

Language Steering for Multilingual In-Context Learning

Neeraja Kirtane¹ Kuan-Hao Huang¹

Abstract

While multilingual large language models have gained widespread adoption, their performance on non-English languages remains substantially inferior to English. This disparity is particularly evident in in-context learning scenarios, where providing demonstrations in English but testing on non-English inputs leads to significant performance degradation. In this paper, we hypothesize that LLMs develop a universal semantic space for understanding languages, where different languages are encoded as distinct directions within this space. Based on this hypothesis, we propose *language vectors*—a training-free language steering approach that leverages activation differences between source and target languages to guide model behavior. We steer the model generations by adding the vector to the intermediate model activations during inference. This is done to make the model’s internal representations shift towards the target language space without any parameter updates. We evaluate our method across three datasets and test on a total of 19 languages on three different models. Our results show consistent improvements on multilingual in-context learning over baselines across all tasks and languages tested. Beyond performance gains, hierarchical clustering of steering vectors reveals meaningful linguistic structure aligned with language families. These vectors also successfully transfer across tasks, demonstrating that these representations are task-agnostic.

1. Introduction

The rapid advancement of multilingual large language models (LLMs) has enabled cross-lingual applications at an unprecedented scale. However, despite their training on diverse language corpora, these models exhibit a persistent

performance gap between English and low-resource languages across various downstream tasks (Zhao et al., 2024). This disparity poses significant challenges for equitable access to language technologies. A key question in this space is how to achieve effective cross-lingual transfer: the ability to apply knowledge learned in high-resource languages (typically English, where abundant high-quality data exists) to improve performance in low-resource languages.

In-context learning (ICL) naturally provides an important scenario for cross-lingual transfer. When demonstrations are provided in a source language (English) where high-quality data exists but the test query is in a low-resource target language, models must internally process the input, recognize the language switch, and transfer the learned task pattern across languages. This cognitively demanding process often fails. The model needs to interpret the source language demonstrations to understand the task structure, then apply this understanding to a query in a different language while maintaining reasoning coherence. This cross-lingual in-context learning gap significantly limits the practical deployment of LLMs in multilingual contexts, particularly for complex reasoning tasks where high-quality demonstrations are primarily available in source languages like English.

Recent work in mechanistic interpretability has revealed that transformer-based models encode distinct patterns for different types of information within their internal representations (Rai et al., 2024; Li et al., 2024; Elhage et al., 2022). Building on these insights, we investigate whether language-specific information is similarly encoded in identifiable activation patterns that can be leveraged to improve the performance of low-resource languages. Specifically, we ask: *Can we extract and apply language-specific steering vectors to change the “language mode” of models, and therefore enhancing multilingual performance without additional training?*

We propose a simple yet effective training-free approach that computes language steering vectors from raw activation differences between source and target language examples. As shown in Figure 1, our method operates in a few-shot in-context learning setting, where we: (1) collect paired examples in the source and the target language, (2) compute the difference in model activations across these pairs to extract a *language-specific steering vector*, and (3) apply this

¹Texas A&M University. Correspondence to: Neeraja Kirtane <kirtane.neeraja@gmail.com>, Kuan-Hao Huang <khhuang@tamu.edu>.

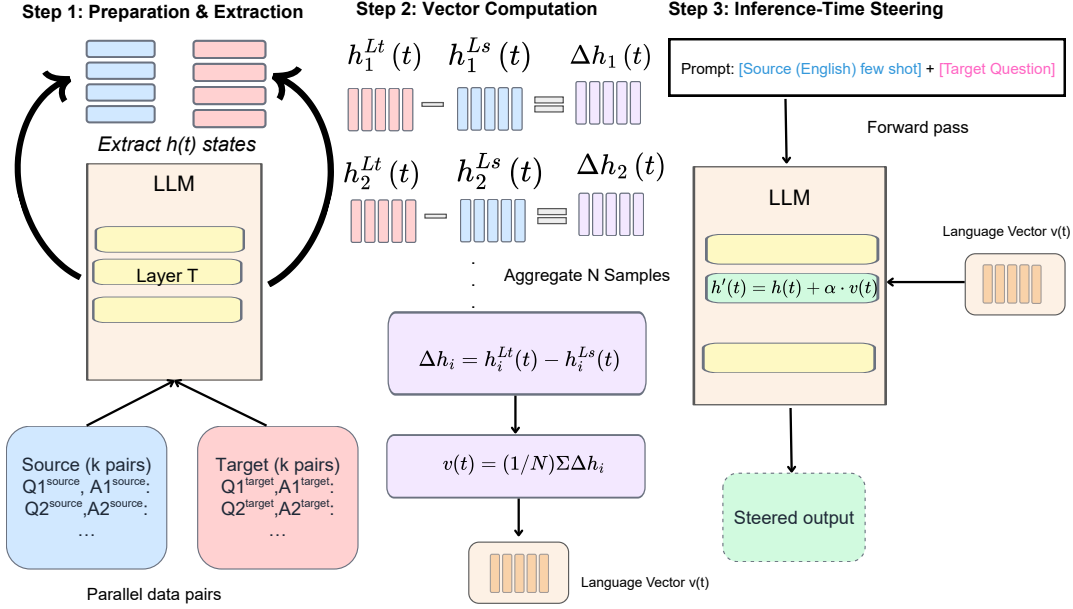


Figure 1. Overview of our language steering approach for multilingual in-context learning. **Step 1 (Extraction)**: We extract hidden states from parallel source and target language question-answer pairs at layer t of the LLM. **Step 2 (Vector Computation)**: The language steering vector $v(t)$ is computed offline as mean of activation differences between target and source language representations across N samples. **Step 3 (Inference-Time Steering)**: During inference with few-shot demonstrations from source language and a target language question, we steer the model by adding $\alpha \cdot v(t)$ to the hidden states at layer t , guiding the model toward target language reasoning patterns.

vector to steer the model’s internal representations during inference on downstream tasks. This language steering vector can help us change the “language mode” of a model from the source to the target language to perform more efficiently. Critically, this approach requires no parameter updates or fine-tuning, making it practical for scenarios where labeled data in the target language is scarce or unavailable.

We evaluate our approach across three diverse datasets, which encompass mathematical reasoning and natural language inference spanning 19 languages over three models. Our results show consistent improvements over baselines, with particularly strong gains on mathematical reasoning. Beyond performance improvements, our analysis reveals that language-specific steering vectors encode meaningful linguistic structure. Closely-related languages produce similar steering vectors. We also observe that middle layers are the most effective in steering, and is also effective for cross-task setups.

The contributions of our work are threefold: (1) We introduce language vectors, a training-free steering approach that leverages activation differences between source and target languages to effectively change the language mode of a model. (2) We provide comprehensive empirical evaluation demonstrating the effectiveness of language vectors across diverse tasks (mathematical reasoning and natural language inference), 19 languages, and three models (3) We show that

these language vectors encode meaningful linguistic structure that can help us better understand multilingual LLMs. The vectors are also effective in cross-task steering. Our code is publicly available.¹

2. Language Steering for In-Context Learning

We consider the multilingual in-context learning (Tu et al., 2025), where each instance consists of few-shot source (English) demonstrations and a test question in the target language. We hypothesize that by adding steering vectors containing language-specific information to internal representations can achieve performance comparable to the setup of using few-shot demonstrations in the target language.

2.1. Problem Formulation and Notations

Let L_s denote the source language and L_t denote the target language. In the multilingual in-context learning setup, the model receives a prompt consisting of k few-shot demonstration examples $\mathcal{F} = \{(q_i^{L_s}, c_i^{L_s}, a_i^{L_s})\}_{i=1}^k$ followed by a test question $q_{\text{test}}^{L_t}$. Each demonstration includes a question $q_i^{L_s}$, chain-of-thought reasoning $c_i^{L_s}$, and answer $a_i^{L_s}$, all in the source language L_s . The model must leverage the task understanding from these source language demonstrations

¹<https://github.com/lab-flair/language-vector>

to generate both the reasoning $c_{\text{test}}^{L_t}$ and answer $a_{\text{test}}^{L_t}$ for the target language.

Given a language model \mathcal{M} with T layers, we denote the hidden state at layer $t \in \{1, \dots, T\}$ and token position j for input text x as $\mathbf{h}_j(t) \in \mathbb{R}^d$, where d is the hidden dimension.

2.2. Steering Vector Computation

Our goal is to compute a language-specific steering vector $\mathbf{v}(t) \in \mathbb{R}^d$ that, when added to hidden states during inference, improves the model’s performance on target language inputs q^{L_t} than the baseline. We aim to achieve this using source language demonstrations along with steering.

We compute language-specific steering vectors by analyzing activation patterns across parallel examples in English (source language L_s) and a target language L_t . Our method extracts these vectors using a few-shot format that captures richer contextual patterns.

Few-shot Format for Steering Vector Computation. To capture richer contextual patterns across multiple examples for a language, we create samples containing k question-answer pairs concatenated together. For each sample i in our compute set $\mathcal{D}_{\text{compute}}$, we format parallel texts as:

$$x_i^{L_s} = \bigoplus_{j=1}^k [q_{i,j}^{L_s}, a_{i,j}^{L_s}] \quad (1)$$

$$x_i^{L_t} = \bigoplus_{j=1}^k [q_{i,j}^{L_t}, a_{i,j}^{L_t}] \quad (2)$$

where $q_{i,j}^{L_s}$ and $q_{i,j}^{L_t}$ are parallel questions in source and target languages respectively, $a_{i,j}^{L_s}$ and $a_{i,j}^{L_t}$ are the answers, \oplus denotes concatenation, and k is the number of question-answer pairs per sample (we use $k = 6$ in our experiments).

Activation Extraction and Vector Computation. We extract hidden states at layer t for both source and target formatted texts. For each sample i , we pass the formatted text through model \mathcal{M} and extract the hidden states at layer t :

$$\mathbf{h}_i^{L_s}(t) = \frac{1}{|x_i^{L_s}|} \sum_{j=1}^{|x_i^{L_s}|} \mathbf{h}_{i,j}^{L_s}(t) \quad (3)$$

$$\mathbf{h}_i^{L_t}(t) = \frac{1}{|x_i^{L_t}|} \sum_{j=1}^{|x_i^{L_t}|} \mathbf{h}_{i,j}^{L_t}(t) \quad (4)$$

where $\mathbf{h}_{i,j}^{L_s}(t)$ denotes the hidden state at token position j of sample i in the source language at layer t , and $|x_i^{L_s}|$ is the total sequence length. We apply mean-pooling across all token positions to obtain a single representation per sample.

The steering vector $\mathbf{v}(t) \in \mathbb{R}^d$ at layer t is then computed as the mean of activation differences:

$$\mathbf{v}(t) = \frac{1}{N} \sum_{i=1}^N \left(\mathbf{h}_i^{L_t}(t) - \mathbf{h}_i^{L_s}(t) \right) \quad (5)$$

where $N = |\mathcal{D}_{\text{compute}}|$ is the number of samples in the compute set. This setup ensures that adding the steering vector shifts activations from the source language distribution (English) toward the target language distribution, effectively steering the model’s internal representations to behave as if processing target language demonstrations.

2.3. Inference-Time Steering

During inference, we apply the pre-computed steering vector to hidden states at specific token positions within the input prompt. Given a test example from $\mathcal{D}_{\text{test}}$ with few-shot demonstrations \mathcal{F} in the source language and a target language question $q_{\text{test}}^{L_t}$, the full input prompt is:

$$\text{prompt} = [\text{system.instruction}] \oplus \mathcal{F} \oplus q_{\text{test}}^{L_t} \quad (6)$$

where $\mathcal{F} = \{(q_1^{L_s}, c_1^{L_s}, a_1^{L_s}), \dots, (q_k^{L_s}, c_k^{L_s}, a_k^{L_s})\}$ consists of k few-shot demonstrations with source language questions $q_i^{L_s}$, source language chain-of-thought reasoning $c_i^{L_s}$, and answers a_i .

Steering Positions. We define a set of token positions \mathcal{P} where steering is applied. We experiment with four configurations. More ablations will be in Section 4.2.

- `on_fewshot`: $\mathcal{P} = \{p : p \in \text{tokens}(\mathcal{F})\}$: steer on all few-shot demonstration tokens.
- `after_fewshot`: $\mathcal{P} = \{p : p = \text{first token after } \mathcal{F}\}$: steer only on the boundary between demonstrations and test question.
- `on_question`: $\mathcal{P} = \{p : p \in \text{tokens}(q_{\text{test}}^{L_t})\}$: steer only on test question tokens.
- `entire`: $\mathcal{P} = \{p : p \in \text{tokens}(\text{prompt})\}$: steer on all prompt tokens.

Modified Forward Pass. We modify the hidden states for each position $p \in \mathcal{P}$ during the forward pass at layer t :

$$\mathbf{h}'_p(t) = \mathbf{h}_p(t) + \alpha \cdot \mathbf{v}(t) \quad (7)$$

where $\mathbf{h}_p(t) \in \mathbb{R}^d$ is the original hidden state at position p and layer t , $\alpha \in \mathbb{R}$ is a scaling hyperparameter controlling steering strength, and $\mathbf{h}'_p(t)$ is the steered hidden state. This modification is applied via forward hooks during generation, requiring no parameter updates or gradient computation. The hyperparameters $(t^*, \alpha^*, \mathcal{P}^*)$ are selected based on validation set performance on \mathcal{D}_{val} .

Table 1. Detailed per-language accuracy results on MGSM, XNLI, and MSVAMP for the Llama-3.1-8b-instruct model. Gray dashes indicate that the dataset does not contain that language. B = Baseline, MFS = Multilingual Few-Shot Baseline (Tu et al., 2025), OR = Oracle. Oracle is the upper bound and is not for a direct comparison **Bold** indicates the best score per language.

Language	MGSM				XNLI				MSVAMP			
	B	MFS	Ours	OR	B	MFS	Ours	OR	B	MFS	Ours	OR
Arabic	—	—	—	—	62.57	60.78	62.87	63.77	—	—	—	—
Basque	32.14	35.70	36.90	52.38	—	—	—	—	—	—	—	—
Bengali	57.14	55.95	61.90	58.33	—	—	—	—	57.49	62.87	59.58	61.38
Bulgarian	—	—	—	—	56.29	61.98	61.68	66.17	—	—	—	—
Catalan	64.29	64.28	69.05	76.19	—	—	—	—	—	—	—	—
Chinese	67.86	67.86	71.43	72.62	59.88	61.38	59.28	61.98	69.76	73.05	71.26	73.35
French	61.90	64.29	70.24	65.48	67.37	68.26	72.75	71.26	71.56	73.05	74.55	73.05
Galician	64.29	69.04	73.81	77.38	—	—	—	—	—	—	—	—
German	66.67	66.67	75.00	71.43	64.07	65.27	66.47	65.27	71.26	70.06	71.26	76.95
Greek	—	—	—	—	68.26	66.17	70.06	72.75	—	—	—	—
Hindi	—	—	—	—	61.38	57.78	64.97	61.98	—	—	—	—
Japanese	55.95	61.90	55.95	63.10	—	—	—	—	63.17	68.26	67.96	70.06
Russian	71.43	71.43	72.62	76.19	58.98	64.37	63.77	60.78	68.86	72.46	72.16	71.56
Spanish	77.38	70.24	76.19	78.57	66.77	66.17	70.66	67.37	74.55	73.95	75.75	74.55
Swahili	55.95	63.10	65.48	66.67	52.40	51.20	55.99	56.89	56.29	59.58	60.78	62.87
Thai	57.14	67.86	61.90	59.52	59.88	64.37	60.78	70.66	59.58	64.67	64.07	66.77
Turkish	—	—	—	—	62.28	57.19	65.87	65.57	—	—	—	—
Urdu	—	—	—	—	55.09	55.39	56.29	57.19	—	—	—	—
Vietnamese	—	—	—	—	64.07	63.17	68.86	68.86	—	—	—	—
Average	61.01	63.19	65.87	68.16	61.38	61.68	64.31	65.04	65.84	68.66	68.60	70.06

3. Experiments

3.1. Setup

We test three instruction-tuned models: *Llama3.1-8b-Instruct*, *Qwen2.5-7b-Instruct* and *Qwen2.5-14b-Instruct* (Grattafiori et al., 2024; Team et al., 2024) and consider three datasets: MGSM (Shi et al., 2022), MSVAMP (Chen et al., 2024) and XNLI (Conneau et al., 2018).

Data Splits. We partition the test set into three equal parts: compute ($\mathcal{D}_{\text{compute}}$) for steering vector calculation, validation (\mathcal{D}_{val}) for hyperparameter and configuration selection, and test ($\mathcal{D}_{\text{test}}$) for final evaluation. Additionally, we sample 6-8 examples from the training set to serve as few-shot demonstrations in our prompts. These few shot prompts are the same across all the experiments for a particular dataset.

Implementation Details. While calculating the vector, we create $N = |\mathcal{D}_{\text{compute}}|$ samples, ensuring each example from the compute set appears at least once across all samples, with remaining slots filled through random sampling with replacement. This few-shot format provides the model with a more context-rich representation that may better capture language-specific patterns. We perform a grid search over steering layers $\ell \in \{5, 10, 15, 20, 25, 30\}$, scaling factors $\alpha \in \{0.5, 1.0, 2.0, 3.0\}$, and the four steering position configurations (on_fewshot, after_fewshot, on_question, entire) on \mathcal{D}_{val} . Only configurations achieving validation accuracy above the source baseline are evaluated on $\mathcal{D}_{\text{test}}$.

Baselines. We compare our method against two baselines: (1) *Source baseline (B)*: few-shot prompts with source questions and source chain-of-thought and answer, representing unsteered cross-lingual transfer performance; (2) *Multilingual few-shot (MFS)*: few-shot examples drawn from multiple languages with their respective chain-of-thought reasoning (Tu et al., 2025), representing the approach of diversifying few-shot examples across languages. To understand what our upper bound would be when evaluating, we also test: *Oracle*: few-shot prompts with target language questions and target answers. Oracle is an upper bound for reference and is not used for direct comparison.

Few-Shot Prompt Structure. During evaluation, our prompt contains few-shot questions in the source language (English) with answer explanations and the final answer. The question to be answered is in the target language. We steer some part of this prompt to understand whether that can help the model change the mode from English to the target language while processing the given question. For our tasks, which are reasoning and inference-based, the final answer remains the same across languages.

3.2. Results

Detailed accuracy results on ($\mathcal{D}_{\text{test}}$) for the Llama3.1-8b-Instruct model are shown in Table 1 across the three datasets. We test on a total of 19 languages across the datasets. Dashes in the table indicate that the language is not present for that data. For every language in a dataset, we compute the

Table 2. Average performance comparison across three datasets and three models. B = Baseline, MFS = Multilingual Few-Shot Baseline (Tu et al., 2025), OR = Oracle. Oracle is the upper bound and is not for a direct comparison.

Model	Performance			
	B	MFS	Ours	OR
MGSM				
LLaMA-3.1-8B-Instruct	61.01	63.19	65.87	68.16
Qwen2.5-7B-Instruct	68.95	69.94	71.83	73.21
Qwen2.5-14B-Instruct	79.40	80.86	84.80	81.23
XNLI				
LLaMA-3.1-8B-Instruct	61.38	61.68	64.31	65.04
Qwen2.5-7B-Instruct	74.81	73.27	75.51	74.53
Qwen2.5-14B-Instruct	72.73	73.91	74.66	76.45
MSVAMP				
LLaMA-3.1-8B-Instruct	65.84	68.66	68.60	70.06
Qwen2.5-7B-Instruct	75.18	76.94	76.68	77.65
Qwen2.5-14B-Instruct	82.10	82.90	83.90	83.03

accuracies for two baselines: having few-shots in source language with target questions, the multilingual few-shot method (Tu et al., 2025).

Our method consistently improves over the source baseline (B) across most languages and datasets, demonstrating that activation steering enables effective cross-lingual transfer without any parameter updates or prompt modifications. The varying magnitude of gains across languages suggests that steering effectiveness depends on language-specific characteristics rather than being uniform.

Clear patterns emerge across datasets. MGSM shows the most consistent and substantial improvements, indicating that structured mathematical reasoning particularly benefits from cross-lingual steering. On MSVAMP, the multilingual few-shot baseline (MFS) proves more competitive, often matching or slightly exceeding our method, which reflects the value of diverse demonstrations for arithmetic word problems. XNLI presents the most mixed results: gains vary considerably across languages, and MFS occasionally outperforms ours, likely due to the greater linguistic and semantic variability in natural language inference tasks. Overall, these findings show that our method offers a robust complement to multilingual few-shot prompting, with especially strong advantages for reasoning-heavy tasks.

Results for Different Models. Table 2 summarizes average performance across all languages for each model-dataset combination. Overall, our method improves over the source baseline (B) in most settings, with particularly strong gains on MGSM across all three models. On MGSM, our method consistently outperforms the baseline, with improvements ranging from approximately 2.9% to 5.4%, indicating

that activation steering is especially effective for structured mathematical reasoning tasks.

On XNLI and MSVAMP, improvements are generally smaller and more variable. Our method outperforms the baseline in several configurations, though the multilingual few-shot baseline (MFS) remains competitive and occasionally achieves higher accuracy, especially on MSVAMP. Despite this variability, our method performs comparably to or better than the baselines in the majority of cases, demonstrating that activation steering provides a robust and model-agnostic approach for improving cross-lingual transfer across different model sizes and task types.

4. Ablation Studies and Analysis

4.1. Cross-Task Analysis

We analyze whether language-specific steering vectors transfer across tasks by conducting systematic evaluation across three datasets: MGSM, MSVAMP, and XNLI, examining all six possible transfer directions. For each direction, we compute steering vectors from the source task and apply them during evaluation on the target task, using only languages common to both datasets. Table 3 presents cross-task transfer results for Qwen-2.5-7b-Instruct.

Successful Transfers.

Five out of six transfer directions demonstrate effective generalization. A→B indicates B dataset was steered using vector calculated by A. MGSM→XNLI achieves 75.96%, while XNLI→MGSM shows transfer at 76.36%. Between math tasks, MGSM→MSVAMP reaches 75.18% and MSVAMP→MGSM achieves 72.22%. All the results are greater than the baseline. High-resource languages like Spanish consistently benefit across successful transfers (83–89% range), while lower-resource languages like Swahili show more variable but generally positive results.

Failed Transfer. MSVAMP→XNLI represents a significant failure case, achieving only 47.98% compared to 75.19% baseline, a 27 percentage point drop. This failure is asymmetric XNLI→MSVAMP works successfully suggesting MSVAMP vectors encode task-specific mathematical patterns detrimental to natural language inference. These results indicate that language-specific steering vectors generally capture task-agnostic cross-lingual representations, but transfer effectiveness depends on task compatibility.

4.2. Analysis on Intervention Methods

To understand which ablation works the best we see the accuracy values across all the methods. As seen in Table 4, our analysis reveals that intervention timing significantly impacts cross-lingual transfer performance. On MGSM with Llama 3.1-8b-instruct, the “On Few-shot” (OF) inter-

Table 3. Cross-task transfer results across all task-vector combinations between MGSM, MSVAMP, and XNLI on Qwen2.5-7b-instruct. For each task pair, we compute steering vectors from one task and evaluate on another task. B = Baseline, MFS = Multilingual Few-Shot Baseline (Tu et al., 2025), OR = Oracle. Oracle is the upper bound and is not for a direct comparison. **Bold** indicates the best score.

Language	MGSM (vector) → XNLI (eval)					XNLI (vector) → MGSM (eval)				
	B	MFS	Ours	CT	OR	B	MFS	Ours	CT	OR
Chinese	78.14	77.84	77.84	78.14	79.34	84.52	88.10	86.90	90.48	79.76
Thai	70.96	73.05	70.06	70.06	72.46	79.76	82.14	79.76	82.14	85.71
Swahili	53.29	47.31	55.09	54.79	54.79	19.05	16.67	26.19	25.00	33.33
Russian	78.14	79.04	79.64	79.94	78.74	85.71	85.71	88.10	88.10	85.71
French	81.44	82.34	82.34	82.63	82.63	78.57	80.95	79.76	79.76	84.52
Spanish	82.63	82.63	83.53	83.23	81.74	84.52	83.33	89.29	88.10	86.90
German	81.74	77.54	82.93	82.98	79.34	83.33	83.33	85.71	85.71	86.90
Average	75.19	74.22	75.92	75.96	75.58	73.64	74.32	76.53	76.36	77.55

Language	MGSM (vector) → MSVAMP (eval)					MSVAMP (vector) → MGSM (eval)				
	B	MFS	Ours	CT	OR	B	MFS	Ours	CT	OR
Bengali	65.57	70.96	67.07	67.66	76.35	58.33	66.67	58.33	58.33	65.48
German	79.04	79.94	80.24	79.64	81.14	83.33	83.33	85.71	77.38	86.90
Spanish	82.04	83.53	83.83	82.04	84.13	84.52	83.33	89.29	90.48	86.90
French	82.34	85.03	84.13	79.04	84.43	78.57	80.95	79.76	77.38	84.52
Japanese	78.74	79.64	79.64	79.34	79.94	71.43	75.00	75.00	76.19	73.81
Russian	79.64	82.93	81.74	80.54	81.44	85.71	85.71	88.10	83.33	85.71
Swahili	49.10	47.90	50.60	48.50	50.30	19.05	16.67	26.19	20.24	33.33
Thai	78.14	79.34	79.34	76.65	78.74	79.76	82.14	86.90	78.57	79.76
Chinese	82.04	83.23	83.53	83.23	82.34	84.52	88.10	86.90	88.10	79.76
Average	75.18	76.94	76.68	75.18	77.65	71.69	73.54	74.34	72.22	75.79

Language	MSVAMP (vector) → XNLI (eval)					XNLI (vector) → MSVAMP (eval)				
	B	MFS	Ours	CT	OR	B	MFS	Ours	CT	OR
Chinese	78.14	77.84	77.84	66.47	79.34	82.04	83.23	83.53	80.84	82.34
Thai	70.96	73.05	70.06	29.04	72.46	78.14	79.34	79.34	76.35	78.74
Swahili	53.29	47.31	55.09	33.83	54.79	49.10	47.90	50.60	48.80	50.30
Russian	78.14	79.04	79.64	43.70	78.74	79.64	82.93	81.74	79.64	81.44
French	81.44	82.34	82.34	50.60	82.63	82.34	85.03	84.13	81.44	84.43
Spanish	82.63	82.63	83.53	58.98	81.74	82.04	83.53	83.83	82.63	84.13
German	81.74	77.54	82.93	53.29	79.34	79.04	79.94	80.24	81.14	81.14
Average	75.19	74.22	75.92	47.98	75.58	75.18	76.94	76.68	75.83	77.65

vention achieves the highest average accuracy at 64.64%, roughly a 4 percent improvement over baseline (60.95%). This indicates that modifying activations after processing source language demonstrations but before the target language question most effectively enables cross-lingual reasoning transfer. The “After Few-shot” (AF) intervention, applied only at the demonstration-question boundary, achieves 63.57% achieves the least accuracy in all the configurations. These results show that optimal steering is both language-dependent and position-sensitive, with few-shot-targeted interventions providing the most reliable improvements.

4.3. Random Vector Baseline

To validate that our steering vectors capture meaningful language-specific information, we compare against random steering vectors generated from a standard normal distribu-

tion. Table 5 shows that random steering achieves surprisingly competitive performance, with only a small gap (0.94–1.29 percentage points) below our method across datasets and models. We hypothesize that the hyperparameter optimization process (searching over layers, alphas, and positions) adapts even random vectors to the target language through validation performance, effectively finding beneficial configurations regardless of initial structure. However, the consistent gap across all models and datasets suggests that explicitly computed language-specific vectors do capture additional meaningful structure beyond what random perturbations can achieve through optimization alone.

4.4. Language Analysis Through Vectors

In our setup, middle layers prove most effective for steering. A detailed analysis of which layer is most effective per

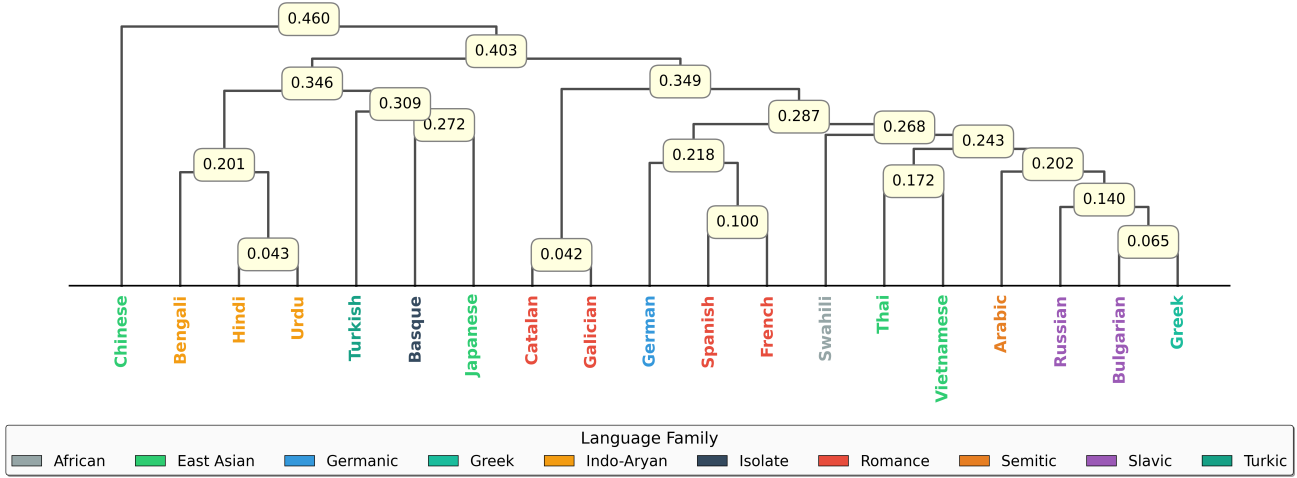


Figure 2. Hierarchical clustering based on cosine distance between their steering vectors at layer 10 of Llama-3.1-8B-Instruct. Edge labels show cosine distances (lower = more similar). Colors indicate language families. This structure suggests that the model’s internal representations encode meaningful language-specific patterns that align with both typological similarity and script characteristics.

Table 4. Per-language accuracy results across different intervention points for the Llama3.1-8b-instruct model on MGSM. Gray dashes indicate missing values. B = Baseline, OF = On Few-shot, AF = After Few-shot, OQ = On Question, ENT = Entire Prompt. **Bold** indicates the best score per language.

Language	B	OF	AF	OQ	ENT
Bengali	57.14	60.71	61.90	60.71	55.95
Catalan	64.29	66.67	69.05	67.86	69.05
Spanish	77.38	73.81	75.00	76.19	72.62
Basque	32.14	35.71	36.90	33.33	36.90
French	61.90	67.86	64.29	67.86	70.24
Galician	64.29	73.81	69.05	66.67	73.81
Russian	71.43	69.05	72.62	72.62	69.05
Swahili	55.95	65.48	60.71	63.10	65.48
Thai	57.14	61.90	54.76	61.90	59.52
Chinese	67.86	71.43	71.43	71.43	71.43
Average	60.95	64.64	63.57	64.16	64.40

language is shown in Table 6. We see that layer 10 is the most effective in a 32 layer model.

To understand relationships between language-specific steering vectors, we performed hierarchical clustering using cosine distance with average linkage as shown in Figure 2. The figure reveals patterns aligned with linguistic relationships.

Genetic Clustering. Closely related languages cluster tightly. Hindi-Urdu form the closest pair (distance ~ 0.043), expected given they share the same linguistic core despite different scripts. Romance languages group together with Catalan-Galician (distance ~ 0.042) merging first, followed by Spanish-French (~ 0.10) and Russian-Bulgarian (distance ~ 0.065).

Language Families. The Indo-Aryan cluster (Ben-

Table 5. Comparison of random steering vs. our method across datasets and models. R = Random Steering.

Dataset	Model	Performance	
		R	Ours
XNLI	LLaMA-3.1-8B-Instruct	63.37	64.31
	Qwen2.5-7B-Instruct	75.32	75.51
	Qwen2.5-14B-Instruct	73.37	74.66

gali, Hindi, Urdu) merges at distance ~ 0.201 . East Asian languages (Thai, Vietnamese) group together (distance ~ 0.172), with Chinese positioned nearby. Turkic and Basque-Japanese show an unexpected early merge (distance ~ 0.272 – 0.309), suggesting typologically distinct languages may require similar steering adjustments.

Notable Patterns. German emerges as an outlier among European languages, merging with other Germanic-Romance clusters only at higher distances (~ 0.287 – 0.343). Despite being Germanic like English (the source language), German requires distinctly different steering strategies, potentially reflecting its morphological complexity. Swahili and Arabic also show unique positioning, suggesting language-specific challenges in cross-lingual transfer. This structure demonstrates that steering vectors encode meaningful linguistic information, capturing regularities beyond surface-level orthography.

5. Related Work

In-Context Learning. In-context learning (ICL) has emerged as a fundamental capability of large language models, enabling them to adapt to new tasks given only a small number of demonstration examples in the input prompt,

Table 6. Best Performing Layer per Language for Llama3.1-8b-Instruct model.

Language	Best Layer	Language	Best Layer
Basque	5	Bengali	20
Bulgarian	10	Catalan	20
Chinese	20	French	5
Galician	5	German	5
Japanese	5	Russian	25
Spanish	30	Swahili	10
Thai	25	Arabic	10
Vietnamese	10	Urdu	10
Hindi	10	Turkish	10
Greek	10		

without any gradient updates. Early work formalized ICL as a paradigm where models perform predictions based purely on context, demonstrating strong few-shot and zero-shot performance across a wide range of tasks (Brown et al., 2020). Subsequent surveys systematized this line of research, categorizing prompting strategies, task types, and theoretical explanations for why ICL emerges in large-scale transformers (Dong et al., 2024).

While much of the early work on ICL focused on English or high-resource languages, more recent studies have begun to explore multilingual and cross-lingual ICL. These works show that when demonstrations are provided in a high-resource language such as English but test queries are in a low-resource language, performance often degrades substantially, highlighting limitations in direct cross-lingual transfer (Tanwar et al., 2023; Winata et al., 2021).

Cross-Lingual Transfer and Alignment. Cross-lingual transfer learning investigates how multilingual language models transfer knowledge across languages (Huang et al., 2021; 2022). Despite being trained on large multilingual corpora, these models frequently exhibit performance gaps in different languages. Recent surveys document persistent cross-lingual disparities and the challenges associated with robust transfer with multilingual models (Qin et al., 2025).

To mitigate these issues, several alignment strategies have been proposed. These include prompt-level techniques such as multilingual prompting and code-switching, as well as representation-level methods that aim to align semantic spaces across languages. Recent work explores cross-lingual alignment to improve in-context learning of multilingual generative models (Li et al., 2024; Zhang et al., 2024), while other approaches introduce cross-lingual in-context pretraining to explicitly encourage better language transfer during inference (Wu et al., 2025; Tu et al., 2025; Ahuja et al., 2023). Despite these advances, most alignment methods still require additional training, parallel data, or specialized prompting, limiting their practicality in low-resource settings.

Multilingual Representations and Interpretability. Understanding how multilingual language models internally represent different languages is an active area of research. Empirical analyses have shown that multilingual models encode language identity and linguistic structure in distinct regions of their latent space, which can significantly affect cross-lingual generalization (Zhao et al., 2024; Gurgurov et al., 2025). These findings suggest that language-specific information is not uniformly distributed, but instead manifests as structured patterns within model activations (Tang et al., 2024; Pokharel et al., 2026).

Recent surveys on multilingual interpretability synthesize evidence that language identity, syntax, and semantics can be disentangled to some extent within transformer representations, and that probing and representation analysis can reveal systematic cross-lingual behaviors (Resck et al., 2025; Joshi et al., 2025). These insights motivate the use of activation-level interventions as a means to directly manipulate language representations without retraining (Nie et al., 2025).

Activation steering and latent space manipulation have recently gained attention as training-free techniques for controlling language model behavior at inference time. These methods modify internal activations to induce desired behaviors such as style control, factual editing, or task adaptation (Turner et al., 2023). Prior work has shown that extracting and applying activation vectors from contextual differences can significantly influence model outputs without modifying model parameters (Liu et al., 2023). More recent work proposes lightweight and automated activation steering methods that operate entirely post-training, demonstrating that simple residual stream interventions can reliably shift model behavior across tasks (Cui & Chen, 2025; Stolfo et al., 2024; Rinsky et al., 2024).

6. Conclusion

We introduce language vectors, a training-free steering approach that improves multilingual in-context learning by leveraging activation differences between languages. Our comprehensive evaluation across three datasets, 19 languages, and three model families demonstrates consistent performance gains, particularly on reasoning tasks. Beyond practical improvements, our analysis reveals that these vectors encode meaningful linguistic structure: they cluster by language families, transfer successfully across tasks, and operate most effectively in middle transformer layers.

In the future, we would like to investigate why certain task transfers fail to better understand the boundary between task-specific and language-specific representations, and explore extending this approach to truly low-resource languages and other multilingual capabilities beyond in-context learning.

7. Impact Statement

This work introduces a training-free method to improve cross-lingual performance of large language models by extracting and applying language-specific steering vectors based on activation differences. By enhancing multilingual in-context learning, our approach aims to narrow performance gaps between high-resource and low-resource languages, potentially enabling more equitable access to advanced AI capabilities for speakers of underrepresented languages.

Improved cross-lingual reasoning and task performance in large language models could benefit a wide range of language technologies, including automated translation, multilingual question answering, educational tools, and information access for diverse populations. These advances may help democratize AI applications by reducing the dependency on large labeled datasets or expensive model fine-tuning, especially in settings where data resources are scarce.

At the same time, advances in multilingual AI raise broader societal considerations. As model performance improves across languages, these technologies may be integrated into high-stakes applications (e.g., legal or medical assistants) that require robust safeguards against hallucination, bias, or misinformation. There is also the risk that such tools could be misused for generating plausible but harmful content in multiple languages if appropriate controls are not considered. Future work should further explore ethical deployment practices, robustness evaluation across language and cultural contexts, and mitigation strategies to ensure responsible use of these technologies.

Overall, we believe this research contributes to both the technical understanding and practical advancement of multilingual AI, while acknowledging the importance of continued attention to safety, fairness, and societal impact as these models are deployed more widely.

References

- Ahuja, K., Diddee, H., Hada, R., Ochieng, M., Ramesh, K., Jain, P., Nambi, A., Ganu, T., Segal, S., Ahmed, M., Bali, K., and Sitaram, S. MEGA: Multilingual evaluation of generative AI. In Bouamor, H., Pino, J., and Bali, K. (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pp. 4232–4267, Singapore, December 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.emnlp-main.258. URL <https://aclanthology.org/2023.emnlp-main.258/>.
- Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.
- Chen, N., Zheng, Z., Wu, N., Gong, M., Zhang, D., and Li, J. Breaking language barriers in multilingual mathematical reasoning: Insights and observations. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pp. 7001–7016, 2024.
- Conneau, A., Rinott, R., Lample, G., Williams, A., Bowman, S., Schwenk, H., and Stoyanov, V. Xnli: Evaluating cross-lingual sentence representations. In *Proceedings of the 2018 conference on empirical methods in natural language processing*, pp. 2475–2485, 2018.
- Cui, S. and Chen, Z. Painless activation steering: An automated, lightweight approach for post-training large language models. *arXiv preprint arXiv:2509.22739*, 2025.
- Dong, Q., Li, L., Dai, D., Zheng, C., Ma, J., Li, R., Xia, H., Xu, J., Wu, Z., Chang, B., et al. A survey on in-context learning. In *Proceedings of the 2024 conference on empirical methods in natural language processing*, pp. 1107–1128, 2024.
- Elhage, N., Hume, T., Olsson, C., Schiefer, N., Henighan, T., Kravec, S., Hatfield-Dodds, Z., Lasenby, R., Drain, D., Chen, C., Grosse, R., McCandlish, S., Kaplan, J., Amodei, D., Wattenberg, M., and Olah, C. Toy models of superposition. *Transformer Circuits Thread*, 2022.
- Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur, A., Schelten, A., Vaughan, A., et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- Gurgurov, D., Trinley, K., Al Ghussin, Y., Bäumel, T., van Genabith, J., and Ostermann, S. 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*, pp. 2911–2937, 2025.
- Huang, K., Ahmad, W. U., Peng, N., and Chang, K. Improving zero-shot cross-lingual transfer learning via robust training. In Moens, M., Huang, X., Specia, L., and Yih, S. W. (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pp. 1684–1697. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.EMNLP-MAIN.126. URL <https://doi.org/10.18653/v1/2021.emnlp-main.126>.

- Huang, K., Hsu, I., Natarajan, P., Chang, K., and Peng, N. Multilingual generative language models for zero-shot cross-lingual event argument extraction. In Muresan, S., Nakov, P., and Villavicencio, A. (eds.), *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, ACL 2022, Dublin, Ireland, May 22-27, 2022, pp. 4633–4646. Association for Computational Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.317. URL <https://doi.org/10.18653/v1/2022.acl-long.317>.
- Joshi, S., Dittadi, A., Lachapelle, S., and Sridhar, D. Identifiable steering via sparse autoencoding of multi-concept shifts. *arXiv preprint arXiv:2502.12179*, 2025.
- Li, C., Wang, S., Zhang, J., and Zong, C. Improving in-context learning of multilingual generative language models with cross-lingual alignment. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pp. 8058–8076, 2024.
- Liu, S., Ye, H., Xing, L., and Zou, J. In-context vectors: Making in context learning more effective and controllable through latent space steering. *arXiv preprint arXiv:2311.06668*, 2023.
- Nie, E., Schmid, H., and Schütze, H. Mechanistic understanding and mitigation of language confusion in english-centric large language models. *arXiv preprint arXiv:2505.16538*, 2025.
- Pokharel, R., Agrawal, A., and Nagar, T. Cross-lingual activation steering for multilingual language models. *arXiv preprint arXiv:2601.16390*, 2026.
- Qin, L., Chen, Q., Zhou, Y., Chen, Z., Li, Y., Liao, L., Li, M., Che, W., and Yu, P. S. A survey of multilingual large language models. *Patterns*, 6(1), 2025.
- Rai, D., Zhou, Y., Feng, S., Saparov, A., and Yao, Z. A practical review of mechanistic interpretability for transformer-based language models. *arXiv preprint arXiv:2407.02646*, 2024.
- Resck, L., Augenstein, I., and Korhonen, A. Explainability and interpretability of multilingual large language models: A survey. In Christodoulopoulos, C., Chakraborty, T., Rose, C., and Peng, V. (eds.), *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pp. 20454–20486, Suzhou, China, November 2025. Association for Computational Linguistics. ISBN 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.1033. URL <https://aclanthology.org/2025.emnlp-main.1033/>.
- Rimsky, N., Gabrieli, N., Schulz, J., Tong, M., Hubinger, E., and Turner, A. Steering llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 15504–15522, 2024.
- Shi, F., Suzgun, M., Freitag, M., Wang, X., Srivats, S., Vosoughi, S., Chung, H. W., Tay, Y., Ruder, S., Zhou, D., et al. Language models are multilingual chain-of-thought reasoners. *arXiv preprint arXiv:2210.03057*, 2022.
- Stolfo, A., Balachandran, V., Yousefi, S., Horvitz, E., and Nushi, B. Improving instruction-following in language models through activation steering. *arXiv preprint arXiv:2410.12877*, 2024.
- Tang, T., Luo, W., Huang, H., Zhang, D., Wang, X., Zhao, W. X., Wei, F., and Wen, J.-R. 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)*, pp. 5701–5715, 2024.
- Tanwar, E., Dutta, S., Borthakur, M., and Chakraborty, T. Multilingual LLMs are better cross-lingual in-context learners with alignment. In Rogers, A., Boyd-Graber, J., and Okazaki, N. (eds.), *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 6292–6307, Toronto, Canada, July 2023. Association for Computational Linguistics. doi: 10.18653/v1/2023.acl-long.346. URL <https://aclanthology.org/2023.acl-long.346/>.
- Team, Q. et al. Qwen2 technical report. *arXiv preprint arXiv:2407.10671*, 2(3), 2024.
- Tu, Y., Xue, A., and Shi, F. Blessing of multilinguality: A systematic analysis of multilingual in-context learning. *arXiv preprint arXiv:2502.11364*, 2025.
- Turner, A. M., Thiergart, L., Leech, G., Udell, D., Vazquez, J. J., Mini, U., and MacDiarmid, M. Steering language models with activation engineering. *arXiv preprint arXiv:2308.10248*, 2023.
- Winata, G. I., Madotto, A., Lin, Z., Liu, R., Yosinski, J., and Fung, P. Language models are few-shot multilingual learners. *arXiv preprint arXiv:2109.07684*, 2021.
- Wu, L., Wei, H.-R., Lin, H., Li, T., Yang, B., Huang, F., and Lu, W. Enhancing LLM language adaptation through cross-lingual in-context pre-training. In Christodoulopoulos, C., Chakraborty, T., Rose, C., and Peng, V. (eds.), *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pp. 27152–27166, Suzhou, China, November 2025.

Association for Computational Linguistics. ISBN 979-8-89176-332-6. doi: 10.18653/v1/2025.emnlp-main.1380. URL <https://aclanthology.org/2025.emnlp-main.1380/>.

Zhang, Z., Lee, D.-H., Fang, Y., Yu, W., Jia, M., Jiang, M., and Barbieri, F. PLUG: Leveraging pivot language in cross-lingual instruction tuning. In Ku, L.-W., Martins, A., and Srikumar, V. (eds.), *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 7025–7046, Bangkok, Thailand, August 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.acl-long.379. URL <https://aclanthology.org/2024.acl-long.379/>.

Zhao, Y., Zhang, W., Chen, G., Kawaguchi, K., and Bing, L. How do large language models handle multilingualism? *Advances in Neural Information Processing Systems*, 37: 15296–15319, 2024.