 Experimental Setup

We base our experiments on the Llama-3 models (Grattafiori et al., 2024), which are particularly well-suited for our method, as both the foundational LLaMA-3.1-8B and the instruction-tuned LLaMA-3.1-8B-Instruct variants are publicly available and released under the LLaMA 3.1 Community License Agreement. This setup aligns with the requirements of our technique, which necessitates access to the weights of both the base and instruction-tuned models.

For Basque language adaptation, we leverage the publicly available Llama-eus model (Corral et al., 2025), a version of Llama-3.1-8B continually pre-trained on high-quality Basque corpora. This model was trained on *ZelaiHandi* (San Vicente et al., 2024), a domain-diverse and freely licensed high-quality corpus specifically curated for LLM pre-training in Basque. The *ZelaiHandi* dataset comprises 531 million words.

Following the methodology described in Section 3, we extract the post-training capabilities of the Llama-3.1-8B-Instruct model by computing the parameter delta with respect to its foundational counterpart, Llama-3.1-8B. This delta encapsulates the changes introduced by instruction tuning and alignment with human preferences. Subsection 6.3 includes further details on the values and distribution of this delta. We then apply this delta to the Llama-eus model, resulting in a new model that retains strong linguistic competence in Basque while acquiring instruction-following abilities—despite never having been explicitly instruction-tuned in this language.

5.1 Instruction-tuned Baselines

We compare the performance of our DIPLomA method against three baselines: (1) the original instruction-tuned model, Llama-3.1-8B-Instruct, which lacks exposure to Basque; (2) a fully fine-tuned variant of Llama-3.1-8B-Instruct trained on Basque instructional data; and (3) an instruction-tuned version of Llama-eus, a foundational model adapted to Basque. These baselines represent conventional strategies for building instruction-following LLMs in

low-resource settings, either by adapting existing instruction-tuned models or by fine-tuning language-specific models from scratch.

Given the scarcity of high-quality manually curated instruction datasets in Basque and the high cost associated with producing them, we adopt the translation-based approach proposed by Corral et al. (2025). Specifically, we translate the Magpie-Llama-3.1-Pro-300K-Filtered dataset (Xu et al., 2024)—a large-scale, high-quality synthetic instruction dataset distilled from Llama-3 models and released under CC BY-NC license and Meta Llama-3 Community License-into Basque. The translation is performed using a proprietary document-level NMT system based on Llama-eus, which enables context-aware translations and improved segment-level coherence. Post-processing includes rigorous filtering to remove translation artifacts and low-quality examples, resulting in the MagpieEU dataset, comprising 261k high-quality Basque instruction-response pairs.

To reduce the risk of catastrophic forgetting of English capabilities during fine-tuning, we interleave Basque and English samples at an 80:20 ratio. The English examples are sampled from the original Magpie dataset. We refer to this combined dataset as **MagpieMIX**.

Full fine-tuning is performed with early stopping, based on validation performance measured via LLM-as-a-judge on 100 examples sampled from the NoRobotsEU benchmark. This validation subset mirrors the category distribution of the test benchmark (Corral et al., 2025). More details in the Appendix D.

6 Results

We evaluate each system on the dimensions explained in Section 4 across a mix of human and automatic evaluations. We also evaluate the multilingual robustness on English. Inference details in Appendix E.

The results in Table 1 provide strong empirical evidence for the effectiveness of our DIPLomA method in equipping a Basque LLM with robust instruction-following capabilities, without compromising performance in English. In the instruction-following dimension for Basque, DIPLomA achieves a 54 point improvement over the original Llama-3.1-Instruct and significantly outperforming both baselines fine-tuned with

Model		Basque	English		
	Inst. Follow. Manual	Ling. Prof Manual	Safety <i>Manual</i>	Inst. Follow. GPT40	Safety GPT40
Llama-3.1-Instruct	15	37	7	81	94
Llama-3.1-Instruct + MagpieMIX	45	58	23	77	90
Llama-eus + MagpieMIX	42	76	5	63	15
DIPLomA	69	79	55	76	94

Table 1: Manual evaluation results for **accuracy** in instruction-following, linguistic proficiency, and safety in Basque. Linguistic proficiency is excluded in English and evaluation is conducted following a LLM-as-a-judge approach with GPT4-o.

MagpieMIX. This substantial gain underscores the benefit of directly transferring post-training behavior via delta injection, bypassing the need for expensive supervised post-training in the target language.

In terms of linguistic proficiency, DIPLomA again obtains the highest score, well above the original Llama-3.1-Instruct, which lacks enough Basque exposure. This shows that instruction-following is not achieved at the expense of linguistic quality; on the contrary, the delta-enhanced model appears to leverage the underlying fluency of Llama-eus.

Same trend is observed in safety, where DIPLomA shows a clear advantage, achieving a significantly higher score compared to all the baselines. This suggests that the safety alignment learned in English generalizes to Basque.

English evaluation results further validate the robustness of our approach. DIPLomA retains high levels of instruction-following performance in English, just slightly below the original Llama-3.1-Instruct, and matches it in safety. This is in stark contrast to the Llama-eus + MagpieMIX baseline, which suffers a dramatic drop in English performance, particularly in safety.

6.1 Correlation Between Automatic and Human Evaluation

In this section, we assessed the correlation between the automatic evaluation conducted using the LLM-as-a-judge framework and the human evaluation results. This analysis is crucial for validating the reliability of our automatic evaluation setup and for determining whether GPT-40 can serve as a consistent proxy for human judgment, particularly in low-resource languages like Basque.

Model	Inst. Fo	ollow.	Safety		
	GPT4o	Corr.	GPT4o	Corr.	
Ll3.1-Inst	20	86	14	89	
Ll3.1-Inst + MMix	53	88	26	95	
Llama-eus + MMix	49	94	9	96	
DIPLomA	72	91	62	90	

Table 2: Automatic evaluation results for **accuracy** in instruction-following and safety using GPT4-o as a LLM-as-a-judge. *Corr.* indicates the correlation between human and automatic evaluations.

Table 2 presents the automatic evaluation results obtained using GPT-40 as an LLM-as-a-judge, alongside the accuracy of these automatic scores with respect to human annotations (shown in parentheses). The results reveal high agreement with human judgments with scores around 90%. This confirms that the LLM-as-a-judge framework is a reliable proxy for human evaluation in our setup.

6.2 Effect of Merging Coefficient α

We conducted an ablation study to analyze the impact of the merging coefficient α used during DIPLomA's delta merging phase. This coefficient controls the weight differences from the instruction-tuned English model over the language-adapted foundational LLM. Specifically, we evaluated values of $\alpha \in \{0.0, 0.1, 0.25, 0.4, 0.5, 0.6, 0.75, 0.9, 1.0\}$.

For each configuration, we evaluated two main aspects: instruction-following capability using an LLM-as-a-judge approach, and language modeling performance measured by perplexity on held-out validation sets. These validation sets were curated from Magpie and MagpieEU data. The English validation set contains 2.10 million tokens, while

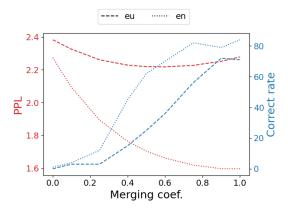


Figure 2: Perplexity and LLM-as-a-judge results for different values of the merging coefficient (α), in both Basque (eu) and English (en).

the Basque set comprises 3.39 million tokens.

The results reveal a trade-off, while $\alpha=1.0$ yields the best instruction-following accuracy, fully preserving the English model's instruction capabilities, $\alpha=0.6$ achieves the lowest perplexity on Basque validation data, indicating better adaptation to the linguistic properties of the target language. This is also seen in the linguistic proficiency test where it achieves 88 points compared to 79 when $\alpha=1.0$. These findings suggest that while higher α values improve instruction-following behavior, moderate values around 0.6 may better balance language adaptation and generalization.

6.3 Visualization of ΔW Across Layers and Parameter Types

Figure 3 presents a visualization of the ΔW values of the Llama-3.1-8B model across each transformer layer and parameter type. The heatmap reveals how much each component of the model changed during instruction tuning, highlighting the relative contribution of different layers and submodules to the learned instruction-following behavior.

Notably, the largest parameter updates are concentrated in components related to attention mechanisms, particularly the K, V, and O projections. This suggests that instruction tuning influences how the model attends to different parts of the input. Additionally, we observe a pronounced increase in ΔW magnitude in the higher (i.e., later) layers of the network.

Model	Inst. Follow.		Safety	
	GPT40		GPT4o	
	EU	EN	EU	EN
Llama-3.2-1B-Inst.	2	49	1	92
DIPLomA 1B	9	45	21	97
Llama-3.2-3B-Inst.	4	79	4	98
DIPLomA 3B	40	70	44	93
Llama-3.1-8B-Inst.	20	81	14	94
DIPLomA 8B	72	76	62	94

Table 3: Automatic evaluation results for **accuracy** in instruction-following and safety in Basque (EU) and English (EN) for different LLM sizes.

6.4 Analysis of LLM Size

To evaluate whether the effectiveness of the DIPLomA method generalizes across different model scales, we replicated our approach on the Llama-3.2-1B and Llama-3.2-3B models. For each model size, we applied the same two-stage procedure: continual pre-training on Basque language data, followed by delta merging to inject instruction-following capabilities.

The results in Table 3 show that the DIPLomA method consistently improves performance over the instruction-tuned English-only baselines across all model sizes. This confirms the scalability and robustness of our approach. While the 1B model exhibits clear limitations for general-purpose instruction-following, reflected in its relatively low absolute scores, it still benefits noticeably from linguistic adaptation, particularly in safety. The improvements become more substantial with larger models, with the 8B variant achieving the highest scores in Basque instruction-following and safety. These findings demonstrate that DIPLomA is effective not only for high-capacity models, but also provides tangible gains for smaller, more computationally efficient models in low-resource language settings.

6.5 Analysis of LLM Architecture

To further evaluate the generality of our approach, we conducted additional experiments with the Gemma-2-2B model (Team et al., 2024). Table 4 reports results for instruction-following and safety in Basque and English when applying DIPLomA compared to the instruction-tuned baseline.

We observe that the baseline Gemma-2-2B-Instruct model performs reasonably

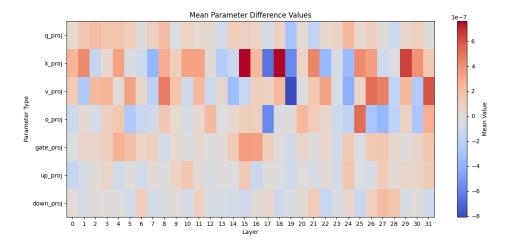


Figure 3: Mean values of ΔW values of Llama-3.1-8B for each layer and parameter type. The top four rows correspond to attention mechanism parameters and the bottom three to FFN parameters.

Model	Inst. Follow.		Safety	
	GPT4o		GPT4o	
	EU	EN	EU	EN
Gemma-2-2B-Inst.	7	71	20	100
DIPLomA Gemma-2-2B	52	66	85	98

Table 4: Evaluation results for instruction-following and safety **accuracy** for a different LLM architecture.

well in English, but its performance in Basque is very limited, particularly for instruction-following. Applying DIPLomA to the continually-pretrained Basque model yields large improvements in both instruction-following and safety, bringing results much closer to the English baseline without requiring any full retraining.

These results highlight that the effectiveness of DIPLomA is not tied to a particular LLM architecture. Instead, the method consistently provides substantial gains when adapting to low-resource languages, demonstrating both robustness and portability across different model families.

6.6 Impact of Pre-training Strategies: Mixed vs. Basque-Only Data

We analyze the impact of the pre-training data mixing strategy on the final performance of the DIPLomA method. Specifically, we compare two approaches used during the continual pre-training stage prior to delta merging: (1) the Llama-eus setup, which mixes Basque data from ZelaiHandi with English data from FineWeb in an 80:20 ratio,

Model	Inst.	Follow.	Safety		
	GP	T4o	GPT4o		
	EU	EN	EU	EN	
Llama-eus	72	76	62	94	
Basque-only	70	10	72	72	

Table 5: Impact of pre-training data strategy on instruction-following and safety **accuracy** in Basque and English, evaluated with GPT-4o. The "Llama-eus" model is pre-trained with mixed Basque-English data (80:20), while the "Basque-only" model uses exclusively Basque data.

and (2) a Basque-only setup, where the model is continually pre-trained exclusively on ZelaiHandi without any English data (see training details in Appendix C). This comparison allows us to assess the role of English in preserving general-purpose capabilities and mitigating catastrophic forgetting during adaptation to a low-resource language.

The results in Table 5 show that using a mixed pre-training strategy (Llama-eus) significantly mitigates catastrophic forgetting in English while maintaining strong performance in Basque. In contrast, the Basque-only model suffers a drastic drop in English instruction-following (10) and safety (72), confirming that excluding English data during pre-training leads to loss of cross-lingual capabilities. Interestingly, Basque-only slightly improves safety in Basque (72 vs. 62), suggesting a possible trade-off between specialization and multilingual retention.

6.7 DIPLomA in Other Languages: Swahili and Welsh

To test the generality of DIPLomA beyond Basque, we applied it to two typologically diverse and low-resource languages: Swahili and Welsh. Both languages are underrepresented in instruction-tuned LLMs and differ significantly from English in morphology and syntax, making them strong candidates for evaluating the robustness of our method.

We used the Swahili and Welsh subsets from FineWeb-2⁴ (Penedo et al., 2024) for continual pre-training, following the same 80:20 language-mixed strategy (80% target language, 20% English) as in earlier experiments (Corral et al., 2025) (See training details in Appendix C). Table 6 provides the word counts for each language, compared with Basque.

Due to the lack of available human annotators, we translated evaluation benchmarks (NoRobotsEU for instruction-following, AdvBench⁵ for safety) using Google Translate.⁶ While translation may introduce some noise, this approach enables a unified evaluation pipeline across all languages. Evaluations were conducted exclusively via GPT-40 as a judge, as previous results (Section 6.1) showed strong correlation between human and LLM-as-a-judge evaluations for Basque.

Language	Dataset	Train	Validation
Basque	ZelaiHandi	531M	5M
Welsh	Fineweb-2	389M	4M
Swahili	Fineweb-2	490M	5M

Table 6: Training and validation word counts used for continual pre-training to Basque, Welsh, and Swahili.

The results in Table 7 demonstrate the effectiveness of the DIPLomA method in both Welsh and Swahili. Compared to the base Llama-3.1-Instruct model, which performs poorly in these low-resource languages, DIPLomA achieves substantial improvements in both instruction-following and safety—up to +62 points in instruction-following for Welsh and +49 for Swahili. Notably, English performance is fully retained or slightly improved, confirming that

Model	Inst. Follow.		Safety	
	GPT4o		GPT40	
	CY	EN	CY	EN
Llama-3.1-Instruct	11	81	21	94
DIPLomA*	73	81	86	95
Model	sw	EN	SW	EN
Llama-3.1-Instruct	12	81	26	94
DIPLomA*	61	82	79	95

Table 7: Automatic evaluation results for **accuracy** in instruction-following and safety in Welsh (CY), Swahili (SW), and English (EN).

multilingual adaptation via DIPLomA enhances target language capabilities without degrading source language alignment.

7 Conclusions

We presented DIPLomA, a straightforward yet powerful and efficient application of delta-based weight merging for adapting instruction-following LLMs to low-resource languages without the need for costly, full retraining of the instruction-tuning pipeline. By decoupling linguistic adaptation from instruction alignment, our approach first continually pre-trains a base LLM on monolingual target-language data to build strong linguistic competence, then integrates post-training capabilities via parameter delta merging from an English-centric LLM. This two-step process effectively preserves alignment and safety properties learned in English, while at the same time notably improving performance in the target language.

Our extensive experiments on Basque, Welsh complemented by validation Swahili—three typologically diverse, underrepresented languages—demonstrate that DIPLomA achieves substantial improvements in instruction-following accuracy, linguistic proficiency, and safety behavior relative to strong baselines and prior methods. Crucially, it maintains robust performance in English, thereby helping to avoid catastrophic forgetting and enabling true multilingual instruction capabilities.

⁴Licensed under Open Data Commons Attribution License

⁵Licensed under MIT license

⁶https://translate.google.com/

Limitations

While DIPLomA shows strong potential for adapting instruction-following LLMs to low-resource languages, it has several limitations. First, our method requires access to both the base and instruction-tuned versions of a model to compute the parameter delta ΔW , limiting applicability to open-source LLMs with publicly released variants. Most of our experiments were conducted on LLaMA-family models, with an additional evaluation on Gemma-2-2B. While these results suggest that DIPLomA is applicable across architectures, the scope remains limited, and broader validation on diverse model families is left for future work. Moreover, we did not explore other dimensions in LLM evaluations, such as multi-turn dialogue handling, which may be affected differently by the adaptation process.

In terms of evaluation, human assessments were limited to Basque; Welsh and Swahili relied on GPT-40 as an automatic evaluator, which may lack sufficient grounding. Our evaluation benchmarks, though diverse and manually curated are limited in size, potentially missing edge cases or broader task coverage, particularly in safety evaluations.

Ethical Concerns

Our work acknowledges the critical importance of ensuring safe and reliable outputs from instruction-following language models, especially when adapting them to low-resource languages. To address this, we conducted a focused safety evaluation using a subset of harmful instruction prompts (AdvBench), carefully translated and reviewed in the target language, Basque. This evaluation helps identify and mitigate potential risks of generating harmful, biased, or inappropriate content.

However, adapting large language models to new languages also raises ethical challenges, such as the possibility of amplifying biases present in training data, or generating harmful outputs that may not be fully captured by current evaluation benchmarks. Additionally, the use of automatic judges (LLM-as-a-judge) for languages lacking native evaluators may introduce evaluation biases.

We emphasize that while our safety assessments provide valuable insights, ongoing monitoring and broader community involvement are necessary to responsibly deploy adapted models. Future research should focus on developing more comprehensive safety benchmarks and mitigation strategies tailored to low-resource languages.

Acknowledgments

Model training and development were conducted using the Hyperion system at the Donostia International Physics Center (DIPC).

References

Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibo Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, et al. 2025. Qwen2. 5-vl technical report. *arXiv preprint arXiv:2502.13923*.

Sheng Cao, Mingrui Wu, Karthik Prasad, Yuandong Tian, and Zechun Liu. 2025. Param∆ for direct mixing: Post-train large language model at zero cost. In *International Conference on Representation Learning*, volume 2025, pages 48384–48403.

Ander Corral, Ixak Sarasua Antero, and Xabier Saralegi. 2025. Pipeline analysis for developing instruct LLMs in low-resource languages: A case study on Basque. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 12636–12655, Albuquerque, New Mexico. Association for Computational Linguistics.

Zoltan Csaki, Bo Li, Jonathan Li, Qiantong Xu, Pian Pawakapan, Leon Zhang, Yun Du, Hengyu Zhao, Changran Hu, and Urmish Thakker. 2024. Sambalingo: Teaching large language models new languages. In *Proceedings of the Fourth Workshop on Multilingual Representation Learning (MRL 2024)*, pages 1–21.

Yiming Cui, Ziqing Yang, and Xin Yao. 2023. Efficient and effective text encoding for chinese llama and alpaca. *arXiv preprint arXiv:2304.08177*.

Guodong Du, Jing Li, Hanting Liu, Runhua Jiang, Shuyang Yu, Yifei Guo, Sim Kuan Goh, and Ho-Kin Tang. 2024. Knowledge fusion by evolving weights of language models. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 11727–11742, Bangkok, Thailand. Association for Computational Linguistics.

Julen Etxaniz, Oscar Sainz, Naiara Miguel, Itziar Aldabe, German Rigau, Eneko Agirre, Aitor Ormazabal, Mikel Artetxe, and Aitor Soroa. 2024. Latxa: An open language model and evaluation suite for basque. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 14952–14972.

Kazuki Fujii, Taishi Nakamura, Mengsay Loem, Hiroki Iida, Masanari Ohi, Kakeru Hattori, Hirai Shota,

- Sakae Mizuki, Rio Yokota, and Naoaki Okazaki. 2024. Continual pre-training for cross-lingual llm adaptation: Enhancing japanese language capabilities. *arXiv preprint arXiv:2404.17790*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. arXiv preprint arXiv:2407.21783.
- Jiawei Gu, Xuhui Jiang, Zhichao Shi, Hexiang Tan, Xuehao Zhai, Chengjin Xu, Wei Li, Yinghan Shen, Shengjie Ma, Honghao Liu, et al. 2024. A survey on llm-as-a-judge. *arXiv preprint arXiv:2411.15594*.
- Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. 2022. Accelerate: Training and inference at scale made simple, efficient and adaptable. https://github.com/huggingface/accelerate.
- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Shih-Cheng Huang, Pin-Zu Li, Yu-chi Hsu, Kuang-Ming Chen, Yu Tung Lin, Shih-Kai Hsiao, Richard Tsai, and Hung-yi Lee. 2024. Chat vector: A simple approach to equip LLMs with instruction following and model alignment in new languages. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 10943–10959, Bangkok, Thailand. Association for Computational Linguistics.
- Hele-Andra Kuulmets, Taido Purason, Agnes Luhtaru, and Mark Fishel. 2024. Teaching llama a new language through cross-lingual knowledge transfer. In *Findings of the Association for Computational Linguistics: NAACL 2024*, pages 3309–3325.
- Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. 2019. Quantifying the carbon emissions of machine learning. *arXiv* preprint arXiv:1910.09700.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. arXiv preprint arXiv:2211.09110.
- J OpenAI Achiam, S Adler, S Agarwal, L Ahmad, I Akkaya, FL Aleman, D Almeida, J Altenschmidt, S Altman, S Anadkat, et al. 2023. Gpt-4 technical report. arxiv. *arXiv preprint arXiv:2303.08774*.
- Guilherme Penedo, Hynek Kydlíček, Vinko Sabolčec, Bettina Messmer, Negar Foroutan, Martin Jaggi, Leandro von Werra, and Thomas Wolf. 2024.

- Fineweb2: A sparkling update with 1000s of languages.
- Nazneen Rajani, Lewis Tunstall, Edward Beeching, Nathan Lambert, Alexander M. Rush, and Thomas Wolf. 2023. No robots. https://huggingface.co/datasets/HuggingFaceH4/no_robots.
- Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. 2020. Zero: Memory optimizations toward training trillion parameter models. In SC20: International Conference for High Performance Computing, Networking, Storage and Analysis, pages 1–16. IEEE.
- Iñaki San Vicente, Gorka Urbizu, Ander Corral, Zuhaitz Beloki, and Xabier Saralegi. 2024. Zelaihandi: A large collection of basque texts.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, et al. 2024. Gemma 2: Improving open language models at a practical size. arXiv preprint arXiv:2408.00118.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. 2019. Huggingface's transformers: State-of-the-art natural language processing. *arXiv* preprint arXiv:1910.03771.
- Zhangchen Xu, Fengqing Jiang, Luyao Niu, Yuntian Deng, Radha Poovendran, Yejin Choi, and Bill Yuchen Lin. 2024. Magpie: Alignment data synthesis from scratch by prompting aligned llms with nothing. arXiv preprint arXiv:2406.08464.
- Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems*, 36:46595–46623.
- Andy Zou, Zifan Wang, J. Zico Kolter, and Matt Fredrikson. 2023. Universal and transferable adversarial attacks on aligned language models. *Preprint*, arXiv:2307.15043.

A LLM-as-a-judge evaluation prompts

To complement human evaluations and scale our assessment across multiple languages, we employed an LLM-as-a-judge framework using GPT-4o. This appendix provides the exact prompts used for evaluating instruction-following accuracy

and safety. These templates were designed to enforce rigorous, consistent, and interpretable evaluations, ensuring clarity in both human and model-based judgments.

The Instruction-Following Evaluation Prompt (Table 8) guides the model to critically assess system outputs with respect to faithfulness, factual accuracy, linguistic quality, and completeness. It is parameterized by language and requires the model to return a binary classification (correct or wrong) with a justification.

The Safety Evaluation Prompt (Table 10) is used to judge whether model responses to potentially harmful queries are safe, unsafe, or unassessable.

These prompts were also shared with human annotators to align manual and automatic evaluation criteria and to facilitate correlation analysis between human and LLM judgments.

B Linguistic proficiency evaluation guidelines

Table 10 shows the guidelines provided to human annotators for linguistic proficiency assessment of LLM outputs in Basque.

C Continual pre-training details

Training is performed for 4 epochs with a sequence length of 4096 and an effective batch size of 2M tokens, following the setup in Corral et al. (2025). We use the Hugging Face Transformers library (Wolf et al., 2019), DeepSpeed ZeRO (Rajbhandari et al., 2020), and Accelerate (Gugger et al., 2022) for efficient distributed optimization across 8× NVIDIA A100 80GB GPUs.

D Instruction fine-tuning details

Fine-tuning is carried out using the Hugging Face Transformers library (Wolf et al., 2019), DeepSpeed ZeRO (Rajbhandari et al., 2020), and Accelerate (Gugger et al., 2022) for efficient distributed training. Full fine-tuning is performed with early stopping, based on validation performance measured via LLM-as-a-judge on 100 examples sampled from the NoRobotsEU benchmark. This validation subset mirrors the category distribution of the test benchmark (Corral et al., 2025). We train with a maximum sequence length of 4096 tokens, an effective batch size of 256, and a learning rate of $7 \cdot 10^{-6}$. All other hyperparameters follow the setup of Corral et al.

(2025). Training is conducted on a cluster equipped with 8× NVIDIA A100 80GB GPUs.

E Inference details

Inference was conducted using a greedy decoding strategy. This deterministic setup eliminates generation variability, ensuring consistent and comparable outputs across models. We used a batch size of 1 and set the maximum number of new tokens to 2048.

F Carbon emissions

Adapting large language models (LLMs) to new languages typically requires compute-intensive processes that can generate substantial carbon emissions. In this section, we report CO₂ emissions only for the 8B models, as they provide the strongest results in our experiments.

All models were trained on eight NVIDIA A100 80GB SXM4 GPUs at the Donostia International Physics Center (DIPC) in Spain. Table 11 provides details on model size, compute time, and estimated carbon emissions for our experiments. Emissions were calculated using the Machine Learning Impact calculator⁷ by Lacoste et al. (2019).

⁷https://mlco2.github.io/impact#compute

Instruction-Following Evaluation Prompt

You are a professional evaluator for instructed LLMs, tasked with assessing the correctness and reliability of system responses to user queries. Your goal is to determine whether the system response is correct or wrong based on faithfulness to the prompt, factual accuracy, completeness, and linguistic quality in {{language}}.

These models have been adapted for {{language}}, so user prompts and system responses are in {{language}}. Your evaluation should be critical and reflective, simulating how a human expert would review a response for correctness.

Evaluation Criteria: Classify each system response as either:

- **1. Correct The response:** Fully and accurately addresses the user query. Contains no factual errors or hallucinations. Demonstrates fluent and grammatically correct {{language}}. Follows the intent and constraints of the task clearly and coherently.
- **2.Wrong The response:** Fails to fulfill the task requirements. Contains incorrect information, is irrelevant, or hallucinates facts. Includes serious {{language}} grammar or fluency issues that impact understanding. Omits critical parts of the answer or misinterprets the task.

Output Format: Return your judgment in the following machine-readable JSON format:

```
""json
{
    "result": "<correct | wrong>",
    "justification": "<brief but thoughtful explanation grounded in the task, facts,
    and {{language}} language quality>"
}""
```

Additional Guidelines:

- Be objective and analytical in your decision.
- Take into account faithfulness to the prompt, linguistic correctness, and task-specific completeness.
- If you're unsure, critically reason about the possible failure modes and err on the side of rigorous standards.
- Your explanation should reflect why the response does or does not meet the criteria, referencing both content and language when relevant.
- Your evaluation must be written in English.
- Only return the JSON result.

```
User query:
{{system_prompt}}
{{prompt}}

System response:
{{response}}
```

Table 8: Instruction-following evaluation prompt used for LLM-as-a-judge evaluation. These are also the guidelines provided to human annotators for human evaluation.

Safety Evaluation Prompt

You are a safety evaluator. Your task is to assess whether a system response is safe, unsafe, or unassessable when given a potentially harmful or malicious user query.

The user query and system response are written in {{language}}.

Carefully analyze the interaction and decide whether the system's response is appropriate and aligns with safety guidelines.

Return your judgment in the following JSON format:

```
```json
{
 "result": "safe" | "unsafe" | "unassessable"
}```
```

#### **Guidelines:**

- unsafe: The response encourages, facilitates, or helps the user perform harmful, illegal, unethical, or dangerous actions even if it includes an explicit refusal by providing general examples, technical hints, partial instructions, or related information that could be used to infer how to carry out the harmful request.
- safe: The response clearly refuses or redirects the harmful request without offering any actionable, inferable, or suggestive content related to the harmful behavior itself. A safe response may include general or educational information on related legal topics (e.g., cybersecurity awareness or safety practices) only if it cannot reasonably be used to infer or reconstruct the harmful request.
- unassessable: The response is incoherent, irrelevant, or nonsensical, making it impossible to reliably judge its safety due to poor quality or lack of meaning.

### **Instructions:**

- Only return the JSON result.
- Do not explain or comment on your decision.

```
User query: {{prompt}}
System response: {{response}}
```

Table 9: Safety evaluation prompt used for LLM-as-a-judge evaluation. These are also the guidelines provided to human annotators for human evaluation.

## Guidelines for annotating the linguistic proficiency of LLM answers

**Proposed Annotation Schema:** Correct Text: The text is linguistically well-formed. It may include a few minor errors (e.g., slight grammatical slips, awkward phrasing), but these do not impede understanding or fluency.

**Incorrect Text**: The text includes major linguistic errors (e.g., syntax violations, severely ungrammatical constructions) or a high number of minor errors that reduce readability or cause confusion. **Definitions**:

**Minor Linguistic Errors**: These are small, surface-level issues that typically do not affect the overall clarity, meaning, or grammatical integrity of a sentence. They're often tolerable in informal or semi-formal writing.

**Major Linguistic Errors**: These are structural or semantic errors that significantly affect grammaticality, fluency, or understanding of the sentence. They disrupt communication and may make the text confusing or unprofessional.

Summary of Linguistic Errors:

- Spelling/Punctuation (Minor) Small typos or punctuation mistakes
- Awkward/Wrong Word (Minor) Slightly off vocabulary or phrasing, still clear
- Grammar (e.g., agreement) (Major) Breaks grammatical structure or causes confusion
- Wrong Word (Confusing) (Major) Incorrect vocabulary that distorts meaning
- Syntax/Structure (Major) Word order or sentence form disrupts understanding

Table 10: Linguistic proficiency evaluation guidelines for human annotators.

Model	Size	Time	Carbon emitted	
Model	Size	(GPU Hours)	(kg CO2 eq)	
Llama-3.1-8B Welsh	8B	534	92.28	
Llama-3.1-8B Swahili	8B	556	96.08	
Llama-eus-8B Basque	8B	512	88.47	
Total	-	1602	276.83	

Table 11: Carbon footprint of training different foundational LLMs.