# INVESTIGATING LANGUAGE-SPECIFIC CALIBRATION FOR PRUNING MULTILINGUAL LARGE LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

#### Abstract

Recent advances in large language model (LLM) pruning have shown state-of-theart (SotA) compression results in post-training and retraining-free settings while maintaining high predictive performance. However, previous research mainly considered calibrating based on English text, despite the multilingual nature of modern LLMs and their frequent use in non-English languages. In this paper, we set out to investigate calibrating the pruning of multilingual language models for monolingual applications. We present the first comprehensive empirical study, comparing different calibration languages for pruning multilingual models across diverse languages, tasks, models, and SotA pruning techniques. Our results offer practical suggestions, for example, calibrating in the target language can efficiently retain the language modeling capability but does not necessarily benefit downstream tasks. Through further analysis of latent subspaces, pruning masks, and individual neurons within pruned models, we find that while pruning generally preserves strong language-specific features, it may fail to retain language-specific neuron activation patterns and subtle, language-agnostic features associated with knowledge and reasoning that are needed for complex tasks.

027 028 029

025

026

006

008 009 010

011

013

014

015

016

017

018

019

021

#### 1 INTRODUCTION

State-of-the-art language models often rely on over-parameterization with millions or billions of 031 parameters, resulting in significant memory and computational demands (Zhang et al., 2017; Allen-032 Zhu et al., 2019; Xu & Du, 2023). Pruning is a model compression technique that removes unim-033 portant weights to reduce memory footprint and inference computation (Gholami et al., 2022; Hoe-034 fler et al., 2021; Kuzmin et al., 2023). Recent pruning methods for language models, such as 035 SparseGPT (Frantar & Alistarh, 2023) and Wanda (Sun et al., 2024), demonstrated SotA performance in a post-training and retraining-free setting (Zhu et al., 2023). The pruning process involves 037 passing a small number of examples, i.e. calibration data, to the model to determine the importance 038 of weights for subsequent pruning (Kuzmin et al., 2022; Frantar & Alistarh, 2023; Kuzmin et al., 2023).

Notably, existing studies typically use calibration data in English and evaluate the post-pruning performance in English. On the other hand, most SotA LLMs are multilingual and frequently employed for tasks in non-English languages (Touvron et al., 2023; Achiam et al., 2023; Jiang et al., 2023). It remains unclear how to calibrate pruning to optimize the post-pruning performance of multilingual LLMs on tasks in non-English languages. In other words, the current literature provides little insight into efficient calibration strategies for multilingual LLMs targeting specific non-English languages. For example, if we aim to prune a multilingual LLM and use the pruned LLM for tasks in German, should we calibrate the compression process with text in English, German, or both?

This research, as the first, presents an extensive study investigating the impact of calibration data language on pruning LLMs for language-specific tasks. Further, we performed a set of analyses at the latent subspace level, the matrix level, and the neuron level, within the pruned models, to investigate the underlying causes for our observations. We summarize our key findings as follows:

052

• Calibrating on the target language consistently yields the lowest perplexity (Section 4.1, but will not guarantee optimal performance on downstream tasks (Section 4.2).

- Pruning generally impairs the language-agnostic features of multilingual models, such as reasoning capability (Section 4.2) and the storage and retrieval of knowledge (Section 4.3), which are essential for downstream tasks (Section 5).
  - Pruning generally affects language-agnostic features, which is potentially associated with reasoning and knowledge storage, more than it impacts language-specific features (Section 5.1); Pruning struggles to consistently identify essential neurons in the attention output and FFN down projection, potentially responsible for reasoning (Section 5.2); pruning also struggles to retain language-specific neurons with low activation probability (Section 5.3).

### 063 2 RELATED WORK

065 2.1 MULTILINGUAL LANGUAGE MODELS

067 Most SotA LMs, such as Llama-3 (Meta AI, 2024) and Phi 3(Abdin et al., 2024), are trained on multilingual data, enabling them to understand and generate text in multiple languages (Huang et al., 068 2023; Holmström et al., 2023; Xu et al., 2024; Meta AI, 2024). Although these multilingual LMs 069 follow the general training paradigm of training monolingual LMs, they are often found to behave differently or show unique features from monolingual models (Xu et al., 2024). For example, Deng 071 et al. (2024) reveals that multilingual LMs are prone to generate unsafe outputs on malicious in-072 structions in non-English languages, i.e. multilingual jailbreak. Wang et al. (2023) also identify that 073 all LLMs tested in their study produce significantly more unsafe responses for non-English queries 074 than English ones. Furthermore, previous work on model explanation finds that multilingual models 075 are different from counterpart monolingual models in terms of their internal processes, i.e., different 076 importance distributions on the same input (Jørgensen et al., 2022; Zhao & Aletras, 2024). A general 077 conclusion drawn from these studies is that findings from monolingual settings, particularly those in 078 English, are unlikely to hold in multilingual settings involving non-English languages.

079 080

055

056

058

060

061

062

064

2.2 CALIBRATION OF POST-TRAINING PRUNING

Unlike sparse training (Yuan et al., 2021; Hoang & Liu, 2023; Zhang et al., 2023) or pruning-aware
training (Liu et al., 2021), which require iterative training to achieve sparsity, post-training pruning (Frantar & Alistarh, 2023; Sun et al., 2024) does not require training but eliminates redundant
weights based on the importance of weights that are calculated with calibration data. The retrainingfree feature of post-training pruning makes it a more efficient approach for LLMs.

Prior research has primarily focused on calculating the importance of weight to optimize the performance of pruned models, examining importance metrics and importance ranking ranges (Frantar & Alistarh, 2023; Sun et al., 2024; Zhang et al., 2024). Two studies have investigated the impact of calibration data, looking at the quantity of calibration data (Zhang et al., 2024) and the sources of the data (Williams & Aletras, 2023). These studies have been limited to English. The effects of pruning on the multilingual capabilities of language models, and the impact of the language of calibration data on performance in other target languages remain unknown.

093 094

#### 3 Methodology

095 096

To investigate the impact of the language of calibration data on the pruned model, we compare the performance (pruning error, Signal-to-noise ratio, perplexity, and downstream tasks) in a range of languages between pruned models and their counterpart full-sized models, where different calibration data are applied.

 Models Our experiments use two SotA LLM families: Llama-3 (Meta, 2024), the open-source
 SotA at the time of writing, and Aya-23 (Aryabumi et al., 2024), renowned for its multilingual pretraining. As our evaluation focuses on instructed generation tasks, we employ the instruction-tuned
 versions of both families (Chrysostomou et al., 2023). Consequently, our experimental setup includes Llama-3-instruction models in 8B and 70B parameter sizes, alongside Aya-instruction models in 8B and 35B parameter sizes.

**Languages** We include seven languages in our study: Arabic (AR), German (DE), English (EN), Spanish (ES), Russian (RU), Swahili (SW), and Chinese (ZH). This selection spans six language

families, four writing systems, and encompasses both high and mid-resource languages. We include
 Swahili as an outlier calibration language, as neither Llama-3 or Aya-23 includes it in their pre training corpora. A summary of languages used in our paper is given in Table 5 in Appendix B.

Pruning and Calibration We construct calibration sets in seven different languages from mC4 (Raffel et al., 2019). Specifically, following Frantar & Alistarh (2023); Sun et al. (2024), we randomly sample 128 sequences of 8192 tokens for each language, from the mC4 split for that language.<sup>1</sup>

We focus on two SotA post-training *pruning methods*, Wanda (Sun et al., 2024) and SparseGPT (Frantar & Alistarh, 2023). Unless stated otherwise, we prune for 50% unstructured sparsity to impose fewer restrictions while keeping all other hyperparameters as in the original paper.

**Evaluation Downstream Tasks** We compare the performance of the pruned models calibrated on different languages using the *perplexity* on a subset of the mC4 validation set, and a selection of downstream tasks in different target languages: MMLU (Hendrycks et al., 2021), MKQA (Longpre et al., 2020), and Belebele (Bandarkar et al., 2023). To better isolate the impact of the calibration language, we also chose mARC and mHellaSwag (Dac Lai et al., 2023) to compare with their original versions, ARC (Clark et al., 2018) and HellaSwag (Zellers et al., 2019) in English for their low sensitivity to the choice of calibration samples (Williams & Aletras, 2023). These tasks primarily assess commonsense reasoning, reading comprehension, and question answering using multiple choice questions. We present further details of each task in Appendix C. Unless stated otherwise, we evaluate in a zero-shot setting. 

Implementation Details We adopt the code from Sun et al. (2024) for implementing model pruning.
We use EleutherAI Evaluation Harness (Gao et al., 2024) for a robust and reproducible evaluation.
We use Huggingface (Wolf et al., 2020) for loading datasets and models. All experiments are conducted with at most two NVIDIA A100 GPUs.

#### 4 Results

4.1 PRUNING RESULTS



Table 1: Language-specific perplexity (PPL), Signal-to-noise ratio (SNR) and pruning error of hidden states for 50% pruned Llama-3 8B, averaged over three pruning runs. The leftmost columns show the model, the pruning technique, and the calibration language. For PPL and pruning errors, the smaller the value (the darker), the better; while for SNR, the greater the value (the lighter), the better. This evaluation uses 128 inputs sampled from the mC4 validation set (1048576 tokens).

 <sup>&</sup>lt;sup>158</sup> <sup>1</sup>In the original implementation of Wanda and SparseGPT, each sample contains 2,048 tokens, matching the maximum context length of Llama 1. We increase the length to 8,192 tokens to match the maximum context of Llama-3. This is to avoid potential context length-dependence pruning behavior, although Sun et al. (2024) reports a negligible return from increased calibration set sizes. For mixing different calibration languages, we mix in equal shares and keep the budget of 128 samples.

Table 1 presents the perplexity (PPL), Signal-to-noise ratio (SNR) and pruning error of hidden states for models pruned to 50% sparsity and calibrated on different languages. SNR and pruning error estimate the pruning accuracy (pruning performance); ideal pruning is able to preserve activations identical to the full model. The formal expression and detailed explanation of SNR and pruning error are presented in the Appendix D.1. Perplexity reflects the general language modeling capability. The column headers denote the evaluation languages, while the row headers indicate the calibration languages.

Overall, our findings reveal that no single calibration language consistently outperforms others across perplexity, pruning error, or SNR metrics. The optimal calibration language varies depending on the target evaluation language. Notably, *calibrating on the target language itself generally yields the best pruning performance and the lowest perplexity*, as evidenced by the diagonal pattern in Table 1. There are a few exceptions. For instance, when evaluating Llama-3 pruned with Wanda on Chinese, Russian calibration performs optimally on perplexity (23.6), slightly outperforming Chinese calibration (24.4).

4.2 DOWNSTREAM TASK PERFORMANCE

## 178179 How to select the calibration language to optimize performance in downstream tasks?

Table 2 shows the models' performance on downstream tasks. The column headers denote the evaluation languages, while the row headers indicate the calibration languages. Given a downstream task in a specific language, we analyze the results column-wise; the lighter the color in a column, the better the performance on the specified downstream task.

First, for both pruning methods, the calibration language affects downstream task performance.
This impact is particularly pronounced when pruning Llama-3 8B model with SparseGPT or Aya-23 models with either Wanda or SparseGPT, as evidenced by the greater color difference in the table, compared to pruning Llama-3 8B with Wanda. For example, when evaluating Llama-3 8B model on Belebele in English, pruning with Wanda results in performance ranging from 58.5 with Spanish calibration to 65.0 with English calibration. In contrast, when pruning Llama-3 8B with SparseGPT the performance varies from 50.0 with Chinese calibration to 71.4 with English calibration.

Unlike the perplexity evaluation in Table 1, calibration with the target language does not reliably result in good performance. For instance, on MMLU, the pruned Llama-3 8B model mainly achieves higher accuracy on evaluation languages other than on its calibration languages. Overall, there are 50 evaluation groups, 25 columns for Wanda and 25 for SpareGPT, in 26 comparison cases. Calibration with the target language yields the best performance, 8 for Wanda and 18 for SparseGPT (e.g.





Arabic on ARC and MMLU). For the remaining cases, calibration with the target language achieves
 the second-best or comparable performance. Therefore, calibrating using the target language typically results in acceptable, though not consistently the best performance. Considering the standard
 deviations shown in Table 7 in Appendix E, *calibration with the target language is a reasonable choice out of downstream performance*.

221 In terms of downstream performance across different languages, the full-sized baseline models ex-222 hibit the strongest performance in English, followed by other Latin languages, with Russian next in 223 line. Arabic and Chinese downstream tasks generally are the most challenging for models. Further, 224 pruning can sometimes alter the original ranking of languages observed in the baseline models. For 225 example, on the Belebele benchmark, the Llama-3 8B model achieves a baseline accuracy of 55.3 226 for Arabic and 44.8 for Chinese, but the pruned models reverse this trend and achieve an accuracy of 45.1 for Arabic and 57.0 for Chinese. That is, pruning can shift which languages the model 227 performs best or worst on. 228

229 *Calibrating on an outlier language or a similar one could benefit downstream tasks in non-English?* 

We include Swahili as an outlier language, which is out-of-domain for the model as it is not included in the pre-training corpora of Llama-3 and Aya-23. In column-wise comparison, the SW cell is consistently the darkest or relatively dark one. That is, in most cases, calibration with Swahili results in the worst or the second-worst performance across different calibration language strategies, given a downstream task.

Furthermore, there is no consistent pattern related to the similarity of calibration-evaluation language pairs. For instance, Latin language pairs such as English-Spanish (calibrating in English and evaluating in Spanish) or pairs from the same language family, like English-German, do not always yield optimal performance. Conversely, pairs with different writing systems, such as Russian-English or Spanish-Arabic, do not consistently perform poorly. On the other hand, calibration with a dissimilar language, i.e. Chinese, often results in particularly low accuracy across many tasks and evaluation languages, as demonstrated by the darker row of "ZH". In summary, we observe *no benefit from calibrating with an outlier language or a similar language*.

243 244 Does the model or the pruning method matter?

245 Between Llama-3 8B and Aya-23 8B, despite their similar decoder-only architecture, Table 2 highlights distinct performance patterns between the two models under investigation. With its larger 246 vocabulary size in the six non-English languages, Aya-23 8B generally outperforms the Llama-3 247 8B model in most evaluation languages and tasks, both before and after pruning. Notably, Aya-23 248 8B experiences less performance drop after pruning but shows less stable results, often perform-249 ing better on languages other than the one used for calibration. Overall, these results suggest that 250 vocabulary size and pre-training data impact baseline performance and accuracy after pruning in a 251 multilingual context. 252

Between Wanda and SparseGPT, Llama-3 8B's performance degrades less after pruning with
SparseGPT pruning, while the performance of the pruned Aya-23 8B model varies, that is, excelling with Wanda or SparseGPT depends on the task. Specifically, SparseGPT produces a distinct
diagonal pattern for Aya-23 8B on ARC and HellaSwag, while Wanda often yields superior performance for calibration with non-target languages on Belebele and MMLU. No pruning technique
consistently performs best in all tasks.

250 259 260

261

#### 4.3 OPEN DOMAIN QUESTION ANSWERING WITHOUT CONTEXT

262 How does pruning impact the knowledge stored in multilingual LLMs?

Table 3 shows the open-domain question-answering performance of Llama-3 8B and Aya-23 8B across various calibration-evaluation language pairs. This "closed-book" task, MKQA, provides no context in the prompt, and relies on the model's internal knowledge to generate answers. To ensure fair cross-languages comparisons, the MKQA dataset is fully parallel and primarily consists of entity-based and structured "atomic" answer types. See Appendix C for further details. First, we observe significant performance differences among the evaluation languages for the full-size models. While Latin languages perform best, Arabic and Chinese perform worst. Further, we notice a significant accuracy drop by pruning, even for performance in English. In summary, *pruning sub-*

## stantially impacts the storage and retrieval of knowledge in a multilingual model across different languages.

273 274

305

306

307

312

313

314

315

316

317 318

319

320

321 322

323

4.4 MULTIPLE CALIBRATION LANGUAGES

Using more languages in the calibration, i.e. bilingual and multilingual calibration set, will benefit the performance of pruned multilingual models?

We repeat the experiments but include more languages in the calibration set. As English is the dominant language in pre-training data, we first test the combination of English and the target language for a bilingual calibration setup. We further experiment with including all seven languages in the calibration, i.e. multilingual setup. For all setups, the total calibration sample number remains the same, i.e. 128.

286 Comparing Table 8 with Table 2, we observe that 287 bilingual and multilingual calibration sets generally 288 outperform monolingual calibration for non-English 289 target languages. For example, for Chinese on ARC, 290 the seven-lingual calibration set, AR-DE-EN-ES-291 RU-SW-ZH, often yields higher accuracy than cali-292 bration with Chinese, i.e. the target language, across 293 models and pruning methods. That is, AR-DE-EN-ES-RU-SW-ZH achieves 28.3 (Wanda, Llama-3 294 8B), 30.0 (SparseGPT, Llama-3 8B), 32.0 (Wanda, 295 Aya-23 8B) and 30.9 (SparseGPT, Aya-23 8B), con-296 sistently higher or comparable to monolingual cal-297 ibration set, i.e. 27.0, 30.1, 30.8 and 31.7 respec-298

			MKQA <sub>[f1]</sub>								
			AR	DE	EN	ES	RU	ZH			
	-	-	8.9	26.8	38.3	27.2	16.4	2.6			
		AR	0.9	6.3	18.4	11.4	5.6	2.0			
Instruct	_	DE	0.3	6.0	19.5	11.5	5.7	2.1			
	abr	EN	0.2	6.3	19.4	11.6	4.5	1.5			
	Vai	ES	0.1	6.6	19.6	11.8	5.4	1.7			
Ä	~	RU	0.4	7.0	19.2	12.0	5.5	2.1			
8		ZH	0.7	4.8	18.0	10.1	3.5	2.1			
a-3		AR	3.5	8.0	19.8	11.3	6.4	0.9			
am	Ы	DE	0.3	10.5	19.8	11.2	6.1	1.6			
П	õ	EN	1.0	9.7	20.2	12.1	5.8	1.5			
	Sparse	ES	0.8	8.4	20.3	12.0	6.3	2.4			
		RU	0.9	9.0	19.9	11.9	7.3	1.8			
		ZH	0.1	7.4	18.1	9.3	5.1	0.7			
-											
	-	-	14.8	28.8	46.8	31.0	22.7	0.3			
	-	- AR	14.8 6.3	28.8 8.8	46.8 16.0	31.0 10.5	22.7 7.3	0.3 0.1			
	-	- AR DE	14.8 6.3 5.8	28.8 8.8 8.4	46.8 16.0 18.0	31.0 10.5 10.9	22.7 7.3 7.8	0.3 0.1 0.1			
	nda	- AR DE EN	14.8 6.3 5.8 5.2	28.8 8.8 8.4 8.1	46.8 16.0 18.0 17.5	31.0 10.5 10.9 10.2	22.7 7.3 7.8 7.0	0.3 0.1 0.1 0.1			
ßB	Vanda	- AR DE EN ES	14.8 6.3 5.8 5.2 5.6	28.8 8.8 8.4 8.1 8.1	46.8 16.0 18.0 17.5 16.8	31.0 10.5 10.9 10.2 10.6	22.7 7.3 7.8 7.0 7.6	0.3 0.1 0.1 0.1 0.1			
.3-8B	Wanda	- DE EN ES RU	14.8 6.3 5.8 5.2 5.6 5.9	28.8 8.8 8.4 8.1 8.1 7.7	46.8 16.0 18.0 17.5 16.8 16.4	31.0 10.5 10.9 10.2 10.6 10.6	22.7 7.3 7.8 7.0 7.6 7.7	0.3 0.1 0.1 0.1 0.1 0.1			
a-23-8B	Wanda	- DE EN ES RU ZH	14.8 6.3 5.8 5.2 5.6 5.9 5.3	28.8 8.8 8.4 8.1 8.1 7.7 8.5	46.8 16.0 18.0 17.5 16.8 16.4 16.8	31.0 10.5 10.9 10.2 10.6 10.6 10.4	22.7 7.3 7.8 7.0 7.6 7.7 6.8	0.3 0.1 0.1 0.1 0.1 0.1 0.1			
Aya-23-8B	Wanda	- DE EN ES RU ZH AR	14.8 6.3 5.8 5.2 5.6 5.9 5.3 5.4	28.8 8.8 8.4 8.1 8.1 7.7 8.5 8.4	46.8 16.0 18.0 17.5 16.8 16.4 16.8 18.0	31.0 10.5 10.9 10.2 10.6 10.6 10.4 10.6	22.7 7.3 7.8 7.0 7.6 7.7 6.8 6.8	0.3 0.1 0.1 0.1 0.1 0.1 0.1 0.1			
Aya-23-8B	PT Wanda .	- DE EN ES RU ZH AR DE	14.8 6.3 5.8 5.2 5.6 5.9 5.3 5.4 4.9	28.8 8.8 8.4 8.1 8.1 7.7 8.5 8.4 7.8	46.8 16.0 18.0 17.5 16.8 16.4 16.8 18.0 19.4	31.0 10.5 10.9 10.2 10.6 10.6 10.4 10.6 8.6	22.7 7.3 7.8 7.0 7.6 7.7 6.8 6.8 6.8 6.4	0.3 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.0			
Aya-23-8B	eGPT Wanda	- AR DE EN ES RU ZH AR DE EN	14.8 6.3 5.8 5.2 5.6 5.9 5.3 5.4 4.9 3.7	28.8 8.8 8.4 8.1 8.1 7.7 8.5 8.4 7.8 7.5	46.8 16.0 18.0 17.5 16.8 16.4 16.8 18.0 19.4 17.3	31.0 10.5 10.9 10.2 10.6 10.6 10.4 10.6 8.6 9.9	22.7 7.3 7.8 7.0 7.6 7.7 6.8 6.8 6.8 6.4 6.3	0.3 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.0 0.1			
Aya-23-8B	arseGPT Wanda	AR DE EN ES RU ZH AR DE EN ES	14.8         6.3         5.8         5.2         5.6         5.9         5.3         5.4         4.9         3.7         3.6	28.8 8.8 8.4 8.1 8.1 7.7 8.5 8.4 7.8 7.5 7.7	46.8 16.0 18.0 17.5 16.8 16.4 16.8 18.0 19.4 17.3 16.6	31.0 10.5 10.9 10.2 10.6 10.6 10.4 10.6 8.6 9.9 8.8	22.7 7.3 7.8 7.0 7.6 7.7 6.8 6.8 6.8 6.4 6.3 5.7	0.3 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.0 0.1 0.2			
Aya-23-8B	SparseGPT Wanda	AR DE EN ES RU ZH AR DE EN ES RU	14.8         6.3         5.8         5.2         5.6         5.9         5.3         5.4         4.9         3.7         3.6         5.1	28.8 8.8 8.4 8.1 7.7 8.5 8.4 7.8 7.5 7.7 7.1	46.8 16.0 18.0 17.5 16.8 16.4 16.8 18.0 19.4 17.3 16.6 18.3	31.0 10.5 10.9 10.2 10.6 10.6 10.4 10.6 <b>8.6</b> 9.9 <b>8.8</b> 9.8	22.7 7.3 7.8 7.0 7.6 7.7 6.8 6.8 6.8 6.4 6.3 5.7 6.3	0.3 0.1 0.1 0.1 0.1 0.1 0.1 0.1 0.0 0.1 0.2 0.4			

Table 3: MKQA F1 accuracy over seven evaluation languages for the Llama-3 8B and Aya-23 8B models pruned with Wanda and SparseGPT for 50% unstructured sparsity.

tively. On the other hand, multilingual calibration pruning often degrades the downstream performance in English, particularly on ARC when pruning Aya-23 with SparseGPT and on MMLU when pruning Llama-3 with Wanda. Interestingly, the best bilingual combination does not necessarily contain the target language but English. For example, for Russian on ARC when pruning with Wanda, calibration with EN-ES and AR-EN leads to the optimal performance for Llama-3 8B and Aya-23 8B.

In summary, *employing bilingual and multilingual calibration sets occasionally improves performance on downstream tasks, compared to monolingual calibration.* However, there is no clear pattern identifying which specific language combinations are most effective for a given target downstream language.



Table 4: Pruned Model Performance averaged over three pruning runs. The lighter, the less drop, the better. Five languages: AR, DE, EN, ES, RU. Seven languages: AR, DE, EN, ES, RU, SW, ZH.

## 324 4.5 IMPACT OF MODEL SIZES

To investigate the scaling impact towards pruning behavior, we repeat experiments in Table 7 on the larger Llama-3 70B and Aya-23 35B models, the results of which are reported in Appendix E. Overall, pruning performance increases with higher baseline accuracy of the full-sized models. However, we observe that *the performance patterns and findings from the smaller models, mentioned above, do not consistently hold true on their bigger counterparts*.

331 332

333

#### 4.6 QUANTIZATION

We further explore the impact of the calibration language in quantization on downstream task perfor-334 mance across different languages. We use GPTQ (Frantar et al., 2023) to quantize weights to 4 bits 335 with a group size of 128 and 8 bits with a group size of 128 (equivalent to 50% sparsity in pruning) 336 on LLaMA-3-8B. We follow our pruning setup for downstream tasks and languages for calibration 337 and evaluation. The results are present in Table 14 and 15 in Appendix F, revealing several key 338 findings that are consistent with our findings on pruning: (1) calibrating with the target language 339 tends to yield better downstream performance; (2) pruning can alter which languages the model per-340 forms best or worst on; and (3) calibrating with an outlier or a linguistically similar language does 341 not provide any notable advantage.

342 343

#### 5 ANALYSIS

344 345

In reviewing Table 1 and Table 7 together, we poist that calibrating on the target language effectively 346 preserves language-related features but not language-agnostic features, such as knowledge memories 347 and reasoning abilities. This implies that the smallest weights or weights of smallest activations 348 contribute to the model's nuanced reasoning or knowledge memorization and retrieval processes. 349 The rationale here is that both SNR and pruning error measure the extent of model change after 350 pruning, which is dominated by large values, that pruning aims to preserve. On the other hand, 351 perplexity reflects the overall language modeling capability. The similar diagonal patterns across the 352 three metrics, SNR, pruning error and perplexity, as shown in Table 1, suggest that calibrating on the 353 target language, i.e. identifying and pruning the smallest values, preserves the next-token prediction 354 capability (perplexity) in the target language, along with the optimal pruning performance in terms 355 of pruning the smallest values (SNR and pruning errors).

However, optimal performance in downstream tasks requires more than just surface-level language
 modeling; it also depends on robust reasoning capabilities. The suboptimal performance observed in
 downstream tasks may be attributed to the degradation of language-agnostic reasoning capabilities.
 This phenomenon suggests that removing the weights or weights of the smallest activations might
 inadvertently compromise the model's ability to process reasoning, underscoring the reasoning role
 of these pruned weights, i.e. the smallest weights or weights of the smallest activations.

362 To validate our hypothesis, we investigate the internal changes of pruned models on three different 363 levels: subspace level, matrix level, neuron level (columns in a matrix, followed by non-linearity). 364 Previous studies on multilingualism of LLMs have sought to separate language-specific features 365 from language-agnostic features, either at the neuron level (Tang et al., 2024; Zhao et al., 2024; 366 Wang et al., 2024) or from an extracted subspace perspective (Xie et al., 2022). To comprehend 367 the underlying mechanisms driving the disparities in pruning metrics and downstream performance when using different calibration languages, we examine the shifts in language-specific and agnostic 368 features and neurons induced by the pruning process. 369

370

371 372

#### 5.1 LANGUAGE-SPECIFIC SUBSPACE REPRESENTATIONS

The core idea of Low-rank Subspace for language-Agnostic Representations (LSAR) by Xie et al. (2022) is to construct a mean embedding matrix  $M \in \mathbb{R}^{d \times L}$  by concatenating L language embeddings. Subsequently, LSAR decomposes M into a vector  $\mu$  representing shared signals across languages and a matrix  $M_s$  specifying a low-rank subspace on which different languages express different linguistic signals. This decomposition process is achieved via singular value decomposition on solving:  $\min_{\mu,M_s,\Gamma} || M - \mu \mathbf{1}^{\top} - M_s \Gamma^{\top} ||_F^2$  s.t.  $\mu \perp \text{Span}(M_s)$  with the orthogonality constraint.  $\Gamma \in \mathbb{R}^{L \times r}$  represents the coordinates of language-specific signals along the subspace's r components, and  $\mathbf{1} \in \mathbb{R}^d$  is a vector of all ones. We can extract language-specific features *a* from a sentence embedding *e* using  $M_s$  by projecting *e* into and back from the low-rank, language signal retaining subspace,  $a = e - s = e - M_s M_s^T e$ . See Appendix D.2 for more details and the experiment setup.

We employ LSAR to decomposite the token-wise averaged embedding output of each layer of SparseGPT-pruned Llama-3 8B into language-specific and language-agnostic features, and estimate their changes after pruning. Figure 1 shows the layer-wise magnitude of differences ( $\Delta$  magnitude) of (a) language-agnostic features and (b) language-specific features after pruning. From a high-level perspective, the greater  $\Delta$  magnitude suggests the greater changes in terms of hidden state after pruning, thus the greater pruning error and the worse pruning performance. Unlike the pruning error presented in Table 1, here we separately analyze the pruning error for language-agnostic and language-specific features. See Figure 5 in Appendix E for layer-wise full plots.

391 First, calibrating on the target language leads to relatively smaller pruning errors in terms of 392 language-specific features as shown in the sub Figure (b) in Figure 1. This phenomenon is ob-393 served across layers, particularly noticeable in layer 32, as pinpointed by stars in Figure 1. This 394 potentially explains the findings in Section 4.1 that calibration on the target language leads to the 395 lowest perplexity, which is associated with a robust language-specific linguistic modeling capability. On the other hand, as indicated by the relatively flat horizontal lines across languages and layers in 396 sub Figure (a), the pruning error on the language-agnostic features remains similar regardless of 397 the calibration languages. This pattern helps explain the sub-optimal downstream task performance 398 observed in Table 7, where no single calibration language consistently yields optimal performance 399 across downstream tasks, including cases where the calibration is performed on the target language. 400 That is, the selection of calibration language is unlikely to impact the language-agnostic features 401 associated with understanding, reasoning, and knowledge retrieval. In contrast, the preservation 402 of language-specific features, such as modeling capabilities related to perplexity, depends on the 403 selection of calibration languages.

Second, the middle layers (as shown in the second and third columns in Figure 1) exhibit greater  $\Delta$  magnitude on language-agnostic feature representations and smaller  $\Delta$  magnitude on languagespecific feature representations. This indicates that pruning errors can be predominantly attributed to the pruning errors on language-agnostic features, with less pruning error arising from language-specific features. Therefore, we conclude *pruning generally affects language-agnostic features—potentially associated with reasoning and knowledge storage—more significantly than it impacts language-specific features.* 

411 412

413 414

415

416

417

418

419

420

421

422

423

424

425

426

427

#### 5.2 PRUNING MASK SIMILARITY

To better understand the internal changes within the model after pruning, we conduct a matrixlevel analysis by calculating the Intersection over Union (IoU) of pruning masks across different



Figure 1: Magnitude difference of language-agnostic and specific features after pruning, averaged over 200 data per language from validation sets. The X-axis represents the evaluation language. The calibration languages are indicated by the color as: AR, DE, ES, EN, RU, SW, ZH. The greater Δ magnitude, the greater pruning error on language-agnostic features (Fig.b) or language-specific features (Fig.c).



Figure 2: Pruning similarities (IoU) between using different calibration languages. Left: IoU of pruning masks for three calibration sets of the same language. Right: IoU between pruning masks for different calibration languages. The higher IoU (indicated as a lighter color), the more similar pruning masks between different calibration languages.

448

449

calibration sets, obtaining a measure of pruning mask similarity. The formal expression and details explaning IoU are presented in Appendix D.3. We use the pruning masks from the Llama-3 8B model pruned with SparseGPT for 50% unstructured sparsity. To reduce calibration set-dependent noise as prevalent in the downstream tasks, we first compute the intersection  $M_I^l$  of pruned neuron indices  $M_i^l$  across three pruned models calibrated with different seeds *i* but in the same language *l*. This intersection represents more stable neuron indices.

The IoU of the left side of Figure 2 depicts the proportion of  $M_I^l$  with respect to all pruned neuron indices. It reveals high pruning mask similarity in the attention query, key and value in the first layer, while the attention output projection varies more significantly, especially after the 20th layer. This suggests that *pruning struggles to consistently identify essential neurons in the attention output projection, partly responsible for the models reasoning capabilities.* 

The right plot compares the IoU of intersected neuron indices  $M_i^l$  from the left plot across calibration in German, English, and Chinese. Notably, the attention query, key, and value in the first layer consistently achieve high IoUs above 0.95 across all languages, indicating that these components handle inputs similarly, irrespective of language differences. However, the attention output and FFN down projection show lower IoU, especially in early layers, with similarity peaking in middle layers (3rd to 15th) before decreasing again in later layers.

This pattern suggests that the attention output and FFN down projection in early and late layers handle language-specific signals, while middle layers process language-agnostic signals, supporting Figure 1. In other words, *early layers focus on language comprehension, middle layers on language-independent reasoning, and later layers on generating language-specific predictions.* This aligns with Zhao et al. (2024), who proposes that LLMs first comprehend queries by converting multilingual inputs into English in the early layers, reason in English in the intermediate layers, and then generate responses aligned with the input language in the final layers.

477 478

479

#### 5.3 LANGUAGE-SPECIFIC NEURONS

This section investigates neuron-level changes after pruning using Language Activation Probability Entropy (LAPE) as introduced by Tang et al. (2024). We focus on neurons of the up projections of the feed-forward layers (FFNs) followed by the non-linearity. A neuron is considered activated when the non-linearity output is greater than zero. LAPE measures the likelihood of individual neurons activating across different language inputs, taking neurons with high activation probability for one or two languages and lower probabilities for all others (i.e. low LAPE score) as language-specific. Details on the LAPE calculation and experiment settings are provided in Appendix D.4.



Figure 3: Language Entropy for FFN neurons of Llama-3 8B and its SparseGPT-pruned versions calibrated for DE. The bottom-right plot visualizes the activation likelihood of neurons on DE. The bottom subfigure shows the mean absolute neuron outputs for DE input. *The lower the LAPE score, the more specialized the neuron is for a particular language.* See Appendix G for the activation probability and neuron output magnitudes across all languages.

506 In Figure 3, each bar and boxplot corresponds to the same group of neurons, comprising 2% of the 507 total neuron population. Neurons in each plot are sorted in ascending order based on their LAPE 508 score in the full-sized model, enabling a comparison of changes in LAPE scores after pruning. We 509 only plot the neuron groups of the 30% lowest LAPE score due to the rapidly diminishing variance 510 for boxplots of higher LAPE scores. This suggests a lower impact of pruning on the activation frequency of high LAPE neurons, i.e. language-agnostic neurons. On the other hand, for the neu-511 ron group with the lowest LAPE scores (the leftmost boxplot), the higher the model-sparsity is, the 512 longer whiskers (higher entropy variance) these neurons have. This indicates that *pruning intro-*513 duces LAPE noise, shifting the LAPE score distribution and creating new language-specific (low 514 LAPE) and agnostic (high LAPE) neurons. The distribution shift of language-agnostic neurons in 515 FFNs may contribute to performance degradation in downstream tasks. This hypothesis aligns with 516 previous causal tracing studies, which have frequently identified FFNs are crucial for knowledge 517 retention and retrieval (Meng et al., 2022; 2023).

518 Moreover, by analysis of boxplots and the bottom-right plot together, we find that the neuron group 519 of the lowest LAPE has the lowest activation probabilities, while neurons with the highest LAPE 520 scores activate more frequently. This indicates that *pruning struggles to retain low activation prob-*521 *ability of language-specific neurons*. We also examine the absolute neuron output magnitude for 522 these neglected neurons. The bottom plot in Figure 3 shows that despite the low activation proba-523 bility, low LAPE neurons tend to have high output magnitudes. Pruned models show no significant 524 differences in output magnitudes, leading us to conclude that pruning retains neuron output magni-525 tudes in FFN modules but may fail to preserve activation frequency.

526 527 528

501

502

503

504 505

#### 6 CONCLUSION

In this paper, we empirically demonstrate that the choice of calibration language influences the downstream performance of pruned models across various evaluation languages. Particularly, calibration with the target language mainly benefits preserving perplexity performance but does not necessarily help maintain the downstream task performance. Our analysis of changes within the model after pruning indicates that calibrating in the test language does not reliably benefit layers associated with semantic understanding. Additionally, we recommend practitioners do not depend on perplexity for accessing the pruned model, or performance in English to estimate the performance on target languages.

- 537
- 538
- 539

## 540 REFERENCES

549

550

- Marah Abdin, Sam Ade Jacobs, Ammar Ahmad Awan, Jyoti Aneja, Ahmed Awadallah, Hany
  Awadalla, Nguyen Bach, Amit Bahree, Arash Bakhtiari, Harkirat Behl, et al. Phi-3 technical report: A highly capable language model locally on your phone. *arXiv preprint arXiv:2404.14219*, 2024.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical
  report. *arXiv preprint arXiv:2303.08774*, 2023.
  - Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via overparameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019.
- Viraat Aryabumi, John Dang, Dwarak Talupuru, Saurabh Dash, David Cairuz, Hangyu Lin, Bharat
  Venkitesh, Madeline Smith, Jon Ander Campos, Yi Chern Tan, Kelly Marchisio, Max Bartolo, Sebastian Ruder, Acyr Locatelli, Julia Kreutzer, Nick Frosst, Aidan Gomez, Phil Blunsom, Marzieh
  Fadaee, Ahmet Üstün, and Sara Hooker. Aya 23: Open weight releases to further multilingual
  progress, 2024. URL https://arxiv.org/abs/2405.15032.
- Lucas Bandarkar, Davis Liang, Benjamin Muller, Mikel Artetxe, Satya Narayan Shukla, Donald Husa, Naman Goyal, Abhinandan Krishnan, Luke Zettlemoyer, and Madian Khabsa. The belebele benchmark: a parallel reading comprehension dataset in 122 language variants. *arXiv preprint arXiv:2308.16884*, 2023.
- George Chrysostomou, Zhixue Zhao, Miles Williams, and Nikolaos Aletras. Lighter, yet more faithful: Investigating hallucinations in pruned large language models for abstractive summarization.
   *arXiv preprint arXiv:2311.09335*, 2023.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and
   Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge.
   *arXiv:1803.05457v1*, 2018.
- Viet Dac Lai, Chien Van Nguyen, Nghia Trung Ngo, Thuat Nguyen, Franck Dernoncourt, Ryan A Rossi, and Thien Huu Nguyen. Okapi: Instruction-tuned large language models in multiple languages with reinforcement learning from human feedback. *arXiv e-prints*, pp. arXiv–2307, 2023.
- Yue Deng, Wenxuan Zhang, Sinno Jialin Pan, and Lidong Bing. Multilingual jailbreak challenges
   in large language models. In *The Twelfth International Conference on Learning Representations*,
   2024. URL https://openreview.net/forum?id=vESNKdEMGp.
- Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot. In *International Conference on Machine Learning*, pp. 10323–10337. PMLR, 2023.
- Elias Frantar, Saleh Ashkboos, Torsten Hoefler, and Dan Alistarh. OPTQ: Accurate quantization
   for generative pre-trained transformers. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=tcbBPnfwxS.
- Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024. URL https://zenodo.org/records/12608602.
- Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of quantization methods for efficient neural network inference. In *Low-Power Computer Vision*, pp. 291–326. Chapman and Hall/CRC, 2022.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Ja cob Steinhardt. Measuring massive multitask language understanding. In *International Confer- ence on Learning Representations*, 2021. URL https://openreview.net/forum?id=
   d7KBjmI3GmQ.

Duc NM Hoang and Shiwei Liu. Revisiting pruning at initialization through the lens of ramanujan graph. *ICLR 2023*, 2023.

Torsten Hoefler, Dan Alistarh, Tal Ben-Nun, Nikoli Dryden, and Alexandra Peste. Sparsity in deep
 learning: Pruning and growth for efficient inference and training in neural networks. *Journal of Machine Learning Research*, 22(241):1–124, 2021.

Oskar Holmström, Jenny Kunz, and Marco Kuhlmann. Bridging the resource gap: Explor ing the efficacy of English and multilingual LLMs for Swedish. In Nikolai Ilinykh, Felix
 Morger, Dana Dannélls, Simon Dobnik, Beáta Megyesi, and Joakim Nivre (eds.), Proceedings
 of the Second Workshop on Resources and Representations for Under-Resourced Languages
 and Domains (RESOURCEFUL-2023), pp. 92–110, Tórshavn, the Faroe Islands, May 2023.
 Association for Computational Linguistics. URL https://aclanthology.org/2023.
 resourceful-1.13.

Haoyang Huang, Tianyi Tang, Dongdong Zhang, Xin Zhao, Ting Song, Yan Xia, and Furu Wei. Not all languages are created equal in LLMs: Improving multilingual capability by cross-lingualthought prompting. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Computational Linguistics: EMNLP 2023*, pp. 12365–12394, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.findings-emnlp.826.
URL https://aclanthology.org/2023.findings-emnlp.826.

- Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
  Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
  Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- Rasmus Jørgensen, Fiammetta Caccavale, Christian Igel, and Anders Søgaard. Are multilingual sentiment models equally right for the right reasons? In Jasmijn Bastings, Yonatan Belinkov, Yanai Elazar, Dieuwke Hupkes, Naomi Saphra, and Sarah Wiegreffe (eds.), *Proceedings of the Fifth BlackboxNLP Workshop on Analyzing and Interpreting Neural Networks for NLP*, pp. 131–141, Abu Dhabi, United Arab Emirates (Hybrid), December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.blackboxnlp-1.11. URL https://aclanthology.org/2022.blackboxnlp-1.11.
- Andrey Kuzmin, Mart Van Baalen, Yuwei Ren, Markus Nagel, Jorn Peters, and Tijmen Blankevoort.
   Fp8 quantization: The power of the exponent. *Advances in Neural Information Processing Systems*, 35:14651–14662, 2022.
- Andrey Kuzmin, Markus Nagel, Mart Van Baalen, Arash Behboodi, and Tijmen Blankevoort. Pruning vs quantization: Which is better? In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL https://openreview.net/forum?id=00U1ZXXxs5.
- Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris
  Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion
  Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav
  Petrov. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl\_a\_00276. URL
  https://aclanthology.org/Q19-1026.
- Shiwei Liu, Tianlong Chen, Xiaohan Chen, Zahra Atashgahi, Lu Yin, Huanyu Kou, Li Shen, Mykola
   Pechenizkiy, Zhangyang Wang, and Decebal Constantin Mocanu. Sparse training via boosting
   pruning plasticity with neuroregeneration. *Advances in Neural Information Processing Systems*,
   34:9908–9922, 2021.
- Shayne Longpre, Yi Lu, and Joachim Daiber. Mkqa: A linguistically diverse benchmark for multilingual open domain question answering, 2020. URL https://arxiv.org/pdf/2007. 15207.pdf.
- Kevin Meng, David Bau, Alex J Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In Alice H. Oh, Alekh Agarwal, Danielle Belgrave, and Kyunghyun Cho (eds.),
   Advances in Neural Information Processing Systems, 2022. URL https://openreview.net/forum?id=-h6WAS6eE4.

648 649 650	Kevin Meng, Arnab Sen Sharma, Alex J Andonian, Yonatan Belinkov, and David Bau. Mass-editing memory in a transformer. In <i>The Eleventh International Conference on Learning Representations</i> , 2023. URL https://openreview.net/forum?id=MkbcAHIYgyS.
652 653	Meta. Introducing Meta Llama 3: The most capable openly available LLM to date — ai.meta.com. https://ai.meta.com/blog/meta-llama-3/, 2024. [Accessed 15-07-2024].
654 655	Meta AI. Meta llama 3. https://ai.meta.com/blog/meta-llama-3/, 2024. Accessed: 2024-07-22.
656 657 658 659	Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. <i>arXiv e-prints</i> , 2019.
660 661 662	Mingjie Sun, Zhuang Liu, Anna Bair, and J Zico Kolter. A simple and effective pruning approach for large language models. In <i>The Twelfth International Conference on Learning Representations</i> , 2024. URL https://openreview.net/forum?id=PxoFut3dWW.
663 664 665 666 667 668	Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Xin Zhao, Furu Wei, and Ji-Rong Wen. Language-specific neurons: The key to multilingual capabilities in large language models. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pp. 5701–5715, Bangkok, Thailand, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.acl-long.309.
669 670 671 672	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> , 2023.
673 674 675 676 677	Hetong Wang, Pasquale Minervini, and Edoardo Ponti. Probing the emergence of cross-lingual alignment during LLM training. In Lun-Wei Ku, Andre Martins, and Vivek Srikumar (eds.), <i>Findings of the Association for Computational Linguistics ACL 2024</i> , pp. 12159–12173, Bangkok, Thailand and virtual meeting, August 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.findings-acl.724.
678 679 680 681	Wenxuan Wang, Zhaopeng Tu, Chang Chen, Youliang Yuan, Jen tse Huang, Wenxiang Jiao, and Michael R. Lyu. All languages matter: On the multilingual safety of large language mod- els. CoRR, abs/2310.00905, 2023. URL https://doi.org/10.48550/arXiv.2310. 00905.
682 683 684	Miles Williams and Nikolaos Aletras. How does calibration data affect the post-training pruning and quantization of large language models?, 2023. URL https://arxiv.org/abs/2311.09755.
685 686 687 688 689 690 691 692	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. Transformers: State-of-the-art natural language processing. In Qun Liu and David Schlangen (eds.), <i>Proceedings of the 2020 Confer-</i> <i>ence on Empirical Methods in Natural Language Processing: System Demonstrations</i> , pp. 38– 45, Online, October 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. emnlp-demos.6. URL https://aclanthology.org/2020.emnlp-demos.6.
693 694 695 696 697 698 699	Zhihui Xie, Handong Zhao, Tong Yu, and Shuai Li. Discovering low-rank subspaces for language- agnostic multilingual representations. In Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (eds.), <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Pro-</i> <i>cessing</i> , pp. 5617–5633, Abu Dhabi, United Arab Emirates, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-main.379. URL https:// aclanthology.org/2022.emnlp-main.379.
700 701	Weihang Xu and Simon Du. Over-parameterization exponentially slows down gradient descent for learning a single neuron. In <i>The Thirty Sixth Annual Conference on Learning Theory</i> , pp. 1155– 1198. PMLR, 2023.

733

- Yuemei Xu, Ling Hu, Jiayi Zhao, Zihan Qiu, Yuqi Ye, and Hanwen Gu. A survey on multilingual large language models: Corpora, alignment, and bias. *arXiv preprint arXiv:2404.00929*, 2024.
- Geng Yuan, Xiaolong Ma, Wei Niu, Zhengang Li, Zhenglun Kong, Ning Liu, Yifan Gong, Zheng Zhan, Chaoyang He, Qing Jin, et al. Mest: Accurate and fast memory-economic sparse training framework on the edge. *Advances in Neural Information Processing Systems*, 34:20838–20850, 2021.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. Hellaswag: Can a machine really finish your sentence? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019.
- Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning requires rethinking generalization. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=Sy8gdB9xx.
- Yingtao Zhang, Jialin Zhao, Wenjing Wu, Alessandro Muscoloni, and Carlo Vittorio Cannistraci.
  Epitopological sparse ultra-deep learning: A brain-network topological theory carves communities in sparse and percolated hyperbolic anns. 2023.
- Yingtao Zhang, Haoli Bai, Haokun Lin, Jialin Zhao, Lu Hou, and Carlo Vittorio Cannistraci. Plugand-play: An efficient post-training pruning method for large language models. In *The Twelfth International Conference on Learning Representations*, 2024. URL https://openreview. net/forum?id=Tr0lPx9woF.
- Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. How do large language models handle multilingualism? *arXiv preprint arXiv:2402.18815*, 2024.
- Zhixue Zhao and Nikolaos Aletras. Comparing explanation faithfulness between multilingual and monolingual fine-tuned language models. In Kevin Duh, Helena Gomez, and Steven Bethard (eds.), Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers), pp. 3226–3244, Mexico City, Mexico, June 2024. Association for Computational Linguistics. URL https://aclanthology.org/2024.naacl-long.178.
  - Xunyu Zhu, Jian Li, Yong Liu, Can Ma, and Weiping Wang. A survey on model compression for large language models. *arXiv preprint arXiv:2308.07633*, 2023.

#### 756 LIMITATIONS А 757

758 759

**Generality of Findings.** Due to resource constraints, we predominantly experimented with the small versions of Llama-3 and Aya-23, and validated our findings with fewer pruning runs on their 760 counterpart large version. Since our results translate between model families, and to bigger model 761 sizes, we assume a certain degree of generalization. Nonetheless, other models trained with different 762 techniques or for other tasks might show different behavior. Given the pace of this research field, it 763 is also unclear whether these results translate to future models.

764

765 **Underrepresented Languages.** Our experiments focused on languages with sufficient model and 766 downstream task support. However, this selection does not encompass all languages of interest, par-767 ticularly mid and low-resource languages that are underrepresented in the pre-training, and challenging to evaluate due to the lack of benchmark support. Future research could benefit from including 768 more languages to explore the interplay between different language families or writing systems and 769 performance after pruning. 770

771 772

773 774

775

776

777

778

779

781 782

783

#### В CALIBRATION AND TEST LANGUAGES

Writing system Language Language family Modern Standard Arabic (AR) Afro-Asiatic Arabic German (DE) Germanic Latin English (EN) Germanic Latin Spanish (ES) Romance Latin Russian (RU) **Balto-Slavic** Cyrillic Swahili (SW) Atlantic-Congo Latin Chinese (simplified) (ZH) Sino-Tibetan Simplified Han

Table 5: Languages used for calibration and testing throughout this paper with their corresponding language family and writing system.

784 785 786

787 788

792

#### С DOWNSTREAM DATASETS

Throughout the paper we used the following widely employed datasets for automated benchmarking. 789 All evaluations were conducted in a zero-shot fashion and employ the chat-template of the respective 790 instruction-tuned model. 791

793 ARC: The AI2 Reasoning Challenge (ARC) dataset introduced by Clark et al. (2018) tests the reasoning and knowledge capabilities through natural, grad-school multiple choice science questions 794 originally authored for standardized human tests. The dataset comprises a total of 7787 questions in 795 english divided into a Challenge set (ARC-C) of hard to answer questions and an Easy set (ARC-E) 796 of questions. 797

798 For evaluation in English, we use the original datasets (e.g. ARC-c & ARC-e), for all other languages the translated version from Dac Lai et al. (2023) is utilized. 799

800

801 Belebele: This carefully curated dataset evaluates 4-way multiple-choice machine reading com-802 prehension among 122 language options, broadly focussing on high-, medium-, and low-resource languages (Bandarkar et al., 2023). Each of the 900 samples is based on an English FLORES-200 803 passage that has been translated into the respective target language by fluent expert speakers. Hence, 804 the dataset is fully parallel, allowing direct performance comparison across all languages. 805

806

807 HellaSwag: The HellaSwag dataset by Zellers et al. (2019) comprises 10042 english samples testing commonsense natural language inference on event descriptions that need to be contin-808 ued/completed in a multiple-choice fashion. Though easily answerable by humans, such paragraph 809 continuation questions still pose a challenge for state-of-the-art LLMs.

MKQA: Multilingual Knowledge Questions and Answers (MKQA) (Longpre et al., 2020) is an open-domain question-answering evaluation set of 10000 samples aligned across 26 languages by human translators. Its question-answer pairs were filtered from the Google Natural Questions dataset (Kwiatkowski et al., 2019), annotating real Google search user questions with answers found on Wikipedia. Given a question, the task is to predict the correct answer or give no answer without any additional context provided. Hence, this dataset tests the knowledge retrieval capabilities of models. For our evaluation, we remove all unanswerable and questions requiring overly long answers for simplicity, yielding a total 6758 remaining samples. 

There are a total of 6,758 evaluation samples included in our experiment, after eliminating those with unanswerable and long, imprecise answers.

MMLU: The Massive Multitask Language Understanding (MMLU) dataset (Hendrycks et al., 2021) is an english benchmark designed to evaluate a model's ability to handle diverse subjects across multiple domains. It contains a total of 14042 question-answer pairs covering 57 task categories, ranging from high school and college-level subjects to professional and specialized knowledge. Each task includes multiple-choice questions, and the dataset measures both the model's factual knowledge and reasoning abilities.

Translated Datasets from Okapi: The Okapi framework, introduced by Dac Lai et al. (2023),
 focuses on instruction tuning LLMs using reinforcement learning from human feedback (RLHF)
 across multiple languages. As part of its resources, it includes translated versions of the ARC,
 HellaSwag, and MMLU datasets, generated using ChatGPT. We leverage these translations to complement the evaluation of the original English datasets in multiple languages.

#### D MATH & SETTINGS

#### D.1 PRUNING ERROR AND SNR

We resort to the definition of Kuzmin et al. (2023) for computing the Pruning Error and SNR. Since originally both of them relate to model weights, which are input and therefore language independent, we compute Pruning Error and SNR with respect to the outputs of each layer, i.e. the hidden states. To cope with different hidden state magnitudes per layer, all hidden states are normalized layer-wise by their average vector norm, yielding the following pruning error:

$$\mathbb{E}[(\boldsymbol{H}^{(k)} - \widetilde{\boldsymbol{H}}^{(k)})^2] = \frac{1}{Nd} \sum_{i=1}^N \sum_{j=1}^d \left(\frac{h_{i,j}^{(k)} - \widetilde{h}_{i,j}^{(k)}}{\mu^{(k)}}\right)^2, \mu^{(k)} = \frac{1}{N} \sum_{i=1}^N ||\boldsymbol{h}_i^{(k)}||_2$$
(1)

with  $\mathbf{H}^{(k)} \in \mathbb{R}^{N \times d}$  denoting the hidden states of all N tokens in layer k before and  $\widetilde{\mathbf{H}}^{(k)} \in \mathbb{R}^{N \times d}$ after pruning, while  $h_{i,j}^{(k)}$  refers to the hidden state of the *j*-th feature element of the *i*-th token in  $H^{(k)}$ . Subsequently, the SNR of a single layer k is computed as

$$SNR_{dB}^{(k)} = 10 \cdot log_{10}(\mathbb{E}(\boldsymbol{H}^{(k)}) / \mathbb{E}[(\boldsymbol{H}^{(k)} - \widetilde{\boldsymbol{H}}^{(k)})^2]), \text{ with } \mathbb{E}(\boldsymbol{H}^{(k)}) = \frac{1}{Nd} \sum_{i=1}^N \sum_{j=1}^d \left(\frac{h_{i,j}^{(k)}}{\mu^{(k)}}\right)^2.$$
(2)

The final model SNR is the average over all layer-wise  $SNR_{dB}^{(k)}$ .

D.2 LOW-RANK SUBSPACE FOR LANGUAGE-AGNOSTIC REPRESENTATIONS (LSAR)

D.2.1 LOW-RANK SUBSPACE COMPUTATION

Embeddings posess language identity components that can be extracted via Low-rank Subspace
 for language-Agnostic Representations (LSAR) as introduced by Xie et al. (2022). They project
 sentence embeddings (i.e. averaged embeddings of a sentence/short input) into a low-dimensional

subspace via a projection matrix  $M_s \in \mathbb{R}^{d \times r}$ .  $M_s$  describes a low-rank subspace of the original multilingual latent space, spanned by r components. The following describes how to obtain  $M_s$ .

First, the mean embeddings  $\mu_l = \frac{1}{n} \sum_{i=1}^{n} e_l^i$  of each language l are computed to then concatenate all  $\mu_l$  column-wise into matrix  $M \in \mathbb{R}^{d \times L}$ . To capture the relationships among the mean embedding vectors in M, it is decomposed into the matrix  $M_s \in \mathbb{R}^{d \times r}$ , spanning the low-rank subspace for linguistic, language-specific signals, and a vector  $\mu$  orthogonal to Span $(M_s)$  representing shared, language-agnostic feature elements. This poses an optimization problem

$$\min_{\boldsymbol{\mu}, \boldsymbol{M}_s, \boldsymbol{\Gamma}} \| \boldsymbol{M} - \boldsymbol{\mu} \boldsymbol{1}^\top - \boldsymbol{M}_s \boldsymbol{\Gamma}^\top \|^2$$
  
s.t.  $\boldsymbol{\mu} \perp \operatorname{Span}(\boldsymbol{M}_s)$ 

with  $1 \in \mathbb{R}^d$  containing all ones, that can be solved optimally using a Singular Value Decomposition (SVD). The following omits the proof and only shows how to compute  $M_s$ , for more details please refer to the original paper by Xie et al. (2022).

 $M_s$  is obtained in a two step process. First, M is approximated in low-dimensional space by centering M with  $M_C = M - \mu' \mathbf{1}^T = M - \frac{1}{d}M\mathbf{1}$  and computing the SVD decomposition of  $M_C$  for the r largest singular values as follows:

$$M'_s, s, \Gamma' = \text{Top-}r\text{SVD}(M_C)$$
 (3)

 $M'_s$  and  $\Gamma'$  are the left and right singular vectors. Then M is reconstructed from  $M'_s$ , s and  $\Gamma$  via

 $M' = \mu' \mathbf{1}^T + M'_{\circ} (\Gamma' \mathbf{I}s)^T$ (4)

with  $\mathbf{I} \in \mathbb{R}^{r \times r}$  as identity matrix. The second step is to enforce the orthogonality constraint  $\mu \perp \text{Span}(M_s)$ . This involves re-centering M' again for SVD preparation to

$$M'_{C} = \mathbf{M}' - \boldsymbol{\mu} = \mathbf{M}' - \frac{1}{||\mathbf{M}'^{-1}\mathbf{1}||^{2}}\mathbf{M}'^{-1}\mathbf{1}.$$
(5)

Afterwards, the SVD for the reconstructed M' is computed, obtaining  $M_s$  as left singular values as follows:

$$M_{s, -}, \Gamma = \text{Top-}r\text{SVD}(M_C') \tag{6}$$

Finally,  $M_s$  is the projection matrix into the low-rank latent space of language-specific signals. Conversely,  $M_s$  allows extracting language-agnostic features for an embedding  $e_l$  by projection onto the null space of  $M_s$ . That is, projecting  $e_l$  into the low-rank subspace, retaining language-specific information in  $\mathbb{R}^r$  only, followed by re-projection into the original embedding space, obtaining the language-specific feature components in  $\mathbb{R}^d$ . Then, obtaining language-agnostic  $a_l$  runs as follows:

$$\boldsymbol{a}_{l} = (\boldsymbol{I} - \boldsymbol{M}_{s}(\boldsymbol{M}_{s}^{T}\boldsymbol{M}_{s})^{-1}\boldsymbol{M}_{s}^{T})\boldsymbol{e}_{l} = \boldsymbol{e}_{l} - \boldsymbol{M}_{s}\boldsymbol{M}_{s}^{T}\boldsymbol{e}_{l}$$
(7)

<sup>908</sup> with *I* as the identity matrix.

#### 910 D.2.2 LSAR EXPERIMENTAL SETTINGS

Since an analysis of token-embeddings via the Pruning Error and SNR does not predict downstream task performance well, we revert to sentence-embeddings or rather prompt-embeddings (i.e. the average of all token-embeddings, omitting special tokens from the chat template) to capture more high-level semantic and syntactic information. For separating language-agnostic from specific feature, we calculate a separate low-rank projection matrix  $M_s^{(m)}$  per model m. Since we use LSAR to obtain pruning error metrics and do not demand any generalization to unseen data,  $M_s^{(m)}$  is computed on the same samples used for evaluation.

## 918 D.3 PRUNING MASK COMPARISON VIA IOU

920 We compare the similarity of pruning masks for different calibration languages by IoU in a two step 921 process. First, we eliminate calibration set dependent noise by computing the intersection  $M_I^l$  of the 922 sets  $M_i^l$  of pruned neuron indices for different calibration sets  $i \in \{0, ..., N\}$  of the same language 923  $l \in \{0, ..., L\}$ . Then,  $M_i^l$  is computed as follows:

$$M_I^l = \bigcap_{i=0}^N M_i^l, \qquad M_U^l = \bigcup_{i=0}^N M_i^l$$
(8)

 $M_I^l$  only contains neuron indices that were pruned for all calibration sets and can therefore be deemed more stable. For comparing the similarity of pruning masks between different calibration languages  $l_1$  and  $l_2$  the IoU is utilized, running as follows:

$$IoU_{l_1,l_2} = \frac{\left| M_I^{l_1} \bigcap M_I^{l_2} \right|}{\left| M_I^{l_1} \bigcup M_I^{l_2} \right|} \tag{9}$$

For validation, in Figure 2 of the main text we also plot the  $IoU_l$  for pruning masks  $M_i^l$  of the same calibration language, visualizing pruning certainty.

#### D.4 LANGUAGE ACTIVATION PROBABILITY ENTROPY (LAPE) SCORE

In our analysis we utilize the LAPE score introduced by Tang et al. (2024) as a measure for languagespecific neuron activations in the FFN. Let the FFN of the *i*-th layer of a Llama-3 model be described by the following formula:

$$\boldsymbol{h}^{(i)} = (SwiGLU(\boldsymbol{h}_{attn}^{(i)} \cdot \boldsymbol{W}_1^{(i)}) \otimes \boldsymbol{W}_2^{(i)}) \cdot \boldsymbol{W}_3^{(i)},$$
(10)

where  $\boldsymbol{h}^{(i)} \in \mathbb{R}^d$  depicts the hidden state and  $\boldsymbol{h}_{attn}^{(i)}$  the attention output of layer *i* of a specific token, with  $\boldsymbol{W}_1^{(i)}, \boldsymbol{W}_2^{(i)} \in \mathbb{R}^{d \times 3,5d}$  as up and gate projection matrices,  $\boldsymbol{W}_3^{(i)} \in \mathbb{R}^{3,5d \times d}$  as down projection matrix and  $\otimes$  as element-wise multiplication.

For simplicity Tang et al. (2024) consider the *j*-th neuron of the FFN to be **activated** if the corresponding output of the SwiGLU activation function is greater than zero. Then, the **activation probability** of the *j*-th neuron in the *i*-th layer for input prompts in language k is defined as:

$$p_{i,j}^{(k)} = \mathbb{E}(\mathbf{1}(SwiGLU(h_{attn}^{(i)} \cdot \boldsymbol{W}_1^{(i)}))_j > 0 \mid \text{language } k), \tag{11}$$

with 1 as the identity function to count the times the neuron is activated and subsequently estimate the likelihood of neuron activation. Repeating this process for all languages and appyling L1 normalization yields the distribution  $\tilde{p}_{i,j} = (p_{i,j}^{(1)}, \dots, p_{i,j}^{(k)}, \dots, p_{i,j}^{(L)}) \div ||(p_{i,j}^{(1)}, \dots, p_{i,j}^{(k)}, \dots, p_{i,j}^{(L)})||_1$ . Finally, this allows for computing the **language activation probability entropy**:

$$LAPE_{i,j} = \sum_{k=1}^{L} \widetilde{p}_{i,j}^{(k)} \cdot \log(\widetilde{p}_{i,j}^{(k)})$$
(12)

A low LAPE score denotes neurons that are activated more often for inputs in specific languagesthan for all other languages.

Unlike Tang et al. (2024), who focus on the lowest 1% of LAPE scores and exclude those with negligible activation probability, our analysis seeks broader trends. Hence, by neglecting any thresholds we relax the original definition.





Table 6: Language-specific perplexity (PPL), Signal-to-noise ratio (SNR) and pruning error for pruned Llama 3 8B and Aya 23 8B models, averaged over three pruning runs. The leftmost columns show the model, the pruning technique, and the calibration language. For PPL and pruning errors, the smaller the value (the darker), the better; while for SNR, the greater the value (the lighter), the better. *Note the diverging color pattern for the Pruning Errors of the Aya-23 8B model pruned with Wanda. Even after repeating the evaluation several times the results look identical. We leave them in for transparency.* 



Table 7: Evaluation task performance of different model sizes (Llama-3 8B/70B and Aya-23 8B/35B) pruned for 50% unstructured sparsity. Due to resource constraints, only the 8B model variants were pruned with multiple seeds. The results are averaged over three pruning runs, with standard deviation given in the subscript. The leftmost columns show the model, the pruning technique, and the calibration language used for pruning. A "-" indicates the unpruned reference model. Each column shows the perplexity score of the pruned models on a specific evaluation language. For evaluation in English we use the original datasets (e.g. ARC-c & ARC-e), for all other languages the translated version from Dac Lai et al. (2023) is utilized. 



Table 10: Downstream task performance of Llama 3.1 pruned by SparseGPT and Wanda towards 50% unstructured sparsity.



RU ZH 34.8 AR DE 69.4

32.3

EN-C

#### 1134 F QUANTIZATION 1135

AR 31.7

28.1 29.2

AR DE EN ES RU SW

GPTQ

DE

33.7



1137

1138 1139 1140

Llama-3-8B- 
 45.0
 53.0
 23.5

 54.8
 46.3
 25.2

 56.9
 47.4
 33.1

 56.4
 60.6
 27.7
 32.5 32.1 55.7 61.2 50.3 59.1 
 35.7
 27.5
 36.0
 56.8
 40.6
 27.1

 38.2
 27.9
 34.8
 58.8
 39.8
 27.8
 24.2 24.9 ZH 34.4 41.2 52.5 1141 1142 Table 14: Downstream performance of Llama-3 8B after GPTQ Quantization to 4-bit weights and with a group-size of 128. 1143

Belebe

58.4

58.4

SW 47.3 ZH MMLU

DE 38.2 EN 58.6

35.

32 1

32.2

34.4 25.7

HellaSwa

52.3

ZH 25.9 AR 36.2 DE 43.2

34.6

			ARC[acc]							Belebele <sub>[acc]</sub>					MMLU <sub>[acc]</sub>					HellaSwag <sub>lacel</sub>							
			AR	DE	EN-E	EN-C	ES	RU	ZH	AR	DE	EN	ES	RU	SW	ZH	AR	DE	EN	ES	RU	ZH	AR	DE	EN	ES	RU
ICL	-	-	31.7	36.6	75.9	48.5	37.9	36.9	34.8	55.3	69.4	58.2	66.3	62.9	47.3	44.8	27.1	38.2	58.6	43.0	30.3	25.9	36.2	43.2	53.3	46.1	41.7
Ę		AR	31.3	36.9	75.8	48.4	38.0	37.3	35.1	56.1	70.2	60.1	68.3	62.9	47.7	48.2	27.6	38.6	59.0	43.6	30.9	26.3	36.2	43.2	53.4	46.1	41.7
ų.		DE	31.5	36.8	75.9	48.5	37.9	36.8	35.6	57.0	70.6	60.2	68.2	62.8	48.3	47.9	27.6	38.7	59.1	43.9	31.0	26.3	36.2	43.3	53.3	46.0	41.7
ė	ø	EN	31.4	36.4	75.9	48.6	38.3	37.2	35.3	55.9	70.0	59.4	66.8	62.3	47.0	46.2	27.4	38.1	59.1	43.0	30.5	26.0	36.2	43.2	53.3	46.1	41.6
2	£	ES	31.1	37.0	75.7	48.5	38.0	37.1	35.2	56.9	70.4	60.6	68.0	63.0	47.9	48.6	27.6	38.2	58.7	43.6	30.5	26.3	36.2	43.3	53.4	46.1	41.8
-81	5	RU	31.1	36.7	75.8	48.3	38.3	37.1	35.3	56.3	70.2	59.6	68.4	62.6	47.6	46.3	27.4	38.3	58.8	43.7	30.6	26.0	36.3	43.3	53.4	46.0	41.7
an		SW	31.3	36.9	75.9	48.5	38.1	37.5	35.5	56.6	70.0	61.0	68.3	63.0	47.3	46.7	27.6	38.6	59.4	44.1	31.0	26.2	36.3	43.3	53.3	46.1	41.7
Ц		ZH	31.5	37.0	75.9	48.5	38.5	37.6	35.2	56.7	70.1	60.8	68.1	63.4	47.0	47.2	27.6	38.4	59.2	43.7	30.7	26.2	36.2	43.3	53.3	46.1	41.6

Table 15: Downstream performance of Llama-3 8B after GPTQ Quantization to 8-bit weights and with a group-size of 128.

1152 1153 1154

1151

- 1155
- 1156 1157
- 1158
- 1159
- 1160

1161

1162

1163

1164 1165

1166

1167 1168

1169 1170

1171 1172

1173

1174 1175

1176

1177 1178

1179

1180 1181

1182

1183

1184 1185

1186



Figure 4: Language-wise mean magnitude of difference between the hidden states of a fullsized Llama3 8B base model and its 50% unstructured sparsity SparseGPT-pruned model for all calibration-evaluation language pairs. The test languages are shown on the x-axis, the calibration languages are color-coded (AR, DE, ES, EN, RU, SW, ZH). The background color indicates the magnitude of the maximum deviation. We use all the 900 samples from the fully parallel Belebele dataset (Bandarkar et al., 2023) as input, concatenating the context passage and question for each sample.

- 1239
- 1240
- 1241



Figure 5: Language-wise mean magnitude of difference between the prompt-wise and layer-wise averaged hidden states of a full-sized Llama3 8B reference model and its 50% unstructured sparsity SparseGPT-pruned model for all calibration-evaluation language pairs. The test languages are shown on the x-axis, the calibration languages are color-coded (AR, DE, ES, EN, RU, SW, ZH). The background color indicates the magnitude of the maximum deviation. We use all the 900 samples from the fully parallel Belebele dataset (Bandarkar et al., 2023) as input, concatenating the context passage and question for each sample.



Figure 6: Language-wise mean magnitude of difference between the prompt-wise and layer-wise averaged *language-agnostic* features extracted with LSAR for an full-sized and pruned (50% unstructured sparsity, SparseGPT) Llama 3 8b model. Both, the LSAR projection matrix and the feature differences, were computed per model over all samples of the Belebele dataset for the seven calibration/test languages. The gray-scaled background represents the y-axis magnitude. The test languages are shown on the x-axis, the calibration languages are color-coded (AR, DE, ES, EN, RU, SW, ZH). The background color indicates the magnitude of the maximum deviation.



Figure 7: Language-wise mean magnitude of difference between the prompt-wise and layer-wise averaged *language-specific* features extracted with LSAR for an full-sized and pruned (50% unstructured sparsity, SparseGPT) Llama 3 8b model. Both, the LSAR projection matrix and the feature differences, were computed per model over all samples of the Belebele dataset for the seven calibration/test languages, concatenating the context passage and question for each sample. The test languages are shown on the x-axis, the calibration languages are color-coded (AR, DE, ES, EN, RU, SW, ZH). The background color indicates the magnitude of the maximum deviation.

1404						
1405						
1406						
1407						
1408						
1400						
1409						
1410						
1411						
1412		Subcomponen	t DE	EN	ZH	
1413	-	Q	0.890	0.880	0.874	
1414		Κ	0.891	0.880	0.875	
1415		V	0.882	0.871	0.865	
1416		attn.out	0.845	0.832	0.825	
1417		ffn.up	0.875	0.863	0.849	
1418		ffn.gate	0.873	0.862	0.847	
1419		ffn.down	0.870	0.852	0.843	
1420			0.070	01002	01010	
1421	Table 16: Average IoU of	all layers for a spe	ecific subc	omponent	. IoUs are	computed for the pruning
1422	masks of three pruning ru	ns on the same cal	libration la	anguage (I	DE, EN or	ZH).
1423						
1424						
1425						
1426						
1427						
1428						
1429						
1430						
1431						
1432						
1433						
1434						
1435						
1436						
1/137						
1/138						
1/130						
1435	-	Subcomponent	EN-DE	ZH-DE	ZH-EN	-
1440	-	Q	0.886	0.855	0.857	-
1441		K	0.886	0.855	0.857	
1442		V	0.879	0.846	0.848	
1443		attn.out	0.867	0.826	0.824	
1444		ffn.up	0.871	0.833	0.835	
1445		ffn down	0.870	0.852	0.855	
1446			0.840	0.808	0.011	
1447	Table 17: Average IoU of	f all layers for a s	pecific sul	ocompone	nt. IoUs a	are computed between the
1448	intersected pruning masks	of three pruning r	uns for the	same cali	bration la	nguage and the intersected
1449	pruning mask of another c	alibration languag	ge (EN-DI	E, ZH-DE	and ZH-E	Ň).
1450						
1451						
1452						
1453						
1454						





pruned (blue) Llama 3 8b model (50% unstructured sparsity, Wanda)

