

Want Better Synthetic Data? Steer It: Activation Steering for Low-Resource Language Generation

Jan Cegin[♣], Daniil Gurgurov[†], Yusser Al Ghussin[†], Simon Ostermann[†]

[♣] Kempelen Institute of Intelligent Technologies, Bratislava, Slovakia

jan.cegin@kinit.sk

[†] German Research Institute for Artificial Intelligence (DFKI), Saarbrücken, Germany

{daniil.gurgurov, yusser.al_ghussin, simon.ostermann}@dfki.de

Abstract

Large language models (LLMs) have become an effective tool for synthetic data generation, including for low-resource languages, where generated data can improve downstream task performance. Current best-performing approaches typically rely on few-shot prompting with target-language examples, which increases inference costs and may reduce diversity through lexical anchoring. In this work, we investigate activation steering as an alternative for low-resource synthetic data generation. We study two steering strategies: *Language Steering*, which targets the linguistic identity of a language, and *Quality Steering*, which captures well-formedness by contrasting human-written and backtranslated text representations. We evaluate these methods across four open-source LLMs, multiple layers, and 11 typologically diverse languages by generating sentiment and topic classification data and finetuning smaller classifiers. Steering is applied in both zero-shot and few-shot prompting settings and compared against non-steered counterparts. Our results show that steering on early layers consistently improves the diversity of generated data while often yielding stronger downstream model performance, particularly for low-resource languages.

1 Introduction

With the emergence of LLMs, various studies have demonstrated their ability to generate label-adherent and well-formed text (Cegin et al., 2023; Ubani et al., 2023). These two capabilities have made them an ideal tool for data *augmentation* and synthetic data *generation*, as demonstrated across domains such as sentiment analysis (Onan, 2023), intent recognition (Cegin et al., 2024a), and topic classification (Piedboeuf and Langlais, 2023). For synthetic data generation, the workflow usually consists of prompting an LLM to produce a set amount of samples with given labels in a particular

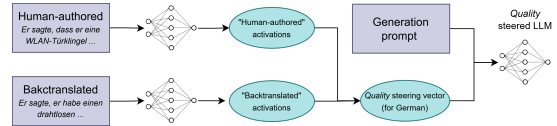


Figure 1: Example showing how *Quality* steering vectors are created. This example shows a contrastive collection of activations for German, from which a steering vector is created. This steering vector is used before generating data from the target LLM and language.

language. These are then often used for fine-tuning of downstream encoder models.

Recent works have extended these approaches to low-resource languages, often improving downstream model performance (Anikina et al., 2025; Pranida et al., 2025; Cegin et al., 2026), where data scarcity is prominent (Joshi et al., 2020). The strongest results generally come from few-shot prompting with examples from the target language (Anikina et al., 2025). However, this increases inference costs and may reduce diversity through lexical anchoring.

In parallel, research on representation engineering or activation steering has shown that modifying internal model activations can steer LLM behavior toward desired semantic or stylistic properties (Zou et al., 2023; Turner et al., 2024). By injecting steering vectors into hidden states, prior work has controlled attributes such as truthfulness (Li et al., 2024), sentiment (Rimsky et al., 2024), and toxicity (Turner et al., 2024) without additional finetuning. Compared to few-shot prompting, activation steering provides an efficient alternative that may better capture various abstract properties (Ostermann et al., 2026). Additionally, it can be combined with any type of prompt, zero- or few-shot (Radevski et al., 2026).

In this paper, we investigate activation steering for low-resource synthetic data generation. We generate synthetic datasets for sentiment and topic classification tasks across 11 typologically diverse

languages and evaluate them through downstream finetuning performance. We study two steering strategies: *Language Steering*, which targets linguistic identity per previous studies (Gurgurov et al., 2026, 2025b; Ghussin et al., 2026a,b), and *Quality Steering*, which isolates well-formedness by contrasting human-written and backtranslated text. To our knowledge, this is the first work to derive and use *Quality* steering vectors. Across four open-source LLMs and multiple layers, we apply these steering vectors to both zero- and few-shot prompts (Anikina et al., 2025) and analyze the resulting data in terms of diversity and representational properties, compared to no-steering baselines. Code can be found at <https://github.com/kinit-sk/steering-synth-gen>.

Our main contributions are:

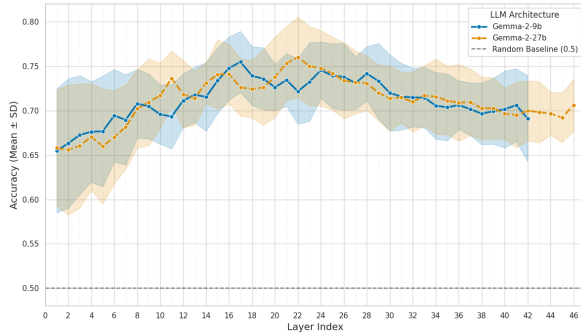
- For the first time, we apply *Language* steering vectors for synthetic data generation, and introduce and utilize *Quality* steering vectors by contrasting human-authored texts with backtranslated texts.
- We analyze the cosine similarity between *Quality* and *Language* steering vectors, demonstrating that their alignment is highly LLM-dependent and most polarized in earlier layers. Crucially, the majority of evaluated languages display strong negative similarity, implying that, in general, *Quality* steering vectors are not tied to a specific language.
- We show that applying *Quality* and *Language* steering to early layers improves LLM performance in generating synthetic data, as measured by downstream performance, while increasing the diversity of the generated data for both zero- and few-shot prompting.

2 Related Work

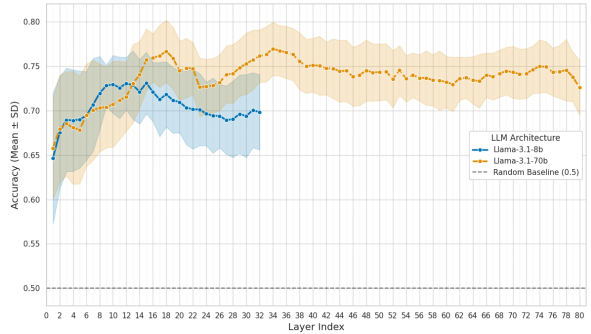
Synthetic data. LLMs are increasingly being used to create semantically new samples adhering to a given label (Ubani et al., 2023; Cegin et al., 2024b). LLM-based augmentation and data generation have been used for a variety of tasks such as automated scoring (Fang et al., 2023), low-resource language generation (Ghosh et al., 2023), intent classification (Sahu et al., 2022), sentiment analysis (Piedboeuf and Langlais, 2023; Onan, 2023), content recommendation (Liu et al., 2024), and health symptoms classifications (Dai et al.,

2023). As LLMs are better generators than classifiers of data in low-resource languages (Pecher et al., 2026), recent studies have focused on generating such label-adhering data in a variety of low-resource languages like Vietnamese (Feng et al., 2021), Marathi (Tran et al., 2026), or various African languages (Belay et al., 2026), and for a variety of tasks and domains such as QA (Namboori et al., 2023), fact-checking (Chung et al., 2025), NER (Liu et al., 2021), or text classification (Glenn et al., 2023). Other studies focused on finding the best approaches for generating data in low-resource languages, such as enhancing generator selection (Cegin et al., 2026) or finding the best prompting techniques (Anikina et al., 2025). As such, the current state-of-the-art technique (Anikina et al., 2025) uses few-shot examples in the prompt itself, leading to increased performance for additional inference costs.

Activation steering. A separate line of work has explored activation steering, where model behavior is modified by intervening directly on intermediate representations during the forward pass rather than through prompting or fine-tuning (Zou et al., 2023; Turner et al., 2024; Ostermann et al., 2026). The most widely used family of methods is difference-based steering, which derives steering vectors from contrastive examples (Turner et al., 2023; Rimsky et al., 2024; Marks and Tegmark, 2023). Beyond this, optimization-based approaches such as ReFT (Wu et al., 2024) and its rank-1 variant (Wu et al., 2025) learn low-rank interventions on hidden states, while dictionary-learning methods use sparse autoencoders to decompose activations into interpretable, individually steerable features (Gao et al., 2024; Cunningham et al., 2024). These techniques have been shown to reliably control a wide range of behaviors, such as sentiment, topic, and style (Turner et al., 2023; Konen et al., 2024), truthfulness and sycophancy (Li et al., 2023; Panickssery et al., 2023), and refusal (Arditi et al., 2024). Particularly relevant to our setting, steering has also been applied to multilingual control: Previous work has identified language-specific neurons (Zhao et al., 2024; Tang et al., 2024; Gurgurov et al., 2025b), SAE features (Chou et al., 2025; Deng et al., 2025), and language vectors (Wong et al., 2026; Gurgurov et al., 2026; Ghussin et al., 2026a) that can be activated or ablated to change the output language while preserving semantics quality. Recent work has further shown that language-vector steering can



(a) Gemma 2 models results.



(b) Llama 3.1 models results.

Figure 2: Linear probing results for human-authored vs. backtranslated textual pairs aggregated over languages.

be used for culturally aware multilingual inference: Ghussin et al. (2026b) reporting improvements for cultural knowledge retrieval. While these studies establish that linguistic identity is encoded as a steerable direction in activation space, they focus on language control or downstream improvement on cultural benchmarks rather than the *generation quality* of the produced text. Our work addresses this gap by using both language and quality steering vectors to produce better synthetic data for low-resource languages.

3 Methodology and Experiments

Our methodology consists of a simple procedure where we (1) collect activations from the residual stream from given layers and create the steering vectors from the collected activations, (2) apply the steering vector to the LLM, (3) generate data with labels in a given language, and (4) finetune a downstream model and evaluate it on human test data. This is done for 11 different languages.

These 11 languages in our study are selected to cover diverse typological properties and ensure overlap with the evaluation datasets. They span Indo-European (Germanic: Danish, German; Slavic: Czech, Slovak, Slovenian), Afro-Asiatic (Semitic: Amharic, Hebrew, Maltese), and Austronesian (Indonesian, Javanese, Sundanese) families (Nordhoff and Hammarström, 2011). This design allows us to test steering robustness across varied morphological systems, from agglutinative to fusional structures, as well as across different writing systems (Ge’ez, Hebrew, Latin). This diversity supports evaluating whether the learned steering manifold generalizes beyond orthographic form.

3.1 Computing Steering Vectors

Language Vectors. We use the FLORES (NLLB Team et al., 2024) and BOUQUET (Andrews et al., 2025) datasets, both of which contain human-authored multilingual texts. FLORES provides professionally translated parallel sentences across all 11 target languages listed in Appendix E, enabling control over semantic content while isolating linguistic identity (Ghussin et al., 2026b). BOUQUET is added to increase lexical diversity and include more naturalistic language. Concatenating both datasets also ensures sufficient token coverage for stable mean activations.

To construct *Language* steering vectors, we employ a one-vs-rest contrastive approach over the residual stream activations of all transformer blocks (Ghussin et al., 2026b). For each of the 11 target languages, we compute a mean activation vector from the dataset by averaging across all non-padding tokens, yielding a language-specific mean representation in the model’s latent space for every layer. The steering vector is then defined as the difference between the target language mean representation and the mean representation of all remaining languages, followed by normalization.

Quality Vectors. We follow the contrastive activation extraction approach of Rimsky et al. (2024), which requires paired examples representing opposing properties (in our case: human vs. generated samples). To construct these pairs, we generate backtranslations for FLORES and BOUQUET using the distilled NLLB model (Koishekenov et al., 2023). We treat the original human-authored text as the desired representation and the backtranslated text as its contrastive counterpart, following prior findings that synthetic text is generally lower quality than human-authored data (Schaffelder and

Table 1: Mean F1 difference across languages, models, steering methods, and layers aggregated over tasks, with **positive** and **negative** differences compared to a **zero-shot prompt** with no steering.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
<i>Language</i>	Gemma 2 27b	L10 ($\alpha = 100.0$)	0.17	0.28	-2.10	2.24	-1.16	2.49	1.84	-0.56	1.52	2.66	2.33	0.88
		L22 ($\alpha = 75.0$)	2.36	-0.22	-0.13	0.73	-0.52	0.77	-2.67	0.38	1.83	0.07	4.92	0.68
		L34 ($\alpha = 75.0$)	4.34	-0.51	-0.43	-1.37	-0.36	2.17	1.96	-1.71	-1.72	-0.07	6.20	0.77
	Gemma 2 9b	L9 ($\alpha = 50.0$)	-3.25	1.48	1.09	0.43	2.29	0.17	-3.63	0.45	4.08	6.88	0.16	0.92
		L20 ($\alpha = 50.0$)	-1.72	-0.25	1.88	-1.04	-3.54	-0.37	-0.86	2.75	6.32	3.76	3.28	0.93
		L31 ($\alpha = 25.0$)	2.30	-0.70	1.69	0.50	1.73	-0.46	-2.22	2.96	-0.59	5.62	2.52	1.21
	Llama3.1 70b	L17 ($\alpha = 2.0$)	3.25	-3.03	-0.91	2.62	0.98	1.19	2.42	0.79	9.67	1.27	-0.57	1.61
		L37 ($\alpha = 3.0$)	3.34	-0.75	1.39	3.43	0.66	0.97	4.10	0.53	4.82	2.22	0.85	1.96
		L57 ($\alpha = 2.0$)	2.97	-0.94	1.69	0.44	1.82	1.89	-0.97	1.49	1.87	3.35	1.56	1.38
	Llama3.1 8b	L7 ($\alpha = 2.0$)	8.39	-1.20	1.18	0.77	-0.77	1.52	-0.43	1.68	4.01	3.01	-0.24	1.63
		L15 ($\alpha = 1.0$)	4.52	1.25	2.13	1.52	-0.38	-0.62	0.57	0.38	2.32	2.11	-2.48	1.03
		L23 ($\alpha = 1.0$)	6.43	1.16	1.91	-0.63	-2.51	1.77	-0.08	2.03	4.36	0.69	-0.59	1.32
<i>Quality</i>	Gemma 2 27b	L10 ($\alpha = 100.0$)	3.01	0.16	-0.02	2.32	-0.91	0.49	-0.80	0.67	1.92	3.66	4.91	1.40
		L22 ($\alpha = 100.0$)	-0.46	-0.31	-1.54	-0.39	0.05	0.98	1.79	1.34	1.04	1.18	3.44	0.65
		L34 ($\alpha = 50.0$)	-0.32	-0.35	-1.03	1.45	-0.60	0.43	-0.94	-1.82	3.83	2.77	3.85	0.66
	Gemma 2 9b	L9 ($\alpha = 75.0$)	-0.01	5.20	1.29	1.75	4.62	1.71	-0.10	1.28	11.70	9.12	4.51	3.73
		L20 ($\alpha = 50.0$)	3.34	0.36	3.31	2.12	-0.14	1.54	-0.42	1.47	1.57	5.04	1.60	1.80
		L31 ($\alpha = 75.0$)	2.25	3.86	1.08	0.88	-0.06	0.89	-0.96	2.05	9.07	1.86	0.72	1.97
	Llama3.1 70b	L17 ($\alpha = 3.0$)	2.58	0.59	0.37	0.25	0.65	0.54	-1.24	2.31	7.04	-0.44	1.65	1.30
		L37 ($\alpha = 1.0$)	4.06	-1.74	-1.65	-1.01	1.44	0.83	-4.45	1.75	4.64	1.71	-0.26	0.48
		L57 ($\alpha = 2.0$)	4.47	0.32	0.72	0.89	0.56	0.14	-1.09	1.30	2.43	-0.29	0.90	0.94
	Llama3.1 8b	L7 ($\alpha = 2.0$)	7.49	2.91	2.08	0.77	0.82	1.45	1.79	3.20	3.14	3.03	-0.11	2.42
		L15 ($\alpha = 3.0$)	9.19	0.64	1.05	-1.57	-1.00	-0.63	0.47	4.06	1.01	1.62	1.20	1.46
		L23 ($\alpha = 2.0$)	10.85	2.01	1.17	-0.35	-3.06	1.37	-0.49	3.18	3.41	2.25	0.20	1.87

Gatt, 2026). Each *Quality* steering vector is created language-specifically from the contrastive activations of human-authored and backtranslated texts from that given language. The aim is to capture subtle representational differences between human-written and synthetic text. An example is shown in Figure 1. All steering vectors are derived from task-agnostic data. Additional details are provided in Appendix F.

To construct *Quality* steering vectors, we use TransformerLens (Nanda and Bloom, 2022) to extract residual stream activations from transformer blocks without further training. We compute token-wise global averages for each example and define the steering vector as the difference between the mean activations of the contrastive pairs. Each of the *Quality* steering vectors is computed for a specific language (from data from that specific language). Details can be found in Appendix D.

3.2 Using Steering Vectors for Generation

To evaluate the effectiveness of steering vectors for synthetic data generation, we apply them at equivalent relative depths across four instruction-tuned LLMs: Gemma-2-9B, Gemma-2-27B, Llama-3.1-8B, and Llama-3.1-70B (Team et al., 2024; AI@Meta, 2024). Steering is applied at approximately 21%, 48%, and 74% of model depth, corresponding to early, middle, and late processing stages in the Transformer hierarchy (Bartoszcze et al., 2025). Early layers should be primarily as-

sociated with the consolidation of syntactic and linguistic identity (Tenney et al., 2019); the middle layers should represent the peak of conceptual abstraction (Geva et al., 2021); and the later layers, where the model transitions from abstract conceptual processing toward task-specific refinement and next-token probability mapping (Belrose et al., 2025). These layers were additionally selected based on stable linear probe performance in distinguishing human-authored from synthetic text (Section 4.1). Further details are in Appendix I.

We further investigate the effect of steering strength α . For Gemma models, we evaluate $\alpha \in \{25, 50, 75, 100\}$, while for Llama models we use $\alpha \in \{1, 2, 3, 4\}$. We observe substantial differences in sensitivity between the model families: Llama models frequently collapse at higher α values (e.g., repetition or empty outputs), whereas Gemma models require larger intervention strengths, remaining stable even at $\alpha = 100$. We attribute this to architectural differences, particularly Gemma-2’s use of logit soft-capping and Query-Key normalization (Team et al., 2024), as well as its larger residual stream norms.

3.3 Finetuning Downstream Models

In our experiments, downstream model performance serves as the primary indicator of synthetic data quality, following prior work (Anikina et al., 2025; Cegin et al., 2026). We evaluate the effect of *Language* and *Quality* steering on synthetic data

Table 2: Mean F1 difference across languages, models, steering methods, and layers aggregated over tasks, with positive and negative differences compared to a **few-shot prompt** with no steering.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Language	Gemma 2 27b	L10 ($\alpha = 25.0$)	1.58	2.48	0.63	1.73	3.25	1.36	-0.96	2.27	-0.39	3.55	-1.76	1.25
		L22 ($\alpha = 50.0$)	-0.64	2.48	0.00	-0.59	1.85	1.04	-1.19	0.00	-0.20	5.51	-1.35	0.63
		L34 ($\alpha = 25.0$)	0.97	2.58	0.34	0.35	-0.26	2.29	-0.69	1.66	-0.75	2.32	-0.87	0.72
	Gemma 2 9b	L9 ($\alpha = 25.0$)	-4.14	0.34	0.24	1.48	-0.54	-1.54	0.51	-1.97	2.06	-1.16	-0.49	-0.47
		L20 ($\alpha = 50.0$)	-0.17	0.96	1.19	0.87	-0.02	-0.21	0.97	-2.13	3.36	-0.47	-1.45	0.26
		L31 ($\alpha = 50.0$)	-2.35	0.25	3.99	-0.77	0.02	0.03	0.20	-1.00	2.44	-1.08	0.49	0.20
	Llama3.1 70b	L17 ($\alpha = 1.0$)	1.85	-0.23	-0.13	0.91	1.32	-0.08	0.32	0.30	-1.20	0.43	2.04	0.50
		L37 ($\alpha = 2.0$)	-0.70	1.28	0.89	0.37	3.46	1.26	0.08	2.21	1.61	0.51	2.85	1.26
		L57 ($\alpha = 1.0$)	0.63	0.68	0.82	0.93	2.96	1.41	-0.24	0.21	0.06	1.42	3.91	1.16
	Llama3.1 8b	L7 ($\alpha = 2.0$)	-1.23	-0.24	2.50	1.23	0.33	1.49	-1.66	2.95	-1.21	1.80	-2.25	0.34
		L15 ($\alpha = 2.0$)	1.39	0.12	-1.09	1.32	-1.07	-0.63	0.70	2.36	-0.38	1.83	-0.90	0.33
		L23 ($\alpha = 3.0$)	-0.14	-0.07	1.77	1.73	0.35	2.61	0.70	0.13	-0.17	1.10	-1.28	0.61
Quality	Gemma 2 27b	L10 ($\alpha = 75.0$)	-0.45	2.18	2.52	1.25	0.72	2.26	-0.28	2.11	2.31	4.73	-1.30	1.46
		L22 ($\alpha = 75.0$)	0.04	3.68	-0.63	2.06	-0.45	2.66	-0.95	0.81	0.99	3.32	-1.72	0.89
		L34 ($\alpha = 25.0$)	0.97	2.58	0.34	0.35	-0.26	2.29	-0.69	1.66	-0.75	2.32	-0.87	0.72
	Gemma 2 9b	L9 ($\alpha = 25.0$)	0.74	0.87	2.75	-0.20	1.75	0.37	1.57	3.19	0.61	0.78	-0.09	1.12
		L20 ($\alpha = 25.0$)	-3.42	-0.72	1.07	-0.05	0.01	1.31	1.13	0.48	-0.24	-0.02	-1.92	-0.22
		L31 ($\alpha = 50.0$)	-1.51	-0.45	2.68	1.31	0.27	0.62	-2.70	-1.20	3.75	-0.67	-0.11	0.18
	Llama3.1 70b	L17 ($\alpha = 1.0$)	1.55	0.68	0.31	1.45	2.83	0.74	1.37	-0.46	1.79	1.12	2.51	1.26
		L37 ($\alpha = 1.0$)	-0.70	0.70	0.80	-0.03	2.84	-0.58	-0.28	-0.44	0.42	-0.09	4.77	0.67
		L57 ($\alpha = 2.0$)	-1.40	1.02	0.46	-0.90	1.25	-0.52	0.45	2.44	1.89	0.43	5.48	0.96
	Llama3.1 8b	L7 ($\alpha = 2.0$)	1.03	1.82	1.70	1.03	0.72	1.00	-1.36	3.13	0.04	-0.67	0.05	0.77
		L15 ($\alpha = 3.0$)	1.49	0.32	-0.08	1.56	-1.37	-0.03	-2.39	3.06	-6.14	-0.39	-0.90	-0.44
		L23 ($\alpha = 4.0$)	4.05	-0.43	-1.72	1.49	-1.53	0.71	-1.08	2.97	-2.02	-0.71	0.20	0.17

for topic classification (multi-class) and sentiment analysis (binary). Due to the limited availability of multilingual benchmarks, we use SIB-200 (Adelani et al., 2024) and the sentiment dataset collection of (Brychcín and Habernal, 2013; Gurgurov et al., 2024, 2025a). For each combination of LLM, language, layer, α , and task, we generate synthetic datasets, producing 50 samples per label for sentiment and 20 samples per label for topic classification, similar to prior setups (Anikina et al., 2025).

For downstream evaluation, we fine-tune XLM-R (Conneau et al., 2019) (*facebookAI/xlm-roberta-base*) with early stopping. We compare steering-based generation against two baselines: (1) zero-shot prompting without demonstrations, and (2) state-of-the-art few-shot prompting with human demonstrations (Anikina et al., 2025). Prompt templates are found in Appendix H, with further implementation details provided in Appendix G.

4 Results and Discussion

4.1 Pre-Experiment: Linear Probing Results

We first conduct a linear probe analysis of human-authored (Flores + BOUQET) vs. back-translated pairs. This diagnostic step is done to justify our use of linear steering vectors for the *quality* approach; if a simple logistic regression can distinguish between these distributions with high accuracy, it confirms that the “quality” concept may be encoded as a specific direction in the latent space that we

can manipulate (Park et al., 2024). The details on how the linear probe is computed can be found in Appendix C.

As seen in the aggregated results over all languages in Figure 2, all tested models across both the Gemma-2 and Llama-3.1 families exhibit accuracies significantly above the 0.5 random baseline from the very first layer. While the probe accuracy peaks in the middle-to-late layers, the “quality” signal is already well-defined by the earlier layers. The Gemma-2-9b model peaks at Layer 17 with an accuracy of approximately 0.75, while the larger Gemma-2-27b model peaks at Layer 22 with a slightly higher accuracy of approximately 0.76. The Llama-3.1-70b model demonstrates superior separability compared to its 8b counterpart, reaching a sustained plateau of approximately 0.77 accuracy between layers 35 and 40.

4.2 Downstream Model Performance Evaluation

We report the results for only the best alpha per layer in terms of downstream model performance over both tasks combined. Additional results about how different alpha values affect downstream model performance can be found in the Appendix K. For statistical tests, we used Mann-Whitney-U (Mann and Whitney, 1947) with $p=0.05$. Full F1 baseline results for both zero- and few-shot settings can be found in Appendix J.

Table 3: Relative diversity metric differences across models, steering methods, and layers, aggregated over tasks and different alphas. Relative **positive** and **negative** changes.

(a) Zero-shot setup						(b) Few-shot setup					
Model + Steering	Layer	LexDiv	EmbDiv	IsoRad.	Homogen.	Model + Steering	Layer	LexDiv	EmbDiv	IsoRad.	Homogen.
Gemma 2 27b (lang.)	L10	0.06	-1.35	-1.02	-1.24	Gemma 2 27b (lang.)	L10	0.33	-3.66	0.38	1.22
	L22	-0.00	-2.16	-0.98	0.55		L22	0.92	-4.92	0.15	0.26
	L34	0.42	-1.17	-0.26	3.06		L34	0.99	-2.61	0.25	0.11
Gemma 2 27b (qual.)	L10	0.83	-0.01	-0.62	1.45	Gemma 2 27b (qual.)	L10	1.37	-3.48	0.33	1.00
	L22	0.30	-1.55	-0.92	-0.51		L22	1.51	-4.45	0.32	0.43
	L34	-0.16	-2.38	-0.60	1.96		L34	0.75	-5.62	0.49	0.65
Gemma 2 9b (lang.)	L9	2.12	23.51	1.81	12.50	Gemma 2 9b (lang.)	L9	-10.57	11.69	0.30	-4.56
	L20	10.72	13.57	1.24	6.07		L20	6.77	11.91	0.09	-7.19
	L31	7.68	16.54	0.99	4.80		L31	4.13	2.74	0.24	-2.80
Gemma 2 9b (qual.)	L9	21.78	52.83	4.09	11.33	Gemma 2 9b (qual.)	L9	10.82	23.05	1.29	0.19
	L20	17.07	18.22	0.10	5.52		L20	5.98	22.17	0.95	0.59
	L31	22.15	32.59	0.66	6.80		L31	4.15	14.53	0.20	-1.90
Llama3.1 70b (lang.)	L17	-29.03	55.32	1.93	13.13	Llama3.1 70b (lang.)	L17	-29.66	43.96	1.37	5.36
	L37	-2.43	-2.94	-0.85	1.60		L37	-1.56	4.48	0.52	2.40
	L57	-0.97	-4.14	-1.61	-0.47		L57	-1.35	4.83	0.29	1.21
Llama3.1 70b (qual.)	L17	14.98	28.64	1.90	6.91	Llama3.1 70b (qual.)	L17	8.15	15.55	0.94	2.62
	L37	10.33	8.37	-0.06	7.26		L37	10.10	12.51	1.56	5.66
	L57	2.40	-6.37	-1.57	1.18		L57	1.39	0.10	0.16	0.37
Llama3.1 8b (lang.)	L7	69.28	68.10	3.24	5.77	Llama3.1 8b (lang.)	L7	43.19	57.04	1.45	-6.83
	L15	44.03	58.97	2.45	3.34		L15	5.40	39.75	0.53	-0.90
	L23	32.80	20.41	2.22	5.27		L23	-6.11	11.37	0.01	-1.06
Llama3.1 8b (qual.)	L7	43.08	35.96	3.35	8.95	Llama3.1 8b (qual.)	L7	14.00	1.84	-0.08	-2.83
	L15	33.73	-14.71	1.21	5.63		L15	-17.68	-44.81	-1.61	3.75
	L23	36.48	-2.27	2.28	7.73		L23	-11.14	-30.27	-1.09	-0.16

Steering vectors applied with zero-shot prompts

Aggregated results across tasks are shown in Table 1, with per-task results provided in Appendix L. Overall, *Quality* steering consistently outperformed *Language* steering in terms of downstream F1 gains, suggesting that steering toward a “human-authored” manifold is more beneficial than steering toward generic linguistic identity alone.

Earlier layers proved the most effective intervention points, particularly for *Quality* steering. Across all task-aggregated experiments, *Quality* steering improved performance in 79.54% of cases for early layers, compared to 68.18% and 70.45% for middle and late layers. *Language* steering showed a similar trend, with gains in 72.73%, 70.45%, and 63.64% of cases, respectively. Statistically significant improvements for *Quality* steering occurred in 44.32% of cases for early layers, versus 27.23% and 31.82% for middle and late layers. For *Language* steering, statistically significant gains were observed in 30.68%, 31.82%, and 29.55% of cases, respectively. Per-task results are provided in Appendix L.

Across languages, Slovak and Slovenian showed the largest gains, while Czech and Danish were comparatively resistant to steering, particularly for Gemma models. Javanese was the only language with more negative than positive outcomes. Among models, Gemma-2-9B and Llama-3.1-8B were the

most responsive to steering, producing the largest absolute F1 improvements. Overall, Llama models benefited more consistently from steering, whereas Gemma-2-27B appeared more conservative, suggesting that larger models may require stronger or more precise interventions. Conversely, Llama3.1 70b demonstrated that while large models can be steered successfully, they are also prone to performance drops if the intervention is poorly aligned with the target layer.

Steering vectors applied with few-shot prompts

The aggregated few-shot results are shown in Table 2, with task-specific results provided in Appendix L. Compared to the zero-shot setting, the advantage of *Quality* steering over *Language* steering becomes less pronounced, and the strong preference for early-layer interventions weakens when demonstrations are included in the prompt. Nevertheless, early-layer *Quality* steering remains the most consistent approach overall.

Quality steering was most effective at early layers, improving downstream performance in 81.82% of cases, compared to 50% and 56.82% for middle and late layers. *Language* steering showed a similar trend, with gains in 72.73%, 70.45%, and 63.64% of cases, respectively. Average gains decreased for most models, except for Llama-3.1-70B, which appeared to benefit from combining

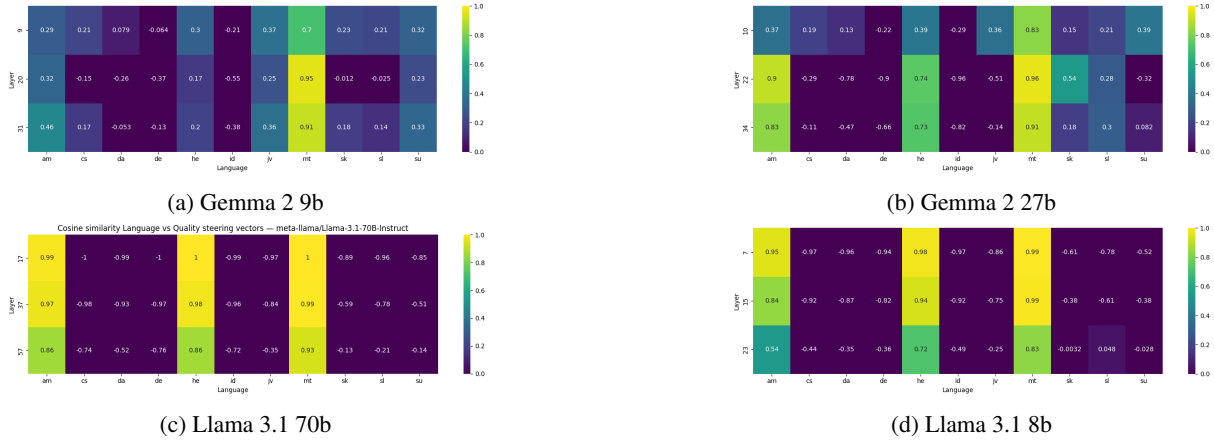


Figure 3: Cosine similarity between quality and language vectors for different LLMs used in this study.

steering with few-shot demonstrations. Statistically significant improvements for *Quality* steering occurred in 31.82% of cases for early layers, versus 14.77% and 18.18% for middle and late layers. For *Language* steering, statistically significant gains were observed in 25%, 13.64%, and 17.05% of cases, respectively. Together with the zero-shot results, these findings indicate that early-layer steering consistently outperforms deeper interventions. Per-task results are provided in Appendix L.

Among models, Gemma-2-27B and Llama-3.1-70B were the most robust, showing predominantly positive performance shifts. In contrast, Gemma-2-9B behaved more inconsistently in the few-shot setting, exhibiting substantially more negative outcomes than in the zero-shot experiments.

Comparison of zero-shot vs. few-shot Comparing zero-shot and few-shot prompting reveals a clear saturation effect. In the zero-shot setting, steering often produced substantial gains, including double-digit F1 improvements, whereas the same configurations yielded smaller improvements under few-shot prompting. This suggests that few-shot demonstrations already shift the model toward the “human-authored” manifold, reducing the additional effect of steering. Nevertheless, *Quality* steering applied to early layers continued to provide consistent improvements.

The strong “early layers are best” trend observed in zero-shot generation also becomes less rigid in the few-shot setting. Larger models, in particular, showed greater robustness and occasionally benefited from steering at deeper layers.

Overall, our results indicate that activation steering is most impactful in zero-shot settings, where it can compensate for the absence of stylistic guid-

ance in the prompt. In few-shot scenarios, steering remains beneficial, though with smaller gains. This shows that, regardless of the original prompt used, steering is very often beneficial for downstream model performance when applied to the used prompting technique, with *Quality* steering offering higher and more consistent increases.

5 Evaluating Diversity of Generated Data

As shown in Section 4.2, both *Quality* and *Language* steering generally improve downstream model performance. To better understand these gains, we additionally analyze the diversity of the generated data, which has previously been linked to improved downstream robustness and performance (Cegin et al., 2026).

We evaluate four diversity metrics: (1) *Lexical diversity*, measuring the ratio of unique character 3-grams; (2) *Embedding diversity*, capturing the average pairwise cosine distance between text embeddings; (3) *Homogeneity*, measuring the structural consistency of the embedding distribution; and (4) *Isocontour radius*, estimating the overall spread of embedding representations. Details of the computation are provided in Appendix H.

We compute these metrics across all datasets. The aggregated results are in Table 3a for zero-shot and in Table 3b for few-shot prompting, with additional visualizations in Appendix N.

In the zero-shot setting, steering generally increases output diversity, particularly for smaller models and when applied to early layers. Most configurations show gains across all diversity metrics, with the notable exception of Gemma-2-27B, which also showed limited responsiveness in downstream evaluation. Increases in lexical and em-

bedding diversity indicate that steering reduces repetitive generations and broadens semantic coverage, while higher homogeneity suggests that this increased diversity remains structurally coherent rather than noisy. Early-layer steering appears especially effective because it influences representations before higher-level semantic processing has stabilized, allowing the model to propagate the intervention throughout the generation. This aligns with the stronger downstream performance improvements observed for earlier-layer steering.

In the few-shot setting, diversity gains become smaller and less consistent, suggesting that demonstrations already constrain the model’s latent space. In some cases, deeper-layer *Quality* steering even reduces diversity, particularly for Gemma-2-27B. Smaller models also tend to show decreases in homogeneity under few-shot steering, whereas larger models such as Llama-3.1-70B often maintain or improve structural consistency. Interestingly, *Language* steering frequently produces larger increases in embedding diversity than *Quality* steering in Llama models, possibly because language vectors are less aligned with the representations already induced by few-shot examples.

Although increased diversity does not always guarantee better downstream performance, our results show that steering vectors consistently encourage the generation of more varied data without explicit diversity optimization. This likely contributes to the improved robustness and downstream effectiveness that we observed.

6 Similarity of *Language* and *Quality* Vectors

We provide cosine similarity between *Language* and *Quality* vectors per LLM in Figure 3.

We see a distinct difference between LLMs in Gemma and Llama families. For Gemma 2 9b, most languages show low to moderate positive similarity, except Maltese. The polarity is much stronger in Gemma 2 27b. In contrast, both Llama models show a remarkably similar, highly polarized pattern for the same languages. Amharic, Hebrew, and Maltese are all nearly perfectly positively correlated, while other languages are nearly perfectly negatively correlated.

The languages appear to cluster into two distinct groups based on vector alignment. For the positive cluster, the *Quality* and *Language* are essentially in the same direction in the model’s latent space. This

explains why Amharic showed strong F1 gains under both steering types in most of our experiments. Notably, all of these languages are Semitic, which could indicate that for specific language groups the model sees *Quality* and *Language* as nearly identical. For the majority of languages, especially European and Southeast Asian ones, the *Quality* direction is often the mathematical opposite of *Language* direction in Llama models.

The polarities (−1.0 or 1.0) are most consistent in the early layers across all models. This reinforces our earlier conclusion that early-layer interventions are more effective because they target these highly defined, clear directions before the representations become more “mixed” or diffuse in later layers. The polarities indicate that the *Quality* direction is often the inverse of the *Language* direction. While both of these steering vectors can help in downstream model performance, this indicates that *Quality* steering vectors are (in general) not tied to a specific language and seem to be language-agnostic. This also explains the volatility of improvements from *Language* steering, as for most of the languages, it is steering the model in the exact opposite direction of the *Quality* steering.

7 Conclusion

In this work, we investigated activation steering for synthetic data generation in low-resource languages. Across four open-source LLMs, 11 languages, and two classification tasks, we showed that both *Language* and *Quality* steering can improve downstream performance while consistently increasing the diversity of generated data. Our results further demonstrate that steering is most effective when applied to earlier transformer layers, particularly in zero-shot settings. This finding is consistent with the cultural steering results of Ghussin et al. (2026b), who likewise observe that the effectiveness of vectors is strongly layer-dependent, with earlier-layer steering yielding the most reliable gains.

Overall, *Quality* steering proved to be the more reliable approach, suggesting that steering models toward a “human-authored” notion is especially beneficial for synthetic data generation. We additionally found that the relationship between *Language* and *Quality* vectors is highly model- and language-dependent, highlighting structured multi-lingual representations in LLM latent spaces.

These findings position activation steering as

an efficient alternative or complement to few-shot prompting for multilingual synthetic data generation, especially in low-resource settings.

Limitations

While our work demonstrates the efficacy of activation steering for low-resource synthetic data generation, several limitations remain to be addressed in future work.

Task and Generative Scope: Our evaluation is bounded by two text classification tasks: sentiment analysis and topic classification. While classification is a standard benchmark for evaluating synthetic data utility (Anikina et al., 2025; Cegin et al., 2026), it does not fully test the generative boundaries of activation steering.

Dependence on Machine Translation Quality: The extraction of our Quality steering vector relies on constructing contrastive text pairs via two rounds of backtranslation using the NLLB-200-600M distilled model. We did not investigate how many specific rounds or which translation models are the best.

Hyperparameter Sensitivity and Model Dependency: Our results highlight that the optimal steering intensity (α) and layer selection are heavily model- and language-dependent. For instance, the optimal scaling factors vary by orders of magnitude between the Gemma and Llama families (e.g., $\alpha = 100.0$ for Gemma 2 27B vs. $\alpha = 1.0$ for Llama 3.1 70B). At present, we lack a predictive, analytical framework to establish the perfect layer and multiplier combination a priori, limiting the immediate plug-and-play deployability of this method.

Acknowledgments

This work was partially funded by the European Union under the project lorAI - Low Resource Artificial Intelligence, GA No. 101136646, by NextGenerationEU through the Recovery and Resilience Plan for Slovakia under the project No. 09I01-03-V04-00006 (DistraceAI), and by the German Federal Ministry of Research, Technology and Space (BMFTR) as part of the project TRAILS (01IW24005).

This work was computationally supported by the Ministry of Education, Youth and Sports of the Czech Republic through the e-INFRA CZ (ID:90254) and by the use of the supercomputer PERUN, with the support of the European Union

from the funds of the Recovery and Resilience Plan of the Slovak Republic within the framework of project No. 17I03-04-P03-00001.

References

- David Ifeoluwa Adelani, Hannah Liu, Xiaoyu Shen, Nikita Vassilyev, Jesujoba O. Alabi, Yanke Mao, Haonan Gao, and En-Shiun Annie Lee. 2024. [SIB-200: A simple, inclusive, and big evaluation dataset for topic classification in 200+ languages and dialects](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 226–245, St. Julian’s, Malta. Association for Computational Linguistics.
- AI@Meta. 2024. [Llama 3 model card](#).
- Pierre Andrews, Mikel Artetxe, Mariano Coria Meglioli, Marta R. Costa-jussà, Joe Chuang, David Dale, Mark Duppenthaler, Nathaniel Paul Ekberg, Cynthia Gao, Daniel Edward Licht, Jean Maillard, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Eduardo Sánchez, Ioannis Tsiamas, Arina Turkatenco, Albert Ventayol-Boada, and Shireen Yates. 2025. [BOUQuET : dataset, benchmark and open initiative for universal quality evaluation in translation](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 27515–27535, Suzhou, China. Association for Computational Linguistics.
- Tatiana Anikina, Jan Cegin, Jakub Simko, and Simon Ostermann. 2025. [A rigorous evaluation of LLM data generation strategies for low-resource languages](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 8282–8303, Suzhou, China. Association for Computational Linguistics.
- Andy Arditi, Oscar Obeso, Aaquib Syed, Daniel Paleka, Nina Panickssery, Wes Gurnee, and Neel Nanda. 2024. [Refusal in language models is mediated by a single direction](#). *Advances in Neural Information Processing Systems*, 37:136037–136083.
- Lukasz Bartoszcze, Sarthak Munshi, Bryan Sukidi, Jennifer Yen, Zejia Yang, David Williams-King, Linh Le, Kosi Asuzu, and Carsten Maple. 2025. [Representation engineering for large-language models: Survey and research challenges](#). *Preprint*, arXiv:2502.17601.
- Tadesse Destaw Belay, Shahriar Kabir Nahin, Israel Abebe Azime, Ocean Monjur, Shamsuddeen Hassan Muhammad, Seid Muhie Yimam, and Anshuman Chhabra. 2026. [Afrilangtutor: Advancing language tutoring and culture education in low-resource languages with large language models](#). *Preprint*, arXiv:2604.20996.
- Nora Belrose, Igor Ostrovsky, Lev McKinney, Zach Furman, Logan Smith, Danny Halawi, Stella Biderman,

- and Jacob Steinhardt. 2025. [Eliciting latent predictions from transformers with the tuned lens](#). *Preprint*, arXiv:2303.08112.
- Tomáš Brychcín and Ivan Habernal. 2013. [Unsupervised improving of sentiment analysis using global target context](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*, pages 122–128, Hissar, Bulgaria. INCOMA Ltd. Shoumen, BULGARIA.
- Jan Cegin, Branislav Pecher, Jakub Simko, Ivan Srba, Maria Bielikova, and Peter Brusilovsky. 2024a. [Effects of diversity incentives on sample diversity and downstream model performance in LLM-based text augmentation](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 13148–13171, Bangkok, Thailand. Association for Computational Linguistics.
- Jan Cegin, Branislav Pecher, Jakub Simko, Ivan Srba, Maria Bielikova, and Peter Brusilovsky. 2024b. [Use random selection for now: Investigation of few-shot selection strategies in llm-based text augmentation for classification](#). *Preprint*, arXiv:2410.10756.
- Jan Cegin, Branislav Pecher, Ivan Srba, and Jakub Simko. 2026. [RoSE: Round-robin synthetic data evaluation for selecting LLM generators without human test sets](#). In *Proceedings of the 19th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5530–5545, Rabat, Morocco. Association for Computational Linguistics.
- Jan Cegin, Jakub Simko, and Peter Brusilovsky. 2023. [ChatGPT to replace crowdsourcing of paraphrases for intent classification: Higher diversity and comparable model robustness](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 1889–1905, Singapore. Association for Computational Linguistics.
- Cheng-Ting Chou, George Liu, Jessica Sun, Cole Blondin, Kevin Zhu, Vasu Sharma, and Sean O’Brien. 2025. [Causal language control in multilingual transformers via sparse feature steering](#). *Preprint*, arXiv:2507.13410.
- Yi-Ling Chung, Aurora Cobo, and Pablo Serna. 2025. [Beyond translation: Llm-based data generation for multilingual fact-checking](#). *arXiv preprint arXiv:2502.15419*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Unsupervised cross-lingual representation learning at scale](#). *CoRR*, abs/1911.02116.
- Hoagy Cunningham, Aidan Ewart, Logan Riggs Smith, Robert Huben, and Lee Sharkey. 2024. [Sparse autoencoders find highly interpretable features in language models](#). In *The Twelfth International Conference on Learning Representations*.
- Haixing Dai, Zhengliang Liu, Wenxiong Liao, Xiaoke Huang, Yihan Cao, Zihao Wu, Lin Zhao, Shaochen Xu, Wei Liu, Ninghao Liu, Sheng Li, Dajiang Zhu, Hongmin Cai, Lichao Sun, Quanzheng Li, Dinggang Shen, Tianming Liu, and Xiang Li. 2023. [Aug-gpt: Leveraging chatgpt for text data augmentation](#). *Preprint*, arXiv:2302.13007.
- Boyi Deng, Yu Wan, Baosong Yang, Yidan Zhang, and Fuli Feng. 2025. [Unveiling language-specific features in large language models via sparse autoencoders](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4563–4608, Vienna, Austria. Association for Computational Linguistics.
- Luyang Fang, Gyeong-Geon Lee, and Xiaoming Zhai. 2023. [Using gpt-4 to augment unbalanced data for automatic scoring](#). *Preprint*, arXiv:2310.18365.
- Steven Y. Feng, Varun Gangal, Jason Wei, Sarath Chandar, Soroush Vosoughi, Teruko Mitamura, and Eduard Hovy. 2021. [A survey of data augmentation approaches for NLP](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 968–988, Online. Association for Computational Linguistics.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. 2024. [Scaling and evaluating sparse autoencoders](#). ArXiv:2406.04093.
- Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. 2021. [Transformer feed-forward layers are key-value memories](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 5484–5495, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Sreyan Ghosh, Chandra Kiran Evuru, Sonal Kumar, S Rameswaran, S Sakshi, Utkarsh Tyagi, and Dinesh Manocha. 2023. [Dale: Generative data augmentation for low-resource legal nlp](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, Sentosa, Singapore.
- Yusser Al Ghussin, Daniil Gurgurov, Tanja Baeumel, Josef van Genabith, Patrick Schramowski, and Simon Ostermann. 2026a. [Multilingual steering by design: Multilingual sparse autoencoders and principled layer selection](#). *Preprint*, arXiv:2605.23036.
- Yusser Al Ghussin, Daniil Gurgurov, Yasser Hamidullah, Josef van Genabith, Cristina España-Bonet, and Simon Ostermann. 2026b. [Dfki-mlt at semeval-2026 task 7: Steering multilingual models towards cultural knowledge](#). *Preprint*, arXiv:2605.23069.
- Parker Glenn, Alolika Gon, Nikhil Kohli, Sihan Zha, Parag Pravin Dakle, and Preethi Raghavan. 2023. [Jetsons at the FinNLP-2023: Using synthetic data and transfer learning for multilingual ESG issue classification](#). In *Proceedings of the Fifth Workshop on*

- Financial Technology and Natural Language Processing and the Second Multimodal AI For Financial Forecasting*, pages 133–139, Macao. -.
- Daniil Gurgurov, Yusser Al Ghussin, Tanja Baeumel, Cheng-Ting Chou, Patrick Schramowski, Marius Mosbach, Josef van Genabith, and Simon Ostermann. 2026. Clas-bench: A cross-lingual alignment and steering benchmark. *arXiv preprint arXiv:2601.08331*.
- Daniil Gurgurov, Mareike Hartmann, and Simon Ostermann. 2024. [Adapting multilingual LLMs to low-resource languages with knowledge graphs via adapters](#). In *Proceedings of the 1st Workshop on Knowledge Graphs and Large Language Models (KaLLM 2024)*, pages 63–74, Bangkok, Thailand. Association for Computational Linguistics.
- Daniil Gurgurov, Rishu Kumar, and Simon Ostermann. 2025a. [GrEmLin: A repository of green baseline embeddings for 87 low-resource languages injected with multilingual graph knowledge](#). In *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 1204–1221, Albuquerque, New Mexico. Association for Computational Linguistics.
- Daniil Gurgurov, Katharina Trinley, Yusser Al Ghussin, Tanja Baeumel, Josef Van Genabith, and Simon Ostermann. 2025b. [Language arithmetics: Towards systematic language neuron identification and manipulation](#). In *Proceedings of the 14th International Joint Conference on Natural Language Processing and the 4th Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics*, pages 2911–2937, Mumbai, India. The Asian Federation of Natural Language Processing and The Association for Computational Linguistics.
- Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The state and fate of linguistic diversity and inclusion in the nlp world. In *Proceedings of the 58th annual meeting of the association for computational linguistics*, pages 6282–6293.
- Yeskendir Koishekenov, Alexandre Berard, and Vasilina Nikoulina. 2023. [Memory-efficient NLLB-200: Language-specific expert pruning of a massively multilingual machine translation model](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3567–3585, Toronto, Canada. Association for Computational Linguistics.
- Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer Schütt, Oliver Bensch, Roxanne El Baff, Dominik Opitz, and Tobias Hecking. 2024. [Style vectors for steering generative large language model](#). *Preprint*, arXiv:2402.01618.
- Yi-An Lai, Xuan Zhu, Yi Zhang, and Mona Diab. 2020. [Diversity, density, and homogeneity: Quantitative characteristic metrics for text collections](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1739–1746, Marseille, France. European Language Resources Association.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2023. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36:41451–41530.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. 2024. Inference-time intervention: Eliciting truthful answers from a language model. *Advances in Neural Information Processing Systems*, 36.
- Linlin Liu, Bosheng Ding, Lidong Bing, Shafiq Joty, Luo Si, and Chunyan Miao. 2021. [MulDA: A multilingual data augmentation framework for low-resource cross-lingual NER](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5834–5846, Online. Association for Computational Linguistics.
- Qijiong Liu, Nuo Chen, Tetsuya Sakai, and Xiao-Ming Wu. 2024. [Once: Boosting content-based recommendation with both open- and closed-source large language models](#). In *Proceedings of the 17th ACM International Conference on Web Search and Data Mining, WSDM '24*, page 452–461, New York, NY, USA. Association for Computing Machinery.
- Henry B Mann and Donald R Whitney. 1947. On a test of whether one of two random variables is stochastically larger than the other. *The annals of mathematical statistics*, pages 50–60.
- Samuel Marks and Max Tegmark. 2023. [The geometry of truth: Emergent linear structure in large language model representations of true/false datasets](#). *Preprint*, arXiv:2310.06824.
- Amani Namboori, Shivam Sadashiv Mangale, Andy Rosenbaum, and Saleh Soltan. 2023. [Gemquad: Generating multilingual question answering datasets from large language models using few shot learning](#). In *NeurIPS 2023 Workshop on Synthetic Data Generation with Generative AI*.
- Neel Nanda and Joseph Bloom. 2022. [Transformerlens](https://github.com/TransformerLensOrg/TransformerLens). <https://github.com/TransformerLensOrg/TransformerLens>.
- NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Hefernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, and 20 others. 2024. [Scaling neural machine translation to 200 languages](#). *Nature*, 630(8018):841–846.
- Sebastian Nordhoff and Harald Hammarström. 2011. [Glottolog/langdoc: Defining dialects, languages, and language families as collections of resources](#). In

- First International Workshop on Linked Science 2011- In conjunction with the International Semantic Web Conference (ISWC 2011).*
- Aytuğ Onan. 2023. [Srl-aco: A text augmentation framework based on semantic role labeling and ant colony optimization](#). *Journal of King Saud University - Computer and Information Sciences*, 35(7):101611.
- Simon Ostermann, Daniil Gurgurov, Tanja Baeumel, Michael A Hedderich, Sebastian Lapuschkin, Wojciech Samek, and Vera Schmitt. 2026. From weights to activations: Is steering the next frontier of adaptation? *arXiv preprint arXiv:2604.14090*.
- Nina Panickssery, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Matt Turner. 2023. Steering llama 2 via contrastive activation addition. *arXiv preprint arXiv:2312.06681*.
- Kiho Park, Yo Joong Choe, and Victor Veitch. 2024. [The linear representation hypothesis and the geometry of large language models](#). *Preprint*, arXiv:2311.03658.
- Branislav Pecher, Jan Cegin, Robert Belanec, Ivan Srba, Jakub Simko, and Maria Bielikova. 2026. [Better as generators than classifiers: Leveraging LLMs and synthetic data for low-resource multilingual classification](#). In *Findings of the Association for Computational Linguistics: EACL 2026*, pages 2840–2857, Rabat, Morocco. Association for Computational Linguistics.
- Frédéric Piedboeuf and Philippe Langlais. 2023. [Is ChatGPT the ultimate data augmentation algorithm?](#) In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 15606–15615, Singapore. Association for Computational Linguistics.
- Salsabila Zahirah Pranida, Rifo Ahmad Genadi, and Fajri Koto. 2025. [Synthetic data generation for culturally nuanced commonsense reasoning in low-resource languages](#). *Preprint*, arXiv:2502.12932.
- Gorjan Radevski, Kiril Gashteovski, Giwon Hong, Carolin Lawrence, and Goran Glavaš. 2026. Compositional steering of large language models with steering tokens. *arXiv preprint arXiv:2601.05062*.
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. 2024. [Steering llama 2 via contrastive activation addition](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522, Bangkok, Thailand. Association for Computational Linguistics.
- Gaurav Sahu, Pau Rodriguez, Issam Laradji, Parmida Atighehchian, David Vazquez, and Dzmitry Bahdanau. 2022. [Data augmentation for intent classification with off-the-shelf large language models](#). In *Proceedings of the 4th Workshop on NLP for Conversational AI*, pages 47–57, Dublin, Ireland. Association for Computational Linguistics.
- Max Schaffelder and Albert Gatt. 2026. [Synthetic eggs in many baskets: The impact of synthetic data diversity on llm fine-tuning](#). *Preprint*, arXiv:2511.01490.
- Tianyi Tang, Wenyang Luo, Haoyang Huang, Dongdong Zhang, Xiaolei Wang, Wayne Xin Zhao, Furu Wei, and Ji-Rong Wen. 2024. Language-specific neurons: The key to multilingual capabilities in large language models. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5701–5715.
- Gemma Team, Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, Léonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ramé, Johan Ferret, Peter Liu, Pouya Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, and 179 others. 2024. [Gemma 2: Improving open language models at a practical size](#). *Preprint*, arXiv:2408.00118.
- Ian Tenney, Dipanjan Das, and Ellie Pavlick. 2019. [BERT rediscovers the classical NLP pipeline](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4593–4601, Florence, Italy. Association for Computational Linguistics.
- Van-Hien Tran, Huy Hien Vu, Hideki Tanaka, and Masao Utiyama. 2026. [Representation-aware prompting for zero-shot Marathi text classification: IPA, Romanization, repetition](#). In *Proceedings of the Second Workshop on Language Models for Low-Resource Languages (LoResLM 2026)*, pages 436–443, Rabat, Morocco. Association for Computational Linguistics.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2023. [Steering language models with activation engineering](#). *Preprint*, arXiv:2308.10248.
- Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2024. [Steering language models with activation engineering](#). *Preprint*, arXiv:2308.10248.
- Solomon Ubani, Suleyman Olcay Polat, and Rodney Nielsen. 2023. [Zeroshotdataaug: Generating and augmenting training data with chatgpt](#). *Preprint*, arXiv:2304.14334.
- Sing Hieng Wong, Hassan Sajjad, and AB Siddique. 2026. [Langfir: Discovering sparse language-specific features from monolingual data for language steering](#). *arXiv preprint arXiv:2604.03532*.
- Zhengxuan Wu, Aryaman Arora, Atticus Geiger, Zheng Wang, Jing Huang, Dan Jurafsky, Christopher D Manning, and Christopher Potts. 2025. [Axbench: Steering llms? even simple baselines outperform sparse autoencoders](#). In *Forty-second International Conference on Machine Learning*.

Zhengxuan Wu, Aryaman Arora, Zheng Wang, Atticus Geiger, Dan Jurafsky, Christopher D. Manning, and Christopher Potts. 2024. [ReFT: Representation finetuning for language models](#). *arXiv:2404.03592*.

Yiran Zhao, Wenxuan Zhang, Guizhen Chen, Kenji Kawaguchi, and Lidong Bing. 2024. How do large language models handle multilingualism? *Advances in Neural Information Processing Systems*, 37:15296–15319.

Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, and 1 others. 2023. Representation engineering: A top-down approach to ai transparency. *arXiv preprint arXiv:2310.01405*.

A Ethical Considerations

Based on a thorough ethical assessment performed on the basis of intra-institutional ethical guidelines and checklists tailored to the use of data and algorithms, we see no ethical concerns pertaining directly to the conduct of this research. Although the production of new data through LLMs bears several risks, such as the introduction of biases, the small size of the produced dataset, sufficient for experimentation, is, at the same time, insufficient for any major machine learning endeavours where such biases could be transferred.

We follow the license terms for all the models and datasets we used (such as the one required for the use of the Llama-3.1 model) – all models and datasets allow their use as part of the research.

We used generative AI for grammatical checks and fixes of the paper, and have verified all the texts changed by the generative AI.

B Computational Resources

Our experiments were run on a computational cluster with H200 GPUs, 2x AMD EPYC 9745 CPUs, and 128 GB of RAM. Our generation experiments took approximately 6-8 minutes per one combination of LLM, steering vector, layer, task, language, and alpha value. Given that we had a total of 2112 combinations, our generation experiments took approximately 260 GPU hours. For the finetuning experiments, we used approximately 240 GPU hours. In total, for our experiments, we used approximately 500 GPU hours.

C Linear Probe Details

To quantify the linear separability of the quality signal within the model’s latent representations, we

implemented a layer-wise probing diagnostic using logistic regression. For each input sequence in our human-authored and backtranslated distributions, we extracted the hidden states from the model’s residual stream. We then applied mean pooling across non-padding tokens to generate representative sequence-level embeddings for every layer. These activations were sanitized by clipping outliers and then L2-normalized to facilitate stable classifier training. Using a stratified 80/20 train-test split, we trained a logistic regression classifier on the embeddings of each layer independently and measured the resulting classification accuracy. As illustrated in Figure 2, the layers exhibiting peak accuracy represent the regions where the quality concept is most robustly and linearly encoded, providing the optimal candidates for activation steering interventions.

D Steering Vectors Computation Details

D.1 Language Steering Vector

To construct the language steering vectors, we employ a one-vs-rest contrastive methodology across the residual stream of all transformer blocks. This approach ensures that the resulting vector captures features specific to the target language rather than shared cross-lingual properties or general model priors.

For each target language $L \in \mathcal{L}$, where \mathcal{L} is our set of 11 typologically diverse languages, we pass the dataset through the model. Using the TransformerLens framework, we extract the residual stream activations for every layer l . We compute the mean activation vector $\mu_{l,L}$ by averaging across all non-padding tokens N_L in the language-specific dataset:

$$\mu_{l,L} = \frac{1}{N_L} \sum_{i,j} \mathbf{a}_{l,i,j} \cdot m_{i,j}$$

. This process results in a language-specific centroid in the latent space for every layer of the model.

To isolate the linguistic identity of language L , we define the steering vector $\mathbf{v}_{l,L}$ as the difference between the language’s mean activation and the mean activation of all other languages in our study:

$$\mathbf{v}_{l,L} = \mu_{l,L} - \frac{1}{|\mathcal{L}| - 1} \sum_{L' \in \mathcal{L}, L' \neq L} \mu_{l,L'}$$

By contrasting the target language against a global multilingual mean, we filter out common semantic

and structural activations, focusing the vector on the unique "fingerprint" of the target language.

Finally, we apply layer-wise L2 normalization to the difference vectors:

$$\hat{\mathbf{v}}_{l,L} = \frac{\mathbf{v}_{l,L}}{\|\mathbf{v}_{l,L}\|_2}$$

This normalization step is critical for maintaining stability during the generation phase, as it ensures that the steering magnitude α remains consistent across different layers and languages regardless of the original activation scales.

D.2 Quality Steering Vector

For each target language $L \in \mathcal{L}$, we utilize two matching corpora representing a binary quality dimension $\mathcal{D} = \{d_{\text{good}}, d_{\text{bad}}\}$, where d_{good} consists of authentic, human-authored text and d_{bad} comprises corresponding backtranslations. We process each corpus through the model independently to capture the token-level post-residual block hidden states $\mathbf{a}_{l,i,j}$ across all layers l . The direct mean activation vector $\mu_{l,L,d}$ for an attribute dimension $d \in \mathcal{D}$ within language L is accumulated over its respective unpadded token count $N_{L,d}$:

$$\mu_{l,L,d} = \frac{1}{N_{L,d}} \sum_{i,j} \mathbf{a}_{l,i,j} \cdot m_{i,j}$$

To extract a clean signal, the raw quality direction $\mathbf{v}_{l,L,\text{qual}}$ is isolated by subtracting the centroid of the corrupted dataset from the centroid of the human-authored dataset:

$$\mathbf{v}_{l,L,\text{qual}} = \mu_{l,L,d_{\text{good}}} - \mu_{l,L,d_{\text{bad}}}$$

By establishing this explicit pairwise vector difference within the same language, task-neutral semantic properties cancel out, leaving a vector that should isolate the geometric directional shift toward more human-like data.

Following the same stability protocol as our language vectors, we enforce a layer-wise L2 normalization step to yield the final quality steering direction:

$$\hat{\mathbf{v}}_{l,L,\text{qual}} = \frac{\mathbf{v}_{l,L,\text{qual}}}{\|\mathbf{v}_{l,L,\text{qual}}\|_2}$$

E Language Abbreviations

Language abbreviations for each language used in the study can be found in Table 4.

Code	Language
am	Amharic
cs	Czech
de	German
da	Danish
he	Hebrew
id	Indonesian
ja	Japanese
mt	Maltese
sk	Slovak
sl	Slovenian
su	Sundanese

Table 4: Language abbreviations.

F Backtranslation Details

For the backtranslation and creation of parallel synthetic data to human-authored data, we employ the NLLB 600M distilled model¹. To ensure enough distinction between the human-authored and synthetic data, we use two rounds of backtranslation from the source language to target language 1, then to target language 2, back to target language 1, and then back to the source language. We use English as target language 1 and Chinese as target language 2 as the two target languages, due to their resourcefulness. We also investigated additional corruptions of the synthetic data (e.g., random token swaps), but decided against it given the sufficient linear probe results from Section 4.1.

G Downstream Fine-tuning of XLM-R

For the downstream evaluation, we fine-tune the XLM-R (Conneau et al., 2019) *FacebookAI/xlm-roberta-base* model with a batch size of 16 and employ early stopping with a patience of 10 epochs to prevent overfitting. We perform hyperparameter optimisation to determine the optimal learning rate and set it to 1e-5. *AdamW* is used as an optimiser. We balance the generated datasets to have the same number of samples per label. We perform finetuning 20 times with different seeds for each generated data. We normalise all inputs by converting them to lowercase and removing punctuation.

H Diversity Metrics

To quantitatively evaluate the geometric and linguistic characteristics of the generated text, we employ four distinct diversity metrics. Before metric computation, all generated texts undergo a standardized normalization pipeline, which normalizes

¹<https://huggingface.co/facebook/nllb-200-distilled-600M>

text into the NFKD variant, strips diacritics and non-textual symbols (such as emojis), and collapses consecutive whitespaces.

For dense vector calculations, text sequences are mapped into a latent embedding space using the Qwen3-Embedding-4B² encoder model.

Lexical diversity: Calculated as the ratio of unique character-level trigrams to total trigrams across the corpus. It captures surface-level variety; higher scores indicate diverse, natural phrasing, while lower scores signify repetitive or formulaic text.

Embedding diversity: Computed by grouping the dataset by task labels, calculating the mean pairwise cosine distance within each group g , and taking the macro-average across all unique labels G . It measures categorical semantic dispersion, showing whether a model generates a broad range of distinct concepts or collapses into a narrow thematic subspace.

Isocontour radius (Lai et al., 2020): Determines the geometric boundary of the generated dataset by measuring the Euclidean distance of all individual embeddings to the global dataset centroid μ . The metric outputs the 90th percentile of this distance distribution. It represents the total volume of the latent operational envelope; a larger radius indicates that the steering intervention has pushed the generations into wider, peripheral regions of the latent space.

Homogeneity (Lai et al., 2020): Evaluates structural density uniformity using a similarity-based Kernel Density Estimation (KDE). It applies an RBF kernel ($\sigma = 0.5$) to the pairwise cosine similarities to compute a localized density score for each point. It tracks structural consistency across the manifold; a lower score indicates a highly uniform, evenly distributed layout, whereas a higher score signals structural imbalance and tight data clustering.

I Prompt Templates and Generation Details

For generating synthetic data, we used this general prompt template for the **zero-shot** setting:

Topic: *"Generate a long wiki-style sentence on the topic of label in language language. Output*

²<https://huggingface.co/Qwen/Qwen3-Embedding-4B>

only the text in language."

Sentiment: *"Generate a review with label sentiment in language language. Output only the text and nothing else."*

For the **few-shot** setting, we used 5 shots in the target language and label, with these prompts:

Topic: *"Generate a long wiki style sentence on the topic of label in language language based on examples. Output only the text in language. Examples: examples"*

Sentiment: *"Generate a review with label sentiment in language language based on examples. Output only the text and nothing else. Examples: examples"*

Sampling parameters used for generation were: *temperature=0.8, top_p=0.9, max_tokens=512, freq_penalty=0.0*. We excluded duplicates to ensure unique samples were collected.

Specific versions of LLMs used for generations were: Llama-3.1-70b-instruct³, Llama-3.1-8b-instruct⁴, gemma-2-27b-it⁵, gemma-2-9b-it⁶. The LLMs were used with *bfloat16* precision.

J Baseline Results for Zero- and Few-shot Setting

We provide the baseline result for topic and sentiment detection for both zero-shot in Table 5 and for the few-shot setting in Table 6.

K Alpha Values Analysis

We provide an analysis of alpha values in all settings in Table 11 for zero-shot and in Table 12 for the few-shot setting.

For generic language vectors, increasing alpha often behaves like a poison pill. As alpha increases, variance widens significantly, and medians frequently trend downward. *Quality* vectors handle high alpha values much better. Instead of collapsing, higher alphas either monotonically improve performance or reach a safe plateau with tight variances. We see the best performance for alpha values at 50 or 75 for Gemma 2 models and 2 or 3 for Llama3.1 models, indicating that the best downstream model performance comes at moderate-to-strong steering.

³<https://huggingface.co/meta-llama/Llama-3.1-70B-Instruct>

⁴<https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

⁵<https://huggingface.co/google/gemma-2-27b-it>

⁶<https://huggingface.co/google/gemma-2-9b-it>

Larger models are more brittle to bad steering, as for both Llama 3.1 70b and Gemma 2 27b, choosing the wrong vector (*Language*) or the wrong layer/alpha combination leads to dramatic, sweeping performance drops, most notably for high alpha values in the Llama3.1 70b model. They require precise steering. On the other hand, Llama 3.1 8b and Gemma 2 9b show remarkably clean, uniform positive shifts across almost all alpha values in their early and middle layers, indicating that smaller models might have more centralized "quality pathways" that are easier to uniformly amplify.

L Downstream Model Performance Per Task

We provide a per-task breakdown for the zero-shot setting for the topic detection task in Table 7 and for the sentiment detection task in Table 8. For the few-shot setting, the topic detection task can be found in Table 9 and the sentiment detection task in Table 10. All of these tables also contain visualizations showing which specific cases had statistically significant increases in downstream model performance.

In terms of the zero-shot setting, we see similar results to those observed in Section 4.2, as the early layers *Quality* steering vector outperforms other approaches, except for the Llama 3.1 70b model, where *Language* vectors applied to earlier layers slightly outperform it. The best performing method is steering at earlier layers, as it offers an increase in downstream model performance in 79.55% of cases for both topic and sentiment detection. Gemma models see a larger increase in sentiment detection task, while Llama models see a larger increase in downstream model performance in the topic detection task.

For the few-shot setting, again, the *Quality* steering vectors outperform the *Language* steering vectors, in particular at the earlier layers. For topic detection, steering via early layer *Quality* vectors results in increased downstream performance in 70.45% of cases for topic detection and 84.09% of cases for sentiment detection. The noisier sentiment detection seems to benefit more from the "human-likeness" introduced by the *Quality* steering vectors.

M Detailed Downstream Model Visualizations

We provide detailed visualizations for different languages, LLMs, steering vectors, and layers in Table 13 for zero-shot and in Table 14 for few-shot.

N Detailed Diversity Visualizations

We provide a detailed diversity visualization for each of our 4 metrics and approaches for both the zero-shot baseline and the few-shot baseline in Tables 15 and 16, respectively.

Table 5: Mean F1 scores across languages and tasks for baseline **zero-shot** setting with no steering.

LLM	Language Task	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Gemma 2 27b	Sentiment	74.53	83.55	95.27	64.10	75.69	84.07	45.47	54.82	69.53	73.23	67.35	71.60
	Topic	48.28	59.26	68.93	64.82	62.99	65.99	60.43	41.33	64.46	63.84	61.01	60.12
Gemma 2 9b	Sentiment	79.52	80.78	94.97	71.50	71.98	87.07	62.53	53.12	66.06	68.64	73.58	73.62
	Topic	49.05	61.48	64.54	62.06	60.26	66.35	60.19	38.14	60.98	63.56	62.40	59.00
Llama3.1 70b	Sentiment	77.21	82.94	91.72	68.91	74.46	86.87	61.79	57.17	71.04	79.69	74.29	75.10
	Topic	50.21	75.82	76.57	73.54	67.18	70.16	49.64	39.01	71.93	76.71	61.85	64.78
Llama3.1 8b	Sentiment	53.66	85.10	94.05	71.18	76.77	87.93	62.10	51.78	82.53	81.89	78.88	75.08
	Topic	33.82	71.06	71.55	75.84	71.24	68.04	62.72	46.41	71.65	69.79	65.18	64.30

Table 6: Mean F1 scores across languages and tasks for baseline **few-shot** setting with no steering.

LLM	Language Task	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Gemma 2 27b	Sentiment	74.93	68.14	91.64	61.12	58.20	80.64	60.01	59.25	75.24	59.76	80.36	69.94
	Topic	56.70	76.64	75.06	79.39	69.56	74.96	67.42	48.30	76.30	73.02	66.14	69.41
Gemma 2 9b	Sentiment	81.69	75.47	87.32	59.17	55.99	87.20	54.91	57.81	72.60	82.43	78.87	72.13
	Topic	56.14	76.33	74.82	76.07	73.53	74.62	62.76	51.02	77.48	74.51	68.82	69.65
Llama3.1 70b	Sentiment	82.72	83.30	94.68	56.36	70.67	77.89	56.34	57.13	65.10	73.90	74.65	72.07
	Topic	62.65	77.92	75.50	80.55	74.85	81.03	69.51	48.74	78.90	76.78	68.74	72.29
Llama3.1 8b	Sentiment	76.53	81.48	85.74	69.01	78.52	87.20	56.67	55.94	80.13	72.16	80.24	74.87
	Topic	52.06	79.78	76.74	76.04	74.93	76.01	69.14	48.15	77.94	76.30	68.11	70.47

Table 7: Mean F1 difference across languages, models, steering methods, and layers for **topic** detection, with **positive** and **negative** differences compared to **zero-shot** prompt with no steering. Cells with * denote statistically significant increases over baseline.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Language	Gemma 2 27b	L10 ($\alpha = 100.0$)	0.84	0.02	-3.76	1.31	-0.36	3.69*	-0.83	0.42	0.67	-0.79	-0.59	0.06
		L22 ($\alpha = 75.0$)	2.60*	0.08	0.21	0.68	0.47	2.14*	-1.31	2.65*	0.20	-1.15	-0.10	0.59
		L34 ($\alpha = 75.0$)	6.49*	0.84	-0.08	-0.54	-0.98	3.36*	-3.19	-0.64	-0.53	-0.41	1.69	0.55
	Gemma 2 9b	L9 ($\alpha = 50.0$)	-4.49	0.72	2.27*	1.93*	1.44*	-0.32	-4.16	2.06*	3.65*	3.30*	-0.74	0.51
		L20 ($\alpha = 50.0$)	-4.87	-3.41	5.19*	-0.35	-0.59	-1.18	-0.75	2.82*	2.44*	-1.22	0.62	-0.12
		L31 ($\alpha = 25.0$)	4.84*	0.02	3.51*	0.90	1.50*	0.35	-4.71	4.41*	-0.59	0.90	0.76	1.08
	Llama3.1 70b	L17 ($\alpha = 2.0$)	1.80	-3.94	-4.28	3.54*	-0.32	1.35	4.75*	0.62	7.77*	-1.37	-3.39	0.59
		L37 ($\alpha = 3.0$)	2.36	-4.22	1.23	4.38*	0.06	1.41	8.30*	0.78	3.42*	0.80	-3.98	1.32
		L57 ($\alpha = 2.0$)	3.09*	-4.31	1.04	1.23	1.02	2.17*	1.93*	0.51	0.20	1.76	-4.40	0.39
	Llama3.1 8b	L7 ($\alpha = 2.0$)	10.57*	-1.55	1.39	0.73	-1.18	2.31	-0.44	1.06	4.79*	3.45*	0.47	1.96
		L15 ($\alpha = 1.0$)	13.84*	3.03*	3.52*	1.19	0.08	-0.22	1.16	-1.72	2.61*	4.58*	-4.44	2.15
		L23 ($\alpha = 1.0$)	7.99*	2.79*	3.48*	0.34	-5.44	3.32*	0.43	0.38	6.00*	1.07	-1.93	1.68
Quality	Gemma 2 27b	L10 ($\alpha = 100.0$)	4.66*	0.18	0.14	0.98	-0.71	0.16	0.27	0.94	-0.01	0.78	-0.80	0.60
		L22 ($\alpha = 100.0$)	3.18*	0.01	-2.85	-1.11	1.23	0.90	-0.91	3.90*	0.58	-2.16	-1.20	0.14
		L34 ($\alpha = 50.0$)	0.84	-0.40	-1.36	1.68	0.42	1.45	-0.24	-0.49	0.35	0.87	-0.50	0.24
	Gemma 2 9b	L9 ($\alpha = 75.0$)	0.44	6.01*	3.07*	2.44*	4.19*	1.09*	0.80	2.18*	6.95*	6.30*	4.15*	3.42
		L20 ($\alpha = 50.0$)	2.68*	2.30*	6.40*	5.08*	-0.05	3.25*	-0.11	3.91*	2.21*	6.20*	1.02	2.99
		L31 ($\alpha = 75.0$)	2.27*	3.89*	2.61*	4.01*	0.22	2.15*	-0.49	4.23*	12.23*	2.88*	0.44	3.13
	Llama3.1 70b	L17 ($\alpha = 3.0$)	2.69*	-0.83	-0.91	0.15	0.10	0.14	-1.07	1.01*	3.57*	-1.23	-3.53	0.01
		L37 ($\alpha = 1.0$)	3.23*	-4.43	-0.45	-3.04	0.93	1.75	-4.45	-1.60	3.01*	0.77	-5.05	-0.85
		L57 ($\alpha = 2.0$)	6.17*	-2.15	1.57	1.91*	0.03	1.12	0.32	1.74	-0.63	0.80	-3.75	0.65
	Llama3.1 8b	L7 ($\alpha = 2.0$)	5.55*	5.18*	4.28*	0.74	0.64	3.59*	3.47*	2.45*	5.00*	5.45*	-1.11	3.20
		L15 ($\alpha = 3.0$)	3.36*	1.14	2.03	-2.44	-1.28	-0.79	1.68	0.81	3.46*	1.17	2.19	1.03
		L23 ($\alpha = 2.0$)	9.53*	2.94*	1.87	0.07	-6.05	2.52*	-0.80	1.46	4.36*	4.44*	1.74	2.01

Table 8: Mean F1 difference across languages, models, steering methods, and layers for **sentiment** detection, with **positive** and **negative** differences compared to a **zero-shot** prompt with no steering. Cells with * denote statistically significant increases over baseline.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Language	Gemma 2 27b	L10 ($\alpha = 100.0$)	-0.50	0.53	-0.44	3.16*	-1.97	1.30	4.51*	-1.54	2.37*	6.12*	5.24*	1.71
		L22 ($\alpha = 75.0$)	2.12*	-0.53	-0.46	0.79	-1.51	-0.59	-4.03	-1.90	3.46*	1.29	9.94*	0.78
		L34 ($\alpha = 75.0$)	2.20*	-1.85	-0.77	-2.20	0.26	0.98	7.11*	-2.78	-2.91	0.26	10.71*	1.00
	Gemma 9b	L9 ($\alpha = 50.0$)	-2.01	2.23	-0.09	-1.08	3.13*	0.66	-3.10	-1.15	4.52*	10.47*	1.07	1.33
		L20 ($\alpha = 50.0$)	1.43	2.90*	-1.43	-1.73	-6.48	0.45	-0.96	2.69*	10.20*	8.74*	5.95*	1.98
		L31 ($\alpha = 25.0$)	-0.24	-1.43	-0.14	0.10	1.96	-1.27	0.27	1.51	-0.58	10.34*	4.27*	1.34
	Llama3.1 70b	L17 ($\alpha = 2.0$)	4.71*	-2.11	2.46*	1.71	2.29	1.02	0.09	0.95	11.57*	3.91*	2.26	2.62
		L37 ($\alpha = 3.0$)	4.32*	2.72*	1.55	2.47	1.25	0.52	-0.11	0.28	6.22*	3.65*	5.68*	2.60
		L57 ($\alpha = 2.0$)	2.85*	2.43	2.35	-0.36	2.62*	1.62	-3.86	2.47	3.53*	4.94*	7.53*	2.37
	Llama3.1 8b	L7 ($\alpha = 2.0$)	6.21*	-0.85	0.97	0.80	-0.37	0.73	-0.41	2.29	3.24*	2.58	-0.95	1.29
		L15 ($\alpha = 1.0$)	-4.80	-0.53	0.73	1.85	-0.84	-1.01	-0.02	2.47*	2.04	-0.36	-0.53	-0.09
		L23 ($\alpha = 1.0$)	4.87*	-0.46	0.33	-1.59	0.42	0.22	-0.58	3.68*	2.71*	0.31	0.75	0.97
Quality	Gemma 2 27b	L10 ($\alpha = 100.0$)	1.37*	0.14	-0.18	3.66*	-1.12	0.82	-1.87	0.40	3.85*	6.54*	10.63*	2.20
		L22 ($\alpha = 100.0$)	-4.11	-0.63	-0.24	0.33	-1.13	1.06	4.49*	-1.22	1.51	4.51*	8.08*	1.15
		L34 ($\alpha = 50.0$)	-1.48	-0.30	-0.70	1.23	-1.62	-0.59	-1.65	-3.16	7.31*	4.67*	8.21*	1.08
	Gemma 2 9b	L9 ($\alpha = 75.0$)	-0.47	4.40*	-0.48	1.05	5.06*	2.33*	-0.99	0.39	16.44*	11.95*	4.87*	4.05
		L20 ($\alpha = 50.0$)	4.00*	-1.58	0.21	-0.83	-0.24	-0.16	-0.74	-0.97	0.93	3.88*	2.19	0.61
		L31 ($\alpha = 50.0$)	2.24*	3.84*	-0.44	-2.26	-0.34	-0.38	-1.44	-0.13	5.90*	0.85	1.00	0.80
	Llama3.1 70b	L17 ($\alpha = 3.0$)	2.48*	2.40*	1.64	0.36	1.20	0.94	-1.41	3.62*	10.51*	0.35	6.84*	2.59
		L37 ($\alpha = 1.0$)	4.90*	0.96	-2.85	1.01	1.94	-0.08	-4.45	5.10*	6.27*	2.65*	4.53*	1.82
		L57 ($\alpha = 2.0$)	2.77*	2.78*	-0.14	-0.13	1.09	-0.85	-2.49	0.85	5.49*	-1.39	5.55*	1.23
	Llama3.1 8b	L7 ($\alpha = 2.0$)	9.42*	0.63	-0.12	0.81	1.00	-0.69	0.11	3.95*	1.28	0.61	0.89	1.63
		L15 ($\alpha = 3.0$)	15.02*	0.15	0.06	-0.69	-0.71	-0.48	-0.73	7.31*	-1.45	2.07	0.21	1.89
		L23 ($\alpha = 2.0$)	12.16*	1.07	0.47	-0.77	-0.07	0.23	-0.18	4.90*	2.46*	0.06	-1.34	1.73

Table 9: Mean F1 difference across languages, models, steering methods, and layers for **topic** detection, with **positive** and **negative** differences compared to a **few-shot** prompt with no steering. Cells with * denote statistically significant increases over baseline.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Language	Gemma 2 27b	L10 ($\alpha = 25.0$)	0.79	1.67	2.26*	-0.96	-0.80	-1.21	-2.48	1.66	-0.03	2.71*	-2.79	0.07
		L22 ($\alpha = 50.0$)	-1.95	0.10	0.49	-0.72	3.00*	0.34	-2.57	-0.11	0.20	2.03	-1.88	-0.10
		L34 ($\alpha = 25.0$)	0.94	-0.30	-0.27	-1.79	1.15	1.06	-2.30	2.13	0.40	2.72*	-2.32	0.13
	Gemma 2 9b	L9 ($\alpha = 25.0$)	-3.88	-0.84	-0.02	3.36*	-0.06	-0.02	-1.57	-1.04	-1.07	1.16	-1.92	-0.54
		L20 ($\alpha = 50.0$)	0.34	1.62	2.44*	-0.31	-0.03	1.23	-0.01	-3.62	0.54	1.11	-3.46	-0.01
		L31 ($\alpha = 50.0$)	-1.10	-0.27	1.33	-0.32	-0.23	0.96	-3.89	-1.65	-2.08	0.78	-0.22	-0.61
	Llama3.1 70b	L17 ($\alpha = 1.0$)	2.84*	-1.26	-0.60	-1.43	0.77	0.14	3.42*	-0.72	0.63	-1.26	0.51	0.28
		L37 ($\alpha = 2.0$)	2.14	2.12	2.28	-1.41	1.44	0.00	2.50*	2.95*	1.06	-0.05	4.27*	1.57
		L57 ($\alpha = 1.0$)	2.70*	1.58	1.57	0.31	3.56*	-0.57	2.40	-1.56	0.15	-0.23	4.32*	1.29
	Llama3.1 8b	L7 ($\alpha = 2.0$)	-1.99	-2.86	-1.20	3.71*	-0.04	1.97	-1.72	3.69*	-2.29	3.01*	-0.38	0.17
		L15 ($\alpha = 2.0$)	2.98*	-2.92	-2.71	2.11	-1.40	-1.23	2.03	1.73	-0.12	2.02	0.41	0.26
		L23 ($\alpha = 3.0$)	-0.11	0.26	0.51	4.02*	0.57	4.20*	1.25	0.03	-0.00	0.41	-0.06	1.01
Quality	Gemma 2 27b	L10 ($\alpha = 75.0$)	-1.63	0.95	2.81*	-0.36	0.89	1.28	-0.06	2.93*	0.93	2.69*	-0.45	0.85
		L22 ($\alpha = 75.0$)	-1.70	2.62*	1.82	0.74	1.66	0.94	-1.32	0.92	0.16	3.16*	-1.89	0.65
		L34 ($\alpha = 100.0$)	-4.22	1.04	1.60	-1.92	-0.41	0.75	-4.16	0.07	0.91	2.67*	-0.09	-0.34
	Gemma 2 9b	L9 ($\alpha = 25.0$)	1.41	1.64	2.47*	1.42	0.76	1.61	-0.28	2.91*	-1.13	1.45	-0.59	1.06
		L20 ($\alpha = 50.0$)	-6.12	-0.27	2.53*	1.22	-0.96	2.11	2.17	-0.33	0.08	1.91	-3.86	-0.14
		L31 ($\alpha = 50.0$)	-0.72	1.26	2.44*	3.37*	-1.01	1.33	0.99	-3.02	0.86	0.32	-2.58	0.29
	Llama3.1 70b	L17 ($\alpha = 1.0$)	5.20*	1.46	0.13	-0.60	2.57*	-0.99	1.53	-1.60	2.54*	1.90	1.64	1.24
		L37 ($\alpha = 1.0$)	0.40	-0.17	1.93	-0.59	0.71	-2.32	0.86	-0.32	-0.72	1.29	4.22*	0.48
		L57 ($\alpha = 2.0$)	0.18	-0.13	1.44	-3.35	3.27*	-3.08	1.10	2.04	2.74*	0.26	4.24*	0.79
	Llama3.1 8b	L7 ($\alpha = 2.0$)	0.84	2.97*	-1.68	2.60*	1.07	1.22	0.69	4.02*	-0.47	-0.50	0.14	0.94
		L15 ($\alpha = 3.0$)	2.21*	-0.15	-1.46	1.95	-2.91	1.06	-0.25	5.30*	-2.54	3.26*	-0.52	0.54
		L23 ($\alpha = 4.0$)	7.47*	0.50	-2.35	2.80*	-2.79	1.22	-0.30	4.30*	-1.89	-1.20	-0.05	0.70

Table 10: Mean F1 difference across languages, models, steering methods, and layers for **sentiment** detection, with **positive** and **negative** differences compared to a **few-shot** prompt with no steering. Cells with * denote statistically significant increases over baseline.

Method	LLM	Layer	am	cs	da	de	he	id	jv	mt	sk	sl	su	Average
Language	Gemma 2 27b	L10 ($\alpha = 25.0$)	2.36	3.30*	-1.01	4.42*	7.30*	3.94*	0.55	2.87*	-0.75	4.40*	-0.73	2.42
		L22 ($\alpha = 50.0$)	0.68	4.86*	-0.49	-0.45	0.71	1.75	0.18	0.12	-0.60	9.00*	-0.83	1.36
		L34 ($\alpha = 25.0$)	1.00	5.47*	0.94	2.48	-1.66	3.52*	0.93	1.19	-1.90	1.92	0.57	1.32
	Gemma 2 9b	L9 ($\alpha = 25.0$)	-4.40	1.52	0.49	-0.40	-1.02	-3.06	2.58*	-2.90	5.19*	-3.48	0.95	-0.41
		L20 ($\alpha = 50.0$)	-0.67	0.30	-0.05	2.04	-0.02	-1.65	1.94	-0.64	6.18*	-2.05	0.56	0.54
		L31 ($\alpha = 50.0$)	-3.60	0.76	6.64*	-1.22	0.27	-0.89	4.28*	-0.35	6.96*	-2.94	1.20	1.01
	Llama3.1 70b	L17 ($\alpha = 1.0$)	0.85	0.79	0.34	3.25*	1.86	-0.31	-2.79	1.32	-3.03	2.13	3.57*	0.73
		L37 ($\alpha = 2.0$)	-3.53	0.45	-0.50	2.15	5.47*	2.52*	-2.35	1.48	2.16	1.07	1.44	0.94
		L57 ($\alpha = 1.0$)	-1.44	-0.22	0.08	1.55	2.36	3.40*	-2.89	1.98	-0.02	3.08*	3.49*	1.03
	Llama3.1 8b	L7 ($\alpha = 2.0$)	-0.47	2.38*	6.19*	-1.25	0.70	1.01	-1.60	2.21	-0.14	0.60	-4.13	0.50
		L15 ($\alpha = 2.0$)	-0.20	3.16*	0.53	0.53	-0.74	-0.02	-0.62	2.99*	-0.65	1.65	-2.21	0.40
		L23 ($\alpha = 3.0$)	-0.17	-0.41	3.04*	-0.56	0.13	1.01	0.15	0.22	-0.34	1.79	-2.51	0.21
Quality	Gemma 2 27b	L10 ($\alpha = 75.0$)	0.74	3.42*	2.86*	2.86*	0.54	3.24*	-0.50	1.30	3.69*	6.78*	-2.15	2.07
		L22 ($\alpha = 75.0$)	1.79	4.75*	-3.07	3.38*	-2.57	4.38*	-0.58	0.70	1.82	3.48*	-1.55	1.14
		L34 ($\alpha = 100.0$)	1.74	0.05	-1.73	3.18*	5.48*	1.63	-1.77	2.34	2.69	4.28*	-1.02	1.53
	Gemma 2 9b	L9 ($\alpha = 25.0$)	0.06	0.10	3.04*	-1.82	2.73*	-0.86	3.42*	3.47*	2.35*	0.11	0.42	1.18
		L20 ($\alpha = 25.0$)	-0.73	-1.18	-0.39	-1.33	0.97	0.52	0.09	1.29	-0.56	-1.95	0.01	-0.29
		L31 ($\alpha = 50.0$)	-2.30	-2.17	2.91*	-0.75	1.55	-0.10	-6.38	0.62	6.64*	-1.67	2.35*	0.06
	Llama3.1 70b	L17 ($\alpha = 1.0$)	-2.10	-0.10	0.48	3.49*	3.18*	2.47*	1.22	0.67	1.04	0.34	3.38*	1.28
		L37 ($\alpha = 1.0$)	-1.81	1.57	-0.33	0.53	4.97*	1.16	-1.42	-0.55	1.56	-1.47	5.32*	0.87
		L57 ($\alpha = 2.0$)	-2.97	2.18	-0.51	1.56	-0.78	2.05	-0.20	2.83*	1.03	0.60	6.72*	1.14
	Llama3.1 8b	L7 ($\alpha = 2.0$)	1.23	0.67	5.08*	0.01	0.37	0.79	-3.41	2.24*	0.55	-0.85	-0.05	0.60
		L15 ($\alpha = 3.0$)	0.77	0.78	1.29	1.17	0.16	-1.11	-4.53	0.82	-9.74	-4.04	-1.27	-1.43
		L23 ($\alpha = 4.0$)	0.63	-1.35	-1.09	0.18	-0.28	0.20	-1.86	1.65	-2.14	-0.23	0.44	-0.35

Table 11: Boxplot visualizations of relative increases/decreases in % for F1 metrics for LLM + steering methods for alpha values in **zero-shot** setting. The X axis in the figures represents the different alpha values for that LLM + steering combination aggregated over languages and tasks, and the Y axis represents the relative difference increase for F1 downstream model performance.

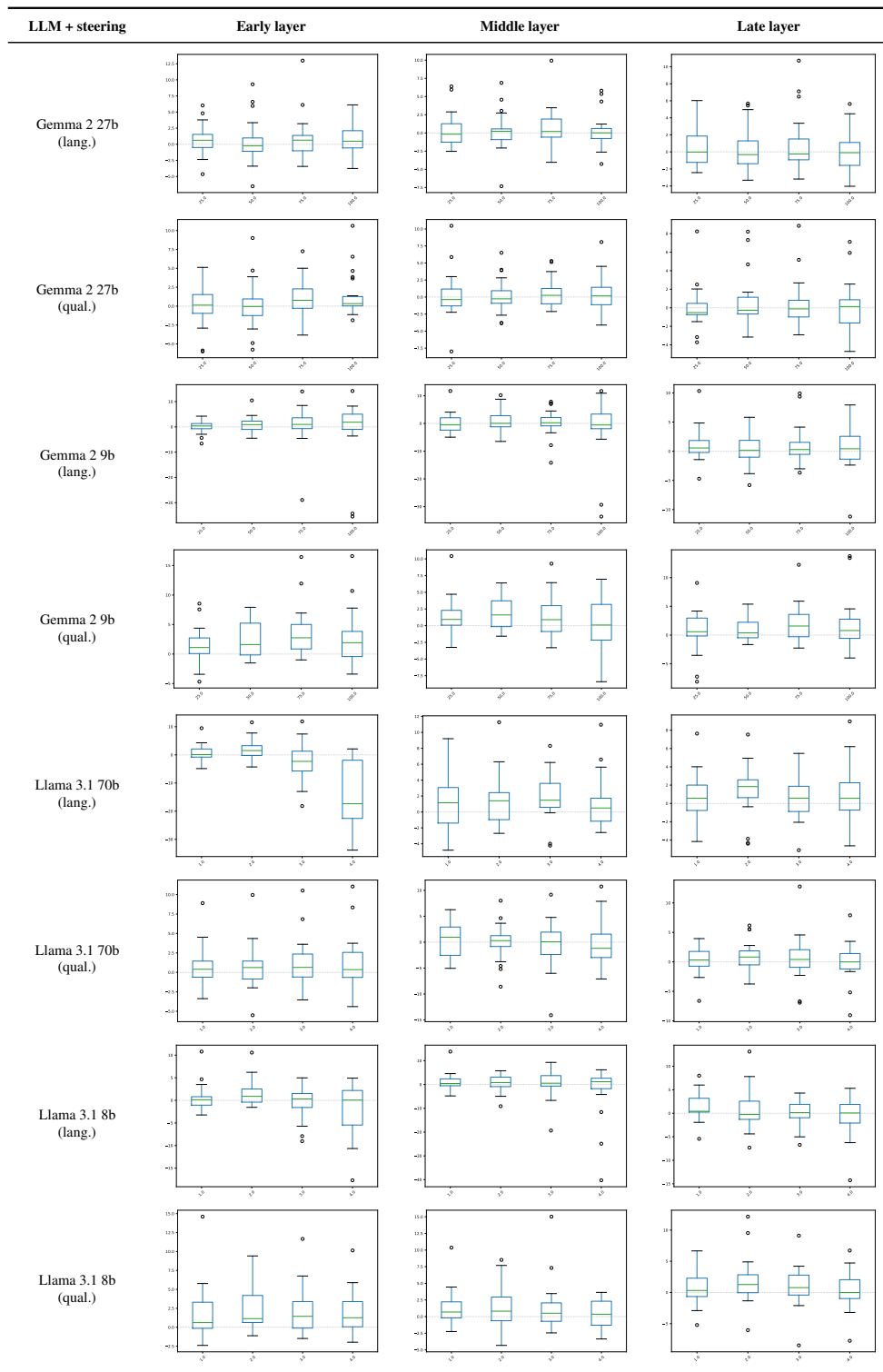


Table 12: Boxplot visualizations of relative increases/decreases in % for F1 metrics for LLM + steering methods for alpha values in **few-shot** setting. The X axis in the figures represents the different alpha values for that LLM + steering combination aggregated over languages and tasks, and the Y axis represents the relative difference increase for F1 downstream model performance.

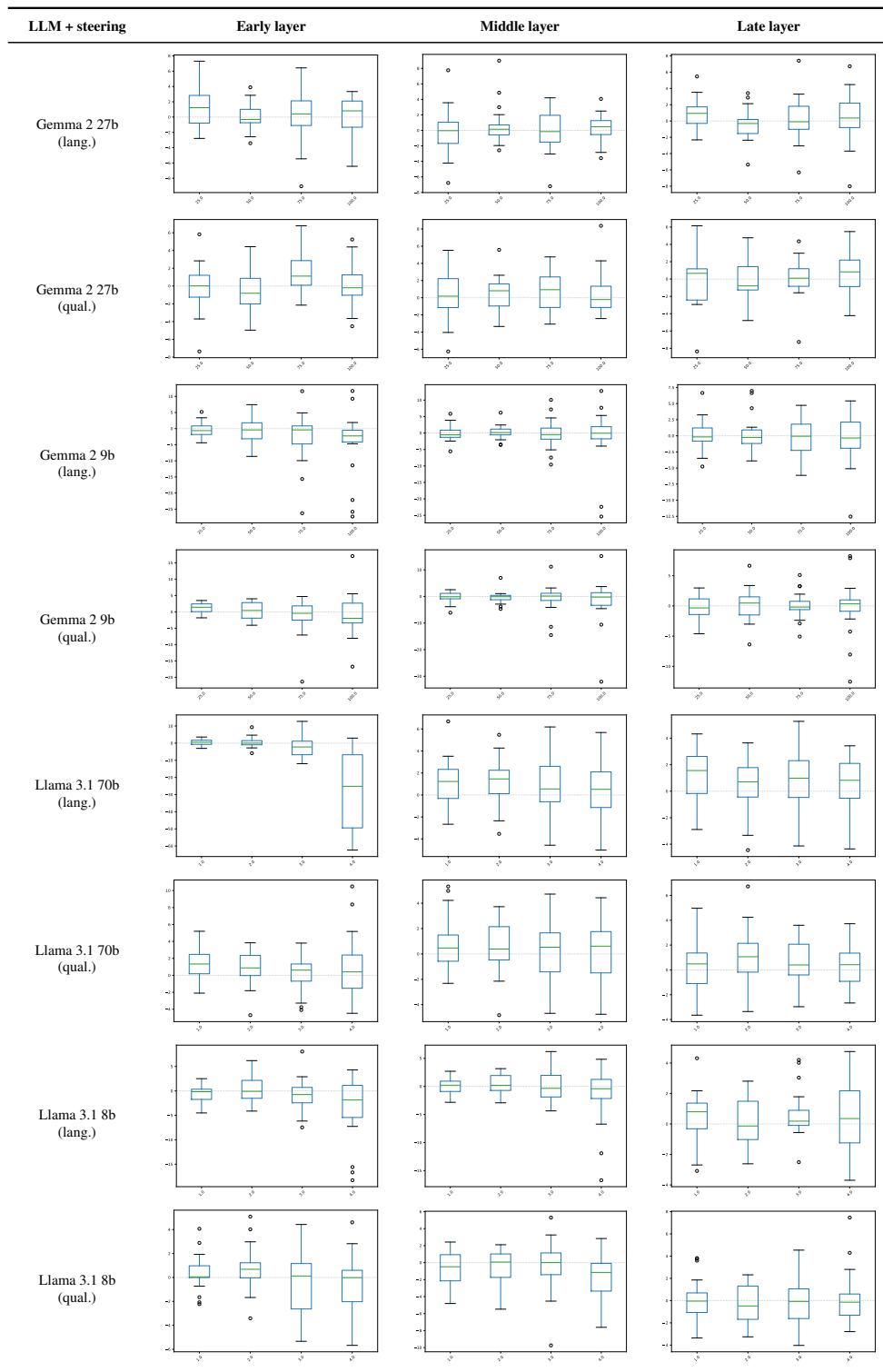


Table 13: Boxplot visualizations of relative increases/decreases in % for F1 metrics for LLM + steering methods for various languages in **zero-shot** setting. The boxplots are for α values from Table 1. The X axis in the figures represents the languages, and the Y axis represents the relative difference increase for F1 downstream model performance.

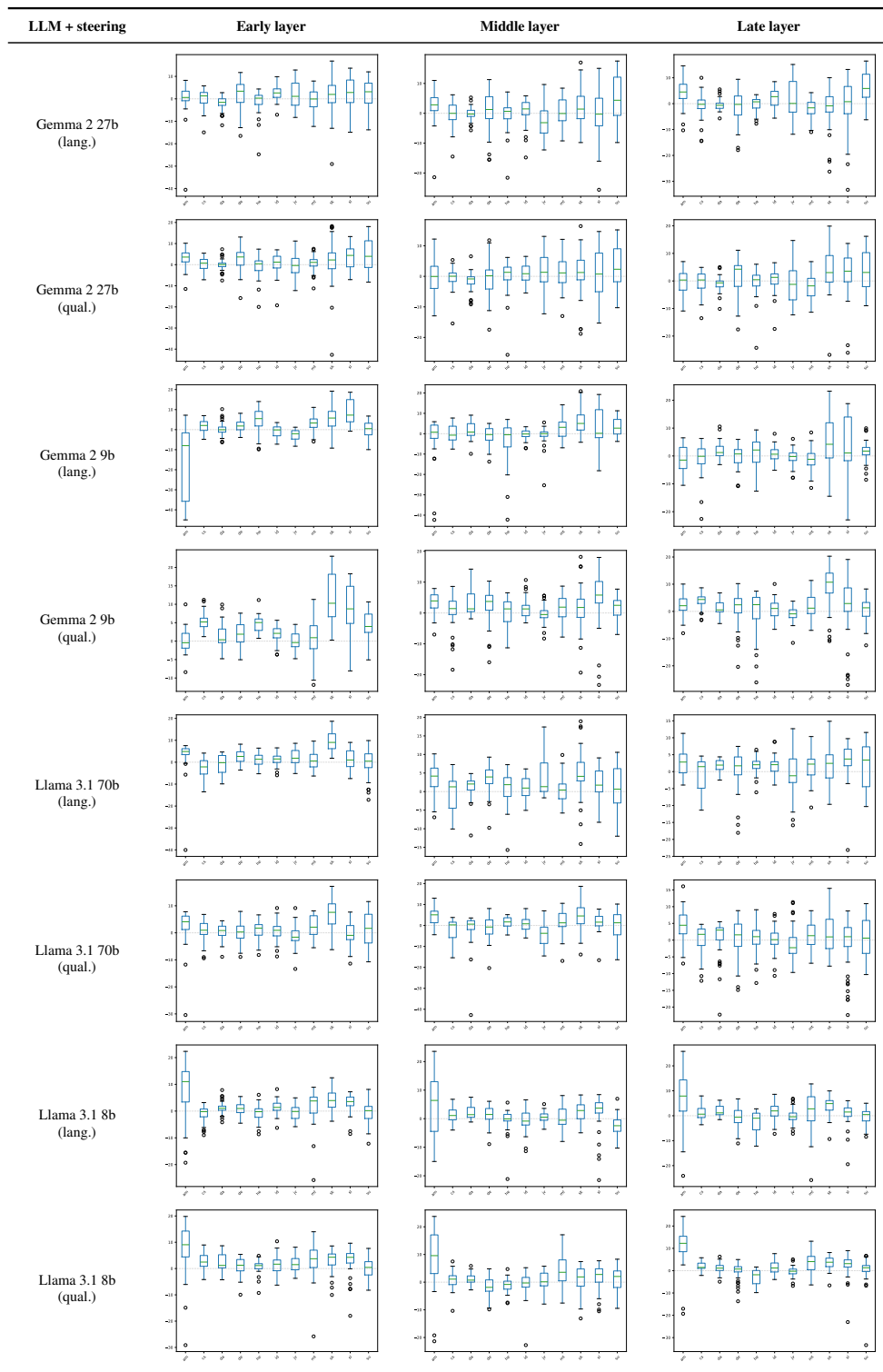


Table 14: Boxplot visualizations of relative increases/decreases in % for F1 metrics for LLM + steering methods for various languages in **few-shot** setting. The boxplots are for α values from Table 1. The X axis in the figures represents the languages, and the Y axis represents the relative difference increase for F1 downstream model performance.

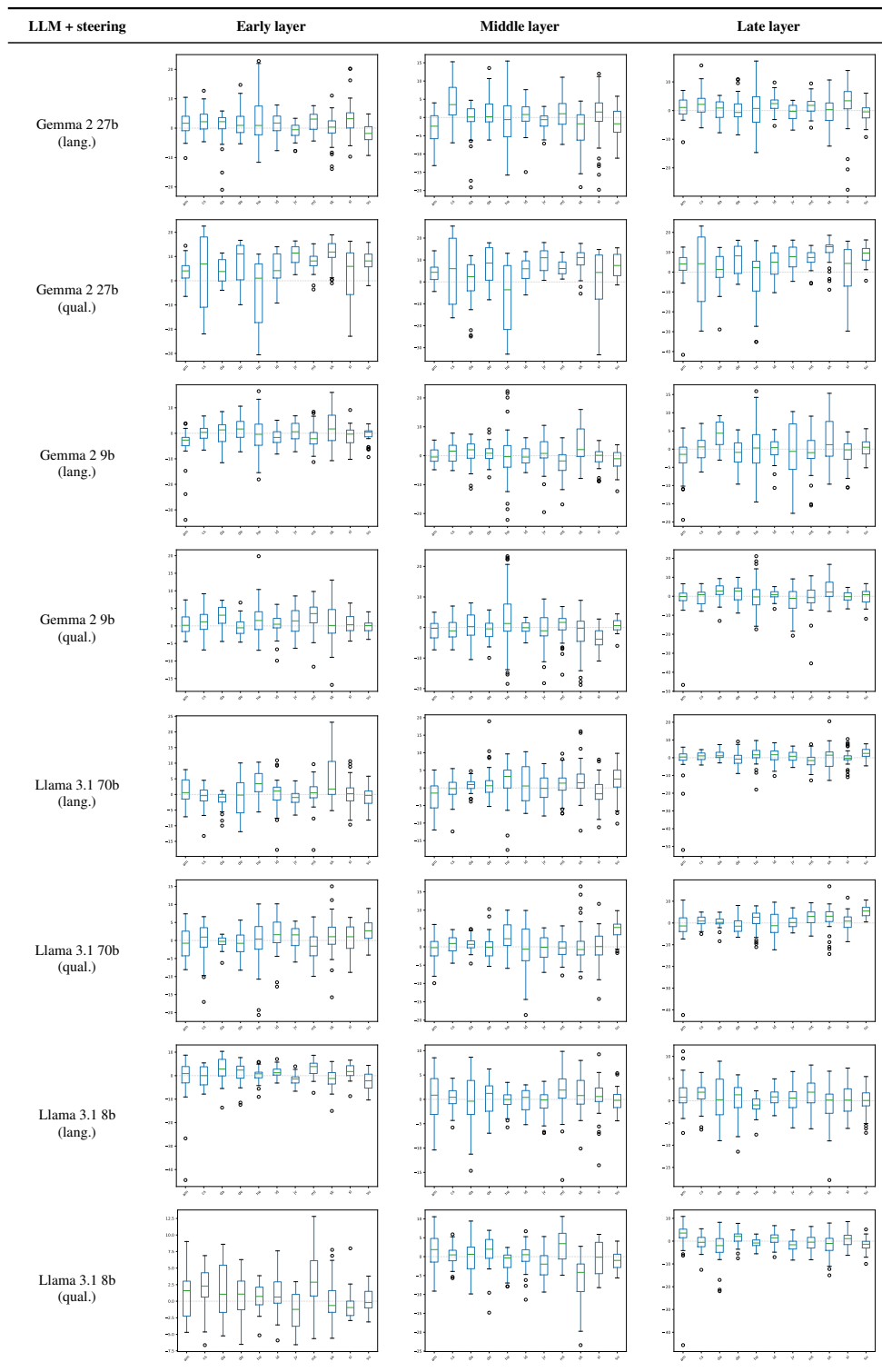


Table 15: Relative increases/decreases in % for diversity metrics aggregated over LLMs and alphas for various languages and layers for **zero-shot** baseline. The X axis in the figures represents the languages, and the Y axis represents the relative difference increase for that given metric.

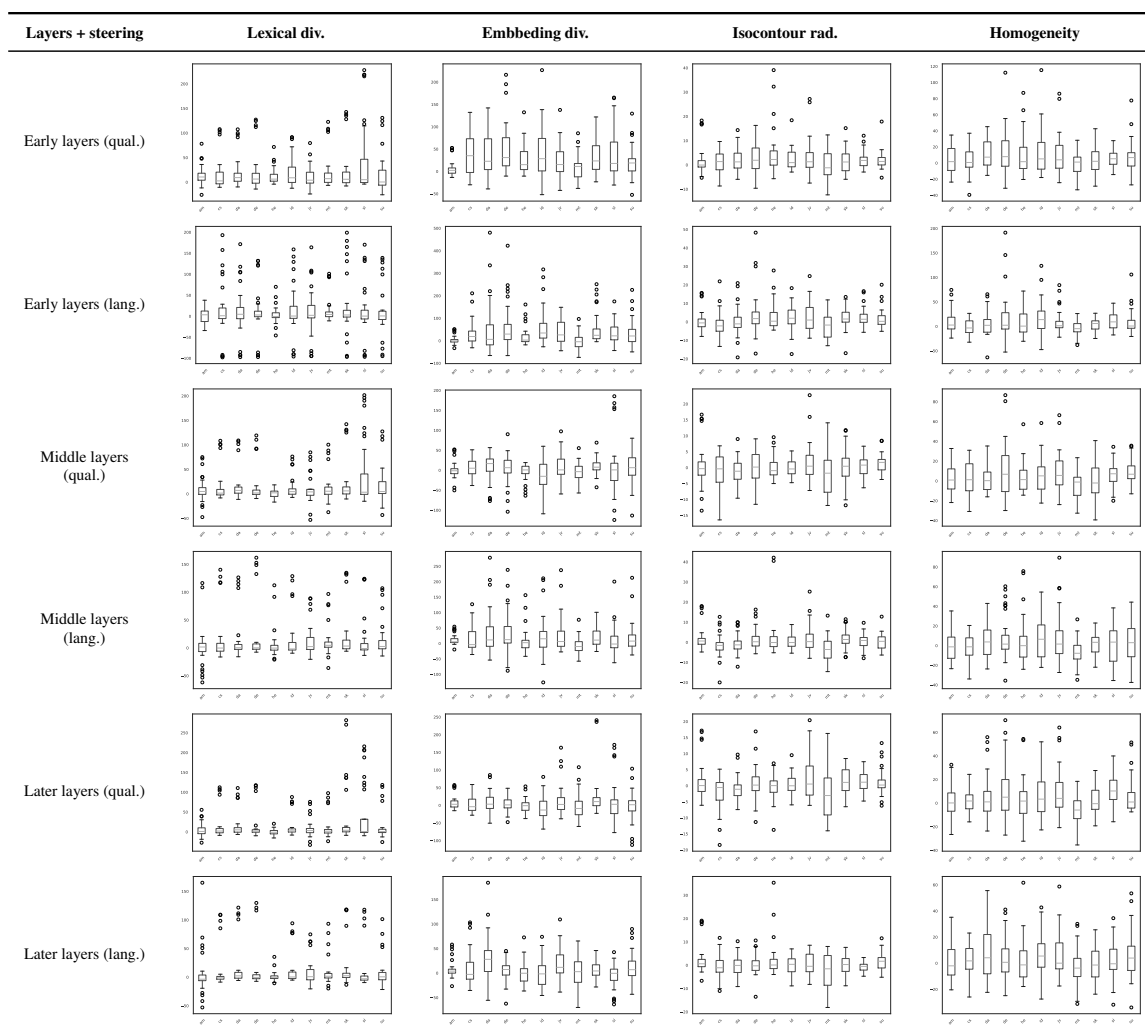


Table 16: Relative increases/decreases in % for diversity metrics aggregated over LLMs and alphas for various languages and layers for **few-shot** baseline. The X axis in the figures represents the languages, and the Y axis represents the relative difference increase for that given metric.

