
Latent Mechanisms of Code-Switching in Large Language Models

Anonymous Authors¹

Abstract

Multilingual large language models can exhibit *unintended code-switching* – unnecessarily alternating between languages during generation. We present a comparative study of three methods that identify language-controlling latents in cross-layer transcoders: activation value-based selection (**ValSel**), activation frequency-based selection (**FreqSel**), and LLM-generated latent annotation-based selection (**AnnSel**). To evaluate the efficacy of these methods in identifying language-controlling latents, we introduce two multilingual code-switching benchmarks designed for fine-grained analysis of language steering across seven languages. Through targeted intervention experiments on **Gemma-2-2B** and **Qwen3-4B**, we find that all three methods effectively manipulate generation language, with **FreqSel** achieving the strongest overall performance, while **AnnSel** offering interpretable latent selection through explicit language annotations. A knock-out analysis suggests the methods select non-overlapping but each-functional latent subsets, indicating redundancy rather than a single canonical language direction.

1. Introduction

Large language models have achieved remarkable multilingual capabilities, but this flexibility comes with an unintended consequence: *code-switching*, where models inappropriately alternate between languages mid-generation (Marchisio et al., 2024). This problem motivated explicit language consistency rewards during training of models such as DeepSeek-R1 (DeepSeek-AI et al., 2025).

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

Recent mechanistic interpretability work has identified neurons and latents that causally influence output language (Tang et al., 2024; Deng et al., 2025), suggesting that targeted interventions could provide fine-grained language control in generation. However, to our knowledge, no systematic comparison exists of language-controlling mechanism discovery, and it remains unclear whether different approaches uncover the same underlying latents or redundant representations.

We compare three approaches that identify language-specific latents in cross-layer transcoders (CLTs) (Ameisen et al., 2025): **ValSel**, based on activation value differences (Deng et al., 2025); **FreqSel**, based on language-specific activation frequency (Andrylie et al., 2025); and **AnnSel** that utilizes LLM-generated annotations to find latents whose descriptions explicitly reference the target language from an extracted circuit graph. Unlike prior work on single-layer SAEs, CLTs provide a unified sparse dictionary across all layers, enabling holistic analysis of language representations in latent space.

2. Preliminaries

Mechanistic Approaches to Language Control

Prior work has investigated language control at multiple levels of abstraction. At the neuron level, Tang et al. (2024) identified language neurons whose selective activation steers output language. At the representation level, Goncharov et al. (2025) extracted language directions via PCA, enabling steering with directional ablation. More recently, sparse autoencoders (SAEs), which address *polysemanticity* by learning an overcomplete sparse basis of monosemantic features (Bricken et al., 2023), have been applied to isolate language-controlling features through activation magnitude (Deng et al., 2025) or frequency patterns (Andrylie et al., 2025).

Cross-Layer Transcoders (CLTs) extend SAEs by replacing MLP blocks across all layers with a jointly-trained sparse dictionary (Ameisen et al., 2025), capturing features spanning model depth. In this paper, we use CLTs trained for **Gemma-2-2B** (Lieberum et al., 2024) and **Qwen3-4B** (Hanna, 2025).

3. Manipulating Generation Language through Latent Interventions

We study three language-specific **latent selection methods**, and four **latent intervention strategies** that steer generation toward a target language.

3.1. Latent Selection Methods

Given a cross-layer transcoder with latent space \mathcal{S} and a collection of languages L , we identify for language $l \in L$ a subset of latents $\mathcal{S}^l \subset \mathcal{S}$ that are specifically associated with the generation of l . We compare three selection criteria based on activation magnitude, activation frequency, and semantic annotation.

3.1.1. VALUE-BASED SELECTION (VALSEL)

Inspired by Deng et al. (2025), this method identifies latents whose activation magnitude is distinctively higher for a target language compared to others. For a latent s and a language-specific corpus D^l , we compute the mean activation:

$$\mu_s^l = \frac{1}{|D^l|} \sum_{t \in D^l} a_s(t) \quad (1)$$

where $a_s(t)$ denotes the activation of latent s for token t . The *l-monolinguality score* of latent s for language l is then defined as:

$$v_s^l = \mu_s^l - \frac{1}{|L| - 1} \sum_{j \in L \setminus \{l\}} \mu_s^j \quad (2)$$

We select the top- K latents with the highest *l-monolinguality scores* for each language l .

3.1.2. FREQUENCY-BASED SELECTION (FREQSEL)

Following Andrylie et al. (2025), this method identifies latents that activate more frequently for a specific language. We compute the *activation probability* for latent s in language l as:

$$P_s^l = \frac{1}{|D^l|} |\{t \in D^l : a_s(t) > 0\}| \quad (3)$$

A latent is considered *specific* to language l if: (1) $P_s^l = \max_{j \in L} P_s^j$ and (2) no language $l' \neq l$ satisfies $P_s^{l'} \geq T \cdot P_s^l$ for a threshold $T \in [0, 1]$. We select the top- K specific latents for l based on P_s^l .

3.1.3. ANNOTATION-BASED SELECTION (ANNSEL)

We propose an approach that supports interpretability by leveraging LLM-generated latent annotations. This method operates in three stages: circuit tracing, path pruning and annotation filtering.

For each sentence in D^l , we trace the attribution graph from input tokens to the predicted output token using Circuit Tracer (Hanna et al., 2025). We truncate sentences at a midpoint ensuring the target token is a valid word in l . Then, we greedily prune the attribution graph by retaining only its high-importance paths from each token embedding to the top predicted logit by maximum bottleneck edge weight. From the pruned graph, we extract all latents contributing to the final token position, but retain only the latents whose annotations explicitly reference the language l (e.g., containing ‘‘Japanese’’ or ‘‘Japan’’ for $l = \text{‘ja’}$), ensuring interpretable, language-related selection.

3.2. Intervention Strategies

For the sets of discovered language-specific latents $\{\mathcal{S}^l\}^{l \in L}$, we apply inference-time interventions to steer generation from a context language l_0 toward a target language l . We evaluate four strategies:

- **Target Amplification:** Set each latent $s \in \mathcal{S}^l$ to its mean nonzero activation computed over D^l . This naturally amplifies the signal for the target language l while avoiding out-of-distribution effects, without affecting other languages.
- **Distractor Zero Ablation:** Set to zero all latents associated with the distractor language l_0 , or the latents of all non-target languages $\bigcup_{j \in L \setminus \{l\}} \mathcal{S}^j$. This suppresses competing signals from other languages.
- **Distractor Direction Ablation:** Ablate the latents of the distractor language \mathcal{S}^{l_0} by *direction ablation*, where the the context language get suppressed by removing the projection of the residual stream activation along a latent’s direction. We do this either for the single layer with most latents (*one-layer*), or for all-layers (*multi-layer*).
- **Combined Interventions:** Simultaneously apply a form of distractor ablation along with target amplification, boosting the signal of the target language while suppressing that of the context language. We test distractor zero ablation + target amplification (**Zero+Amp**), and distractor one-layer direction ablation + target amplification (**1L+Amp**).

4. Experiments

All the experiments are conducted on the transcoders of Gemma-2-2B-pt¹ from Gemma-scope (Lieberum et al., 2024) and Qwen3-4b².

For the activation value-based ValSel we use $K = 50$ (to parallel the 2-latent per layer choice in (Deng

¹<https://hf.co/mwhanna/gemma-scope-transcoders>

²<https://hf.co/mwhanna/qwen3-4b-transcoders>

Table 1. Data format and examples of the two evaluation datasets. The sentence frame (in sans-serif) is in the context language l_0 , while the target language l is in **bold**. The examples show selections of context/target languages (l_0, l).

	Context	Target
Antonyms	The opposite of “[adj]” is “	[antonym]
ex1 (en, fr)	The opposite of “ <i>grand</i> ” is “	<i>petit</i>
ex2 (de, ja)	Das Gegenteil von “ 大きい ” ist “	小さい
Enumerations	The [enum] are: [item 1], [item 2],	[item 3], [item 4], ...
ex3 (en, es)	The four seasons are: <i>el verano, el otoño,</i>	<i>el invierno, la primavera</i>
ex4 (fr, zh)	Les jours de la semaine sont: 星期日, 星期一, 星期二,	星期三, 星期四, ...

et al., 2025) for the 26 layers of Gemma-2-2B). For the frequency-based selection FreqSel, we consider latents that are active in more than 98% of a language’s corpus and set $T = 0.8$ following (Andrylie et al., 2025). For AnnSel, we utilize the Circuit Tracer tool (Hanna et al., 2025) and Neuronpedia annotations³.

To extract the latent sets, we use FLORES+ (NLLB Team et al., 2024) in seven languages: English (en), French (fr), German (de), Spanish (es), Chinese (zh), Korean (ko), and Japanese (ja). We apply the three selection methods separately to both CLTs and each language subset. We report in the Appendix (Table 6) the count of latents extracted per method and language.

4.1. Datasets

We introduce **Antonyms** and **Enumerations**, two novel evaluation datasets designed to measure the effectiveness of language-steering interventions.⁴ **Antonyms** task requires producing a target antonym in a desired language l given a sentence context in a different language l_0 , creating a controlled code-switching scenario, while **Enumerations** asks for a continuation of a listing in the desired language, testing open-ended generation. The data format and some examples of the two datasets are shown in Table 1.

Dataset construction. We collect 100 antonym pairs per language for the seven languages (en, fr, de, es, zh, ja, ko) for **Antonyms** and 7 ordered categories for **Enumerations**. The complete lists of adjectives and words are provided in Appendix H and I. For each of the $7 \times 6 = 42$ ordered pairs of context and target languages (l_0, l) where $l_0 \neq l$, we generate a test example per antonym pair and per enumeration category.

Evaluation protocol. Given a context, we measure whether the model assigns a higher probability to the target antonym or enumeration in language l compared to its translations in other languages. For each example, we record the *logit* assigned to the target token of the

antonyms for **Antonyms** and the length-normalized *logprob* for **Enumerations** across all seven languages. A successful intervention should increase the target logit or logprob post-intervention relative to the baseline before the intervention.

4.2. Findings

4.2.1. INTERVENTION STRATEGY EFFECTIVENESS

To assess language manipulation efficacy, we intervene on the extracted latents enforcing a target language l in a sentence with context language l_0 using **Antonyms**, where the target is a single token. We measure the target token’s *top-1 logit margin* which is the separation of the logit of the target adjective from its best competitor: $m_l = \max_{i \in L \setminus \{l\}} a_i - a_l$, where a_l is the logit of the target adjective in language l . A more negative m_l value indicates a higher probability for target language l . We measure the change in m_l due to the intervention:

$$\Delta m_l = -(m_l^{after} - m_l^{before}) \quad (4)$$

hence a larger positive change indicates a more effective intervention. Summary results for Gemma-2-2B are in Table 2, with detailed results in the Appendix Tables 7 to 9, as well as results for Qwen3-4B (Tables 10 to 12).

The most effective intervention is **Zero+Amp**, combining zero ablation of distractor language l_0 with target language amplification, followed by **1L+Amp** that instead uses single-layer directional ablation of l_0 . Zero ablation seems to outperform the more powerful direction ablation possibly due to latent overlap of related languages, which weakens the target language when direction-ablating the context language.

4.2.2. SELECTION METHOD EFFICACY

We now apply the most effective intervention (**Zero+Amp**) to all selection methods and languages to contrast their performance per target language. Table 3 shows summary results for Gemma-2-2B on **Antonyms** and **Enumerations** (detailed results in Tables 13 to 15 in the Appendix), and results for Qwen3-4B are in Table 16.

³<https://www.neuronpedia.org>

⁴The datasets will be publicly released.

Table 2. Intervention effectiveness per method using **Antonyms**, computed as the total change in top-1 logit margin under the intervention (Equation 4) averaged over all target languages (bigger is better) for **Gemma-2-2B**.

	Distractor (l_0) zero ablation	$L \setminus \{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
ValSel	-8.02	22.39	1.99	13.98	69.57	79.83	<u>69.73</u>
FreqSel	2.02	21.39	-2.95	-51.62	90.83	99.07	<u>96.16</u>
AnnSel	-0.20	11.90	0.30	64.05	87.00	96.72	<u>89.33</u>

Table 3. Selection method efficacy per language under the **Zero+Amp** intervention for **Gemma-2-2B**. *Top*: post-intervention total change in top-1 logit margin for **Antonyms**; *Bottom*: total change in normalized logProb of the target sequence for **Enumerations**.

		es	en	zh	de	ja	fr	ko	Avg
Antons.	ValSel	14.19	10.16	16.06	<u>13.56</u>	4.63	<u>9.01</u>	<u>12.18</u>	11.40
	FreqSel	8.95	4.34	34.50	7.85	<u>8.65</u>	6.45	28.30	14.15
	AnnSel	<u>13.86</u>	<u>4.87</u>	<u>24.52</u>	19.67	11.12	13.54	9.11	<u>13.81</u>
Enums.	ValSel	<u>1.00</u>	<u>-0.08</u>	<u>1.13</u>	0.93	<u>1.54</u>	0.72	<u>1.87</u>	<u>1.02</u>
	FreqSel	2.08	-0.33	2.12	2.69	2.52	1.40	3.92	2.06
	AnnSel	0.90	0.30	1.05	<u>1.83</u>	1.47	<u>0.85</u>	0.38	0.99

In total across all settings and models, all the selection methods are effective. **FreqSel** achieves the highest total absolute logit and logProb change across the tasks, with significant performance gains in Asian languages (especially Chinese and Korean). **AnnSel**, on the other hand, shows decent performance across all languages but is hindered by its weak performance in Korean which might be justified by the small number of extracted latents for Korean (Table 6).

4.2.3. IS THERE REDUNDANCY IN LATENT LANGUAGE REPRESENTATION?

To study to what extent the latent sets identified by the three selection methods are equivalent we conduct a knock-out experiment that involves a pair of methods and the same target language. We utilize the amplification intervention on the latents of the first method, then we ablate the projection of the residual stream change vector on the representative direction of the latents of the second method. We aim by this experiment to remove the influence of the second methods' latents completely to see if the first method latents are still able to produce positive logit margin change.

We report in Table 4 the total logit margin change aggregated over the seven language per method pair for **Gemma-2-2B** and **Qwen3-4B** on **Antonyms**. The generally positive numbers indicate a form of robustness, however the asymmetry of the results for the same pair of methods require further investigation. We provide in the Appendix (Section J) additional observations related to the cosine similarity of the latent set rep-

resentative residual-stream directions, and to latent annotations.

Table 4. Aggregate knock-out experiment results over all languages for **Gemma-2-B** and **Qwen3-4B**, where the latents of one selection method for the target language are amplified and the latents of the second method are used for direction ablation.

	Amplification method	Ablation method	Total Δm_l
Gemma-2-2B	AnnSel	FreqSel	5.20
	AnnSel	ValSel	38.82
	FreqSel	AnnSel	<u>90.63</u>
	FreqSel	ValSel	92.16
	ValSel	AnnSel	49.52
	ValSel	FreqSel	-1.19
Qwen3-4B	AnnSel	FreqSel	1.28
	AnnSel	ValSel	7.58
	FreqSel	AnnSel	155.29
	FreqSel	ValSel	85.82
	ValSel	AnnSel	<u>107.55</u>
	ValSel	FreqSel	33.86

5. Conclusion

We presented a comparative study of three methods for identifying language-specific latents in cross-layer transcoders: value-based (**ValSel**), frequency-based (**FreqSel**), and annotation-based (**AnnSel**). Through experiments on the controlled **Antonyms** and **Enumerations** tasks in seven languages, we demonstrated that all three methods can effectively steer generation language. **FreqSel** achieves the strongest overall performance, particularly for Asian languages, while **AnnSel** provides the added benefit of interpretable latent selection.

Future work could look into generalizing the findings of this work to larger models and to more languages. Real-world code-switching scenarios may present different challenges to the studied controlled settings. Deeper understanding of to what extent different methods identify similar or different latents has implications for both understanding multilingual representations and developing robust language-steering interventions.

References

- Ameisen, E., Lindsey, J., Pearce, A., Gurnee, W., Turner, N. L., Chen, B., Citro, C., Abrahams, D., Carter, S., Hosmer, B., Marcus, J., Sklar, M., Templeton, A., Bricken, T., McDougall, C., Cunningham, H., Henighan, T., Jermyn, A., Jones, A., Persic, A., Qi, Z., Ben Thompson, T., Zimmerman, S., Rivoire, K., Conerly, T., Olah, C., and Batson, J. Circuit tracing: Revealing computational graphs in language models. *Transformer Circuits Thread*, 2025. URL <https://transformer-circuits.pub/2025/attribution-graphs/methods.html>.
- Andrylie, L. M., Rahmanisa, I., Ihsani, M. K., Wicaksono, A. F., Wibowo, H. A., and Aji, A. F. Sparse autoencoders can capture language-specific concepts across diverse languages, 2025. URL <https://arxiv.org/abs/2507.11230>.
- Bricken, T., Templeton, A., Batson, J., Chen, B., Jermyn, A., Conerly, T., Turner, N., Anil, C., Denison, C., Askell, A., Lasenby, R., Wu, Y., Kravec, S., Schiefer, N., Maxwell, T., Joseph, N., Hatfield-Dodds, Z., Tamkin, A., Nguyen, K., McLean, B., Burke, J. E., Hume, T., Carter, S., Henighan, T., and Olah, C. Towards monosemanticity: Decomposing language models with dictionary learning. *Transformer Circuits Thread*, 2023. <https://transformer-circuits.pub/2023/monosemantic-features/index.html>.
- DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang, P., Bi, X., Zhang, X., Yu, X., Wu, Y., Wu, Z. F., Gou, Z., Shao, Z., Li, Z., Gao, Z., Liu, A., Xue, B., Wang, B., Wu, B., Feng, B., Lu, C., Zhao, C., Deng, C., Zhang, C., Ruan, C., Dai, D., Chen, D., Ji, D., Li, E., Lin, F., Dai, F., Luo, F., Hao, G., Chen, G., Li, G., Zhang, H., Bao, H., Xu, H., Wang, H., Ding, H., Xin, H., Gao, H., Qu, H., Li, H., Guo, J., Li, J., Wang, J., Chen, J., Yuan, J., Qiu, J., Li, J., Cai, J. L., Ni, J., Liang, J., Chen, J., Dong, K., Hu, K., Gao, K., Guan, K., Huang, K., Yu, K., Wang, L., Zhang, L., Zhao, L., Wang, L., Zhang, L., Xu, L., Xia, L., Zhang, M., Zhang, M., Tang, M., Li, M., Wang, M., Li, M., Tian, N., Huang, P., Zhang, P., Wang, Q., Chen, Q., Du, Q., Ge, R., Zhang, R., Pan, R., Wang, R., Chen, R. J., Jin, R. L., Chen, R., Lu, S., Zhou, S., Chen, S., Ye, S., Wang, S., Yu, S., Zhou, S., Pan, S., Li, S. S., Zhou, S., Wu, S., Ye, S., Yun, T., Pei, T., Sun, T., Wang, T., Zeng, W., Zhao, W., Liu, W., Liang, W., Gao, W., Yu, W., Zhang, W., Xiao, W. L., An, W., Liu, X., Wang, X., Chen, X., Nie, X., Cheng, X., Liu, X., Xie, X., Liu, X., Yang, X., Li, X., Su, X., Lin, X., Li, X. Q., Jin, X., Shen, X., Chen, X., Sun, X., Wang, X., Song, X., Zhou, X., Wang, X., Shan, X., Li, Y. K., Wang, Y. Q., Wei, Y. X., Zhang, Y., Xu, Y., Li, Y., Zhao, Y., Sun, Y., Wang, Y., Yu, Y., Zhang, Y., Shi, Y., Xiong, Y., He, Y., Piao, Y., Wang, Y., Tan, Y., Ma, Y., Liu, Y., Guo, Y., Ou, Y., Wang, Y., Gong, Y., Zou, Y., He, Y., Xiong, Y., Luo, Y., You, Y., Liu, Y., Zhou, Y., Zhu, Y. X., Xu, Y., Huang, Y., Li, Y., Zheng, Y., Zhu, Y., Ma, Y., Tang, Y., Zha, Y., Yan, Y., Ren, Z. Z., Ren, Z., Sha, Z., Fu, Z., Xu, Z., Xie, Z., Zhang, Z., Hao, Z., Ma, Z., Yan, Z., Wu, Z., Gu, Z., Zhu, Z., Liu, Z., Li, Z., Xie, Z., Song, Z., Pan, Z., Huang, Z., Xu, Z., Zhang, Z., and Zhang, Z. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025. URL <https://arxiv.org/abs/2501.12948>.
- Deng, B., Wan, Y., Yang, B., Zhang, Y., and Feng, F. Unveiling language-specific features in large language models via sparse autoencoders. In Che, W., Nabende, J., Shutova, E., and Pilehvar, M. T. (eds.), *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 4563–4608, Vienna, Austria, July 2025. Association for Computational Linguistics. ISBN 979-8-89176-251-0. doi: 10.18653/v1/2025.acl-long.229. URL <https://aclanthology.org/2025.acl-long.229/>.
- Fiotto-Kaufman, J., Loftus, A. R., Todd, E., Brinkmann, J., Pal, K., Troitskii, D., Ripa, M., Belfki, A., Rager, C., Juang, C., Mueller, A., Marks, S., Sharma, A. S., Lucchetti, F., Prakash, N., Brodley, C., Guha, A., Bell, J., Wallace, B. C., and Bau, D. Nnsight and ndif: Democratizing access to open-weight foundation model internals, 2025. URL <https://arxiv.org/abs/2407.14561>.
- Goncharov, A., Kondusov, N., and Zaytsev, A. Language steering in latent space to mitigate unintended code-switching, 2025. URL <https://arxiv.org/abs/2510.13849>.
- Hanna, M. Qwen3-4b transcoders, 2025. URL <https://huggingface.co/mwhanna/qwen3-4b-transcoders>. Huggingface.
- Hanna, M., Piotrowski, M., Lindsey, J., and Ameisen, E. circuit-tracer. <https://github.com/safety-research/circuit-tracer>, 2025. The first two authors contributed equally and are listed alphabetically.
- Lieberum, T., Rajamanoharan, S., Conmy, A., Smith, L., Sonnerat, N., Varma, V., Kramar, J., Dragan, A., Shah, R., and Nanda, N. Gemma scope: Open

- 275 sparse autoencoders everywhere all at once on gemma
276 2. In Belinkov, Y., Kim, N., Jumelet, J., Mo-
277 hebbi, H., Mueller, A., and Chen, H. (eds.), *Pro-*
278 *ceedings of the 7th BlackboxNLP Workshop: Ana-*
279 *lyzing and Interpreting Neural Networks for NLP*,
280 pp. 278–300, Miami, Florida, US, November 2024.
281 Association for Computational Linguistics. doi:
282 10.18653/v1/2024.blackboxnlp-1.19. URL [https://](https://aclanthology.org/2024.blackboxnlp-1.19/)
283 aclanthology.org/2024.blackboxnlp-1.19/.
284
- 285 Marchisio, K., Ko, W.-Y., Bérard, A., Dehaze, T., and
286 Ruder, S. Understanding and mitigating language
287 confusion in llms. *arXiv preprint arXiv:2406.20052*,
288 2024.
- 289 NLLB Team, Costa-jussà, M. R., Cross, J., Çelebi, O.,
290 Elbayad, M., Heafield, K., Heffernan, K., Kalbassi,
291 E., Lam, J., Licht, D., Maillard, J., Sun, A., Wang,
292 S., Wenzek, G., Youngblood, A., Akula, B., Bar-
293 rault, L., Gonzalez, G. M., Hansanti, P., Hoffman,
294 J., Jarrett, S., Sadagopan, K. R., Rowe, D., Spruit,
295 S., Tran, C., Andrews, P., Ayan, N. F., Bhosale,
296 S., Edunov, S., Fan, A., Gao, C., Goswami, V.,
297 Guzmán, F., Koehn, P., Mourachko, A., Ropers,
298 C., Saleem, S., Schwenk, H., and Wang, J. Scal-
299 ing neural machine translation to 200 languages.
300 *Nature*, 630(8018):841–846, 2024. ISSN 1476-4687.
301 doi: 10.1038/s41586-024-07335-x. URL [https://](https://doi.org/10.1038/s41586-024-07335-x)
302 doi.org/10.1038/s41586-024-07335-x.
303
- 304 Tang, T., Luo, W., Huang, H., Zhang, D., Wang, X.,
305 Zhao, W. X., Wei, F., and Wen, J.-R. Language-
306 specific neurons: The key to multilingual capabilities
307 in large language models. In *Proceedings of the 62nd*
308 *Annual Meeting of the Association for Computational*
309 *Linguistics (Volume 1: Long Papers)*, pp. 5701–5715,
310 2024.
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329

330 A. AnnSel implementation details

331 The implementation of AnnSel can be decomposed into three steps: attribution graph construction, path pruning,
332 and annotation-based latent selection.

334 A.1. Attribution graph construction

335 Since attribution graphs explain how the model predicts the next token, given a sentence from Flores+, we
336 randomly cut the sentence in the middle so that the next token that the model should predict is a valid word in
337 the target language, i.e. not a number or punctuation.

338 We used implementation from Circuit Tracer (Hanna et al., 2025) with parameters: `max_n_logits =`
339 `5`, `desired_logit_prob = 0.95`, `max_feature_nodes = None`. It means the attribution graphs are computed
340 to explain the top n logits, where $n = \min(5, m)$ and m is the smallest integer that satisfies $\sum_{i=1}^m prob_i \geq 0.95$.
341 We then pruned graph with `node_threshold = 0.8`, `edge_threshold = 0.98`, meaning keeping nodes that explain
342 80% of total influence on the output logits and edges that explain 98% of total influence on the output logits.

345 A.2. Path Pruning

346 After computing the attribution graph, we pruned the graphs to prevent unimportant features from being included
347 in the language supernodes. From each token to the top logit, we select important paths as follows. We first order
348 the edges by the edge weight in descending order and keep choosing edges that are incident to already chosen
349 nodes, starting from the token. When we reach the top logit, we can form a path from token to logit. Then, we
350 remove that chosen path and start again. We keep doing this until either the edge weight chosen is less than or
351 equal to 0.1 or we choose 75 distinct paths. We do this for every token position and obtain a set of paths.

352 After obtaining important paths, we further choose more important paths by setting threshold on the weights
353 of the first edge and last edge since they represent the importance of the chosen paths in terms of the token
354 and logit respectively and thus directly involved in the information flow from token to logit. We choose paths
355 whose first edge has more than 0.5 edge weight and last edge has more than 0.25 edge weight. Then, out of the
356 chosen features, we extracted all the features in the last token position and counted how many times each of
357 them appear.

360 A.3. Annotation-based latent selection

361 After obtaining the set of features and the number of times they appear in the important paths, we first extract
362 the description of each feature from Neuronpedia. We then remove all the features that does not mention the
363 language name or highly relevant country name, e.g. ["Japanese", "japanese", "Japan", "japan"] for Japanese.

364 Then, what we have is a list of features that includes language name or country name in its description and their
365 frequency. Let us denote the list by ℓ and the frequency of feature f by $freq_f$. We then choose f that satisfies
366 $freq_f \geq 0.1 \max_{k \in \ell} freq_k$.

370 B. Applying SAE feature selection methods on CLT

371 We applied two methods, ValSel and FreqSel, from (Deng et al., 2025) and (Andrylie et al., 2025) respectively
372 on the CLT. Since the original paper executes these methods on SAE, in order to justify the effectiveness of the
373 methods on CLT, we conducted two experiments proposed in (Deng et al., 2025) and (Andrylie et al., 2025).

376 B.1. ValSel and code-switching experiment

377 (Deng et al., 2025) confirms the specificity of the language specific features by comparing the activation values on
378 normal context and code-switching context. Following their experiments, we used the code-switching dataset of
379 theirs and plotted the activation values of the top-50 features found and used for steering in our experiments.
380 Keep in mind that we plotted the activations for all the languages used in our experiment except German because
381 (Deng et al., 2025) did not include German in the experiments. Also, since we conducted the experiments on the
382 CLT, the extracted features are not distributed uniformly across layers. The trend is the same as what is claimed
383

in (Deng et al., 2025) and the activation is generally highest in the original context and noun (Lang A Prefix + Lang A Noun), somewhat active in original context and modified noun (Lang A Prefix + Lang B Noun), and close to zero in modified noun (Lang B Noun).

B.2. FreqSel and text-generation experiment

(Andrylie et al., 2025) confirms the downstream effectiveness of the features by steering and doing unconditional text-generation. Their steering is $\mathbf{x}_{new} = \mathbf{x} + \alpha z_j^{max} \mathbf{d}^j$, where $\alpha \in \mathbb{R}$ is a scaling factor and z_j^{max} is the maximum value of the feature of interest \mathbf{d}^j in the multilingual corpora. We used the same multilingual corpora used to determine the language specific features when deciding on z_j^{max} . The model was prompted with [BOS] and decoded with top-p sampling. Some examples are shown in Figure 5. The output is steered toward the target languages although the sentences might not be perfectly natural.

Table 5. Examples of unconditional text generation given [BOS] ($p = 0.0$ means greedy decoding)

Intervened Language	α	p	Output
de	0.8	0.9	package Produktion Kleidung refroidissement nôtre enfans hindurchholen nôtre bissunter Bedürf Flüssigkeit geben Flüssigkeit Kleidung
en	0.1	0.0	How to get the value of a variable in a function in python\n\n
es	0.8	0.9	How convert Activity anuales in española y y argentina\n pre acuerdos efectivos
fr	0.4	0.8	Que efficaces financières russes russes efficaces financières financières efficaces financières efficaces financières efficaces
ja	0.3	0.0	How toget指定したフォルダ内のファイル一覧をスクスクスクスク
ko	0.5	0.0	How to 3D 렌샷 것이다 렌샷 것이다 렌
zh	0.5	0.9	The特平在道路上有什用呢？什不把柏的

C. Latent counts

Table 6 shows the number of extracted token per method and language pair for Gemma-2-2B and Qwen3-4B. Note that ValSel extracts 50 latents by design, while the other two methods return a variable number of latents that relates to the CLT training data.

Table 6. Count of extracted latents per method and language for Gemma-2-2B and Qwen3-4B.

		en	fr	de	es	zh	ko	ja
Gemma	ValSel	50	50	50	50	50	50	50
	FreqSel	101	423	579	509	268	938	234
	AnnSel	17	41	59	19	41	7	45
Qwen	ValSel	50	50	50	50	50	50	50
	FreqSel	11	438	500	299	136	366	293
	AnnSel	12	16	13	10	16	8	13

D. Ablation intervention details

We tried several types of ablations. The one mentioned earlier is direction ablation.

D.1. Distractor zero ablation

Using Circuit Tracer (Hanna et al., 2025), we find the language supernodes as described above. In context zero ablation, we ablate the distractor language A by setting all the latents in the language A supernode to zero. In non-target zero ablation, we ablate all the languages except the target language B by setting all the latents in the language C supernode to zero, where C is all the languages except B.

440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494

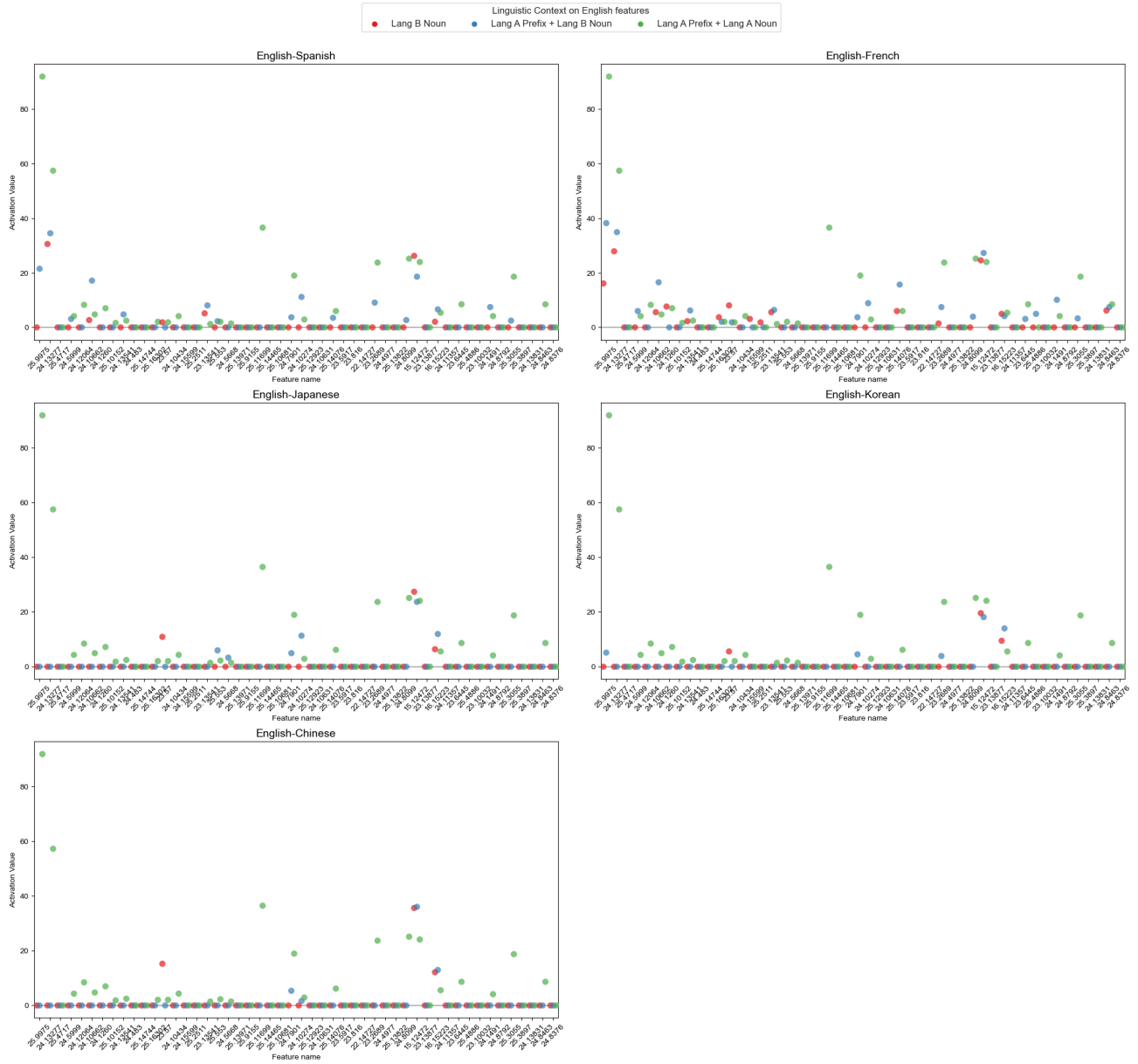


Figure 1. Activation values of English features on other language context and nouns

495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549



Figure 2. Activation values of French features on other language context and nouns

550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604



Figure 3. Activation values of Spanish features on other language context and nouns

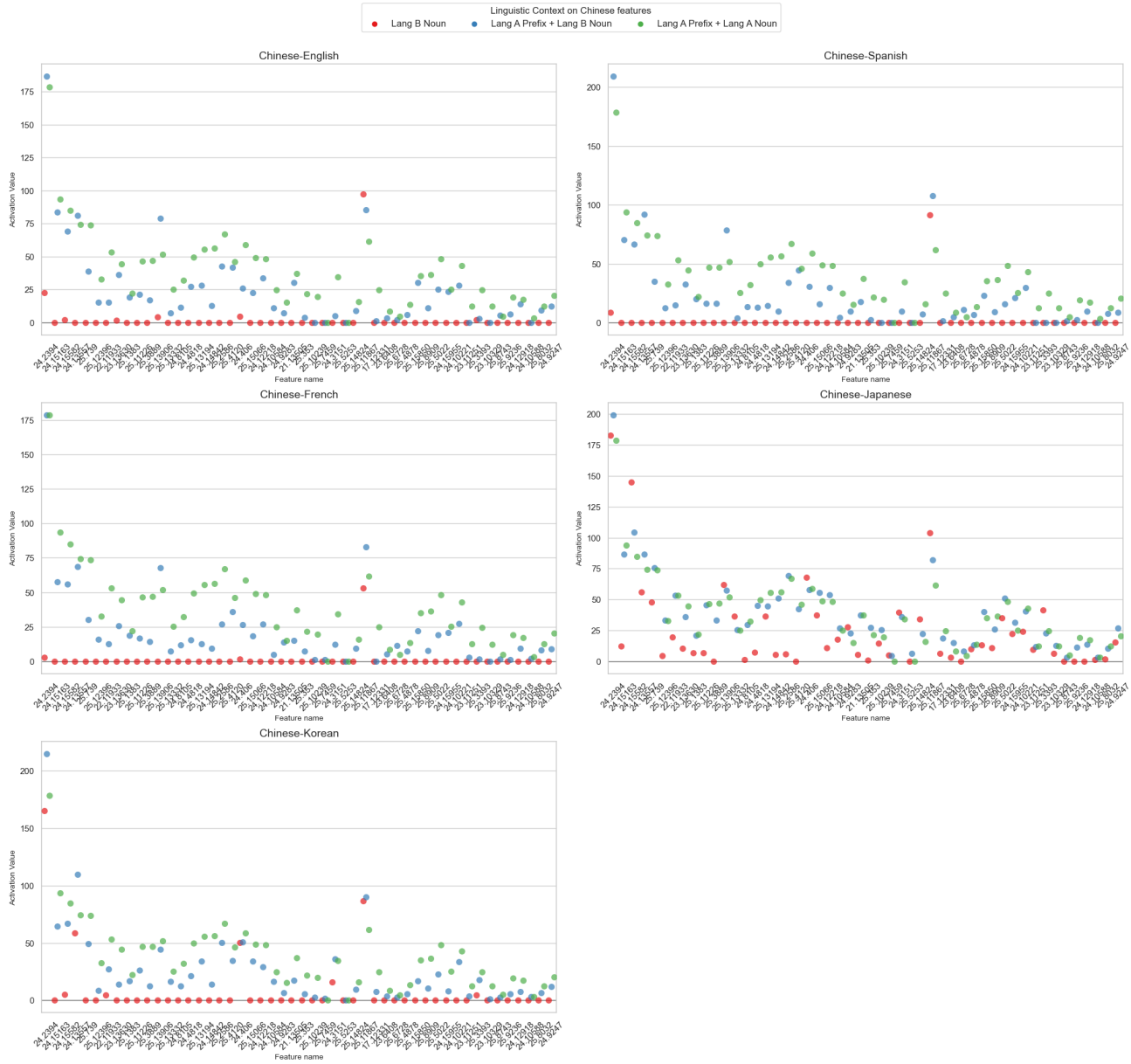


Figure 4. Activation values of Chinese features on other language context and nouns

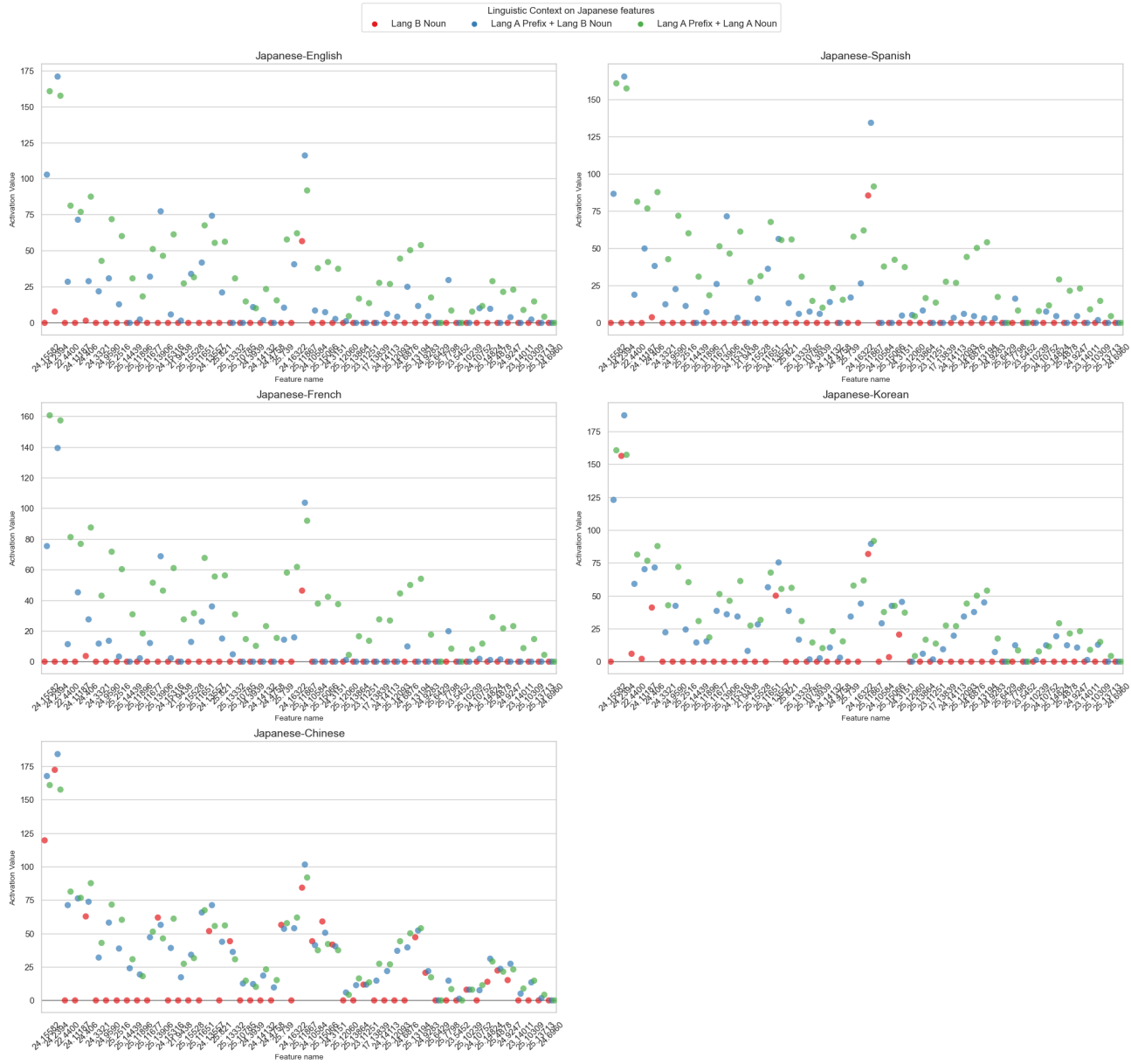


Figure 5. Activation values of Japanese features on other language context and nouns

715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769



Figure 6. Activation values of Korean features on other language context and nouns

D.2. One-layer direction ablation

The language features identified span across several layers, mostly in the last couple layers. Therefore, we chose the layer with the largest number of features for each language node with tie-breaking by the later layer, and did a direction ablation on that layer. Specifically, for each prompt, we run the attribution process to extract the activation values for the features in that layer and calculate the negative values we should insert into the features in the language node to completely get rid of the direction of the target features. Distractor one-layer direction ablation negates only one language, and one-layer direction ablation negates all the languages except one language.

The limited effectiveness of this approach might be due to the skip connections in CLTs and getting rid of the effects of one layer does not prevent other layers from influencing the logits through skip connections.

D.3. Multi-layer direction ablation

Instead of just getting rid of influence in one layer, we get rid of influence in all the layers. For all the features in a language supernode, we get the decoding weight matrices from CLT. Then, for each layer, we identify which direction to ablate using the weight matrices extracted based on which language features are in that layer. We finally do ablation in the activation space by extracting the projection onto that direction from the original activation. Since circuit tracer is a tool to manipulate the models in latent space, we used NNSight (Fiotto-Kaufman et al., 2025) instead to do ablation in activation space. Distractor (multi-layer) direction ablation intervenes on only one language, and (multi-layer) direction ablation intervenes on all the languages except one language.

Direction ablation works when the prompt language is either Chinese, Korean, or Japanese and when the adjective language is either German or Spanish. It does not work, for example, when the adjective language is French. When the prompt language is Chinese, Korean, or Japanese with French as the adjective language, the direction ablation pushes English and Spanish higher than French while successfully pushes the prompt languages down. This could suggest that a large part of French direction overlaps with Spanish and English directions.

Note that because multi-layer direction ablation operates in activation space and is implemented in NNSight (Fiotto-Kaufman et al., 2025), it is hard to combine it with non-distractor amplification, which operates in latent space and is implemented in circuit tracer (Hanna et al., 2025). Hence, we do not have multi-layer direction intervention.

E. Zero+Amp Detailed Results

Tables 13 to 15 show the logit difference results of ValSel, FreqSel and AnnSel, respectively, before and after the **Zero+Amp** intervention. The numbers in the table are for *top-1 logit margin* to the target adjective, which is the separation in the logit space of the target token to its best competitor. A value ≤ 0 (green) indicates that the target adjective is the most likely token, while a positive value (red) means it lags behind other undesired tokens.

Specifically, rows show the context language and columns show the adjective language. For each language pair, the ‘before’ value shows the logit margin before any intervention, while ‘after’ is its value after the latents amplification is applied. The desired outcome is that the value becomes smaller (greener) after intervention. Red cells indicate failures.

While the three methods show reasonable efficacy in enforcing Spanish and Chinese, FreqSel is weaker in enforcing German and English, and ValSel shows fragility when the context language is Chinese (Table 13). The situation of French adjectives is unique where ValSel and FreqSel fail for some context languages while AnnSel shows better robustness. On the other hand, AnnSel shows fragility in Korean (where only 7 latents were identified (Table 6) — far fewer than for any other language, while FreqSel shows the most robust performance.

F. Detailed Intervention Strategy Comparison

Next, we provide detailed results for the different interventions under Gemma-2-2B and the **Antonyms** task in Tables 7 to 9.

Table 7. Intervention effectiveness per target language for ValSel latents of Gemma-2-2B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	-3.18	3.12	-0.46	-3.15	14.52	13.6	13.92
en	7.49	5.06	2.77	52.07	5.88	10.16	4.24
es	0.25	3.06	-0.06	-1.57	13.35	14.19	15.08
fr	-0.77	1.15	-0.28	-9.65	8.74	9.01	8.71
ja	-4.61	-1.57	-0.41	-11.59	3.11	4.62	5.09
ko	-0.08	12.13	-1.61	-9.4	10.38	12.19	7.58
zh	-7.12	-0.56	2.04	-2.73	13.59	16.06	15.11
Total	-8.02	22.39	1.99	13.98	69.57	79.83	<u>69.73</u>

Table 8. Intervention effectiveness per target language for FreqSel latents of Gemma-2-2B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	-3.45	0.8	-1.67	0.58	5.08	7.86	6.29
en	5.86	5.43	3.6	-0.21	-0.55	4.36	1.94
es	-0.7	3.39	-0.6	-5.12	6.68	8.95	7.09
fr	-1.51	0.65	-0.78	-5.33	4.45	6.45	5.09
ja	-0.97	1.38	-1.45	-27.91	8.65	8.65	9.49
ko	4.82	6.62	-3.3	0.06	29.54	28.31	27.6
zh	-2.03	3.12	1.25	-13.69	36.98	34.49	38.66
Total	2.02	21.39	-2.95	-51.62	90.83	99.07	<u>96.16</u>

Table 9. Intervention effectiveness per target language for AnnSel latents of Gemma-2-2B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	-3.87	0.68	-0.04	-3.62	20.15	19.68	20.6
en	4.35	4.83	0.82	43.07	1.12	4.88	1.74
es	-0.04	1.98	-0.15	6.6	13.91	13.86	14.05
fr	-1.01	0.62	0.02	-3.71	13.69	13.54	13.95
ja	-2.21	1.5	-1.01	-4.97	11.17	11.12	10.53
ko	6.76	5.97	0.44	12.19	4.27	9.13	4.81
zh	-4.18	-3.68	0.22	14.49	22.69	24.51	23.65
Total	-0.20	11.90	0.30	64.05	87.00	96.72	<u>89.33</u>

G. Detailed Selection Method Comparison

We show detailed results for Gemma-2-2B under the three selection methods and the **Zero+Amp** intervention with **Antonyms** in Tables 13 to 15.

Table 16 presents summary results for selection method efficacy for Qwen3-4B latents for **Antonyms** and **Enumerations**.

Table 10. Intervention effectiveness per target language for ValSel latents of Qwen3-4B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	5.24	0.83	-4.32	-21.71	4.58	8.47	4.99
en	1.05	6.6	-2.81	2.16	-14.3	-6.17	-14.71
es	3.35	1.7	-1.4	-20.7	2.79	7.4	4.44
fr	3.51	-0.23	-2.79	-16.11	2.85	7.51	4.05
ja	3.33	-4.71	-6.52	-19.02	-6.02	-1.2	-4.12
ko	18.59	2.5	-2.11	-5.64	3.44	11.81	7.29
zh	5.47	8.84	2.67	17.1	7.51	17.74	13.9
Total	<u>40.54</u>	15.53	-17.28	-63.92	0.85	45.56	15.84

Table 11. Intervention effectiveness per target language for FreqSel latents of Qwen3-4B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	-2.92	0.15	-0.28	-24.44	4.69	6.35	5.08
en	6.67	8.39	1.19	-9.42	-2.01	6.11	-1.55
es	1.41	3.67	0.35	-21.37	2.85	4.92	2.93
fr	-0.96	2.51	-0.3	-16.97	5.05	6.88	5.18
ja	-3.3	0.18	-3.63	-37.57	-0.27	0.37	-1.25
ko	3.99	7.29	-0.25	-30.15	19.97	21.29	20.5
zh	11.66	11.86	7.94	-25.09	26.47	29.76	27.71
Total	16.55	34.05	5.02	-165.01	56.75	75.68	<u>58.60</u>

Table 12. Intervention effectiveness per target language for AnnSel latents of Qwen3-4B

	Distractor (l_0) zero ablation	$L/\{l\}$ zero ablation	Distractor (l_0) one-layer direction ablation	Distractor (l_0) multi-layer direction ablation	Target (l) amplification	Zero+Amp	1L+Amp
de	-0.67	0	0.02	-2.51	6.73	6.49	6.89
en	4.62	4.9	0.6	8.84	0	0	0
es	1.2	0	0.04	-2.58	6.31	7.28	6.41
fr	-0.71	0	-0.02	-6.21	6.5	6.35	6.56
ja	0.4	0	-0.18	-13.1	1.76	1.83	1.72
ko	0.12	0	-2.42	-5	5.6	7.7	6.08
zh	3.46	0	3.08	-0.4	6.39	9.78	9.73
Total	8.42	4.9	1.12	-20.96	33.29	39.43	<u>37.39</u>

Table 13. Logit difference for ValSel and Gemma-2-2B before and after the intervention using **Antonyms**

	es		en		zh		de		ja		fr		ko	
	before	after	before	after	before	after	before	after	before	after	before	after	before	after
de	-1.79	-2.93	-4.57	-5.19	-2.97	-4.88	-	-	-2.62	-2.69	-0.24	-0.98	-2.38	-2.21
en	-2.29	-3.29	-	-	-3.77	-4.41	-1.45	-2.35	-2.96	-2.71	-0.10	-0.66	-2.87	-3.27
es	-	-	-4.15	-5.83	-3.77	-4.56	-1.30	-2.68	-2.27	-2.58	-0.08	-2.03	-2.30	-2.25
fr	-1.22	-3.19	-4.13	-5.08	-1.89	-4.22	-1.02	-2.26	-2.48	-2.65	-	-	-2.22	-2.18
ja	0.95	-1.87	-3.54	-5.90	0.34	-4.45	0.87	-1.73	-	-	1.65	0.25	1.69	-3.29
ko	1.54	-2.37	-3.28	-4.40	-0.44	-6.04	0.80	-2.53	-1.98	-2.54	1.75	-0.61	-	-
zh	2.62	-0.73	-2.29	-5.72	-	-	3.38	-0.77	2.25	-1.52	2.80	0.80	5.79	-1.27
Total	-0.19	-14.38	-21.96	-32.12	-12.50	-28.56	1.28	-12.32	-10.06	-14.69	5.78	-3.23	-2.29	-14.47
Change		+14.19		+10.16		+16.06		+13.6		+4.63		+9.01		+12.18

Title Suppressed Due to Excessive Size

Table 14. Logit difference for FreqSel and Gemma-2-2B before and after the intervention using **Antonyms**

	es		en		zh		de		ja		fr		ko	
	before	after	before	after	before	after	before	after	before	after	before	after	before	after
de	-1.79	-1.59	-4.57	-4.97	-2.97	-7.58	-	-	-2.62	-3.13	-0.24	-0.18	-2.38	-5.82
en	-2.29	-1.55	-	-	-3.77	-7.99	-1.45	-1.17	-2.96	-3.34	-0.10	0.16	-2.87	-5.73
es	-	-	-4.15	-4.62	-3.77	-8.23	-1.30	-1.03	-2.27	-3.21	-0.08	-0.26	-2.30	-5.61
fr	-1.22	-1.75	-4.13	-4.70	-1.89	-7.79	-1.02	-1.04	-2.48	-3.30	-	-	-2.22	-5.90
ja	0.95	-1.45	-3.54	-4.06	0.34	-7.06	0.87	-1.05	-	-	1.65	0.19	1.69	-3.99
ko	1.54	-1.85	-3.28	-4.26	-0.44	-8.35	0.80	-1.30	-1.98	-3.13	1.75	-0.59	-	-
zh	2.62	-0.95	-2.29	-3.69	-	-	3.38	-0.98	2.25	-2.60	2.80	0.01	5.79	-3.54
Total	-0.19	-9.14	-21.96	-26.30	-12.50	-47.00	1.28	-6.57	-10.06	-18.71	5.78	-0.67	-2.29	-30.59
Change		+8.95		+4.34		+34.5		+7.85		+8.65		+6.45		+28.3

Table 15. Logit difference for AnnSel and Gemma-2-2B before and after the intervention using **Antonyms**

	es		en		zh		de		ja		fr		ko	
	before	after	before	after	before	after	before	after	before	after	before	after	before	after
de	-1.79	-3.04	-4.57	-5.32	-2.97	-5.67	-	-	-2.62	-3.63	-0.24	-1.67	-2.38	-2.54
en	-2.29	-3.11	-	-	-3.77	-6.46	-1.45	-3.01	-2.96	-3.69	-0.10	-1.04	-2.87	-2.57
es	-	-	-4.15	-4.84	-3.77	-6.52	-1.30	-3.37	-2.27	-3.62	-0.08	-2.26	-2.30	-2.50
fr	-1.22	-3.06	-4.13	-4.88	-1.89	-5.91	-1.02	-3.23	-2.48	-3.74	-	-	-2.22	-2.49
ja	0.95	-1.62	-3.54	-4.69	0.34	-6.00	0.87	-2.71	-	-	1.65	-0.58	1.69	-2.07
ko	1.54	-2.36	-3.28	-3.33	-0.44	-6.46	0.80	-3.55	-1.98	-3.53	1.75	-1.61	-	-
zh	2.62	-0.86	-2.29	-3.77	-	-	3.38	-2.52	2.25	-2.97	2.80	-0.60	5.79	0.77
Total	-0.19	-14.05	-21.96	-26.83	-12.50	-37.02	1.28	-18.39	-10.06	-21.18	5.78	-7.76	-2.29	-11.40
Change		+13.86		+4.87		+24.52		+19.67		+11.12		+13.54		+9.11

Table 16. Selection method efficacy per language under the **Zero+Amp** intervention for Qwen3-4B. Post-intervention total change in top-1 logit margin for **Antonyms**; Bottom: average change in normalized logProb of the target sequence for **Enumerations**.

		es	en	zh	de	ja	fr	ko	Avg
Antons.	ValSel	7.36	-5.31	17.56	6.38	<u>-1.17</u>	7.49	<u>11.29</u>	<u>6.23</u>
	FreqSel	4.90	6.09	29.87	4.13	-7.34	6.90	20.86	9.34
	AnnSel	<u>7.35</u>	<u>5.12</u>	9.84	<u>4.95</u>	2.23	<u>6.95</u>	6.79	6.17
Enums.	ValSel	<u>0.89</u>	<u>-0.73</u>	<u>0.25</u>	0.95	0.40	0.97	<u>0.61</u>	<u>0.48</u>
	FreqSel	1.50	0.28	1.41	1.90	1.26	1.49	1.81	1.38
	AnnSel	0.73		0.17	0.82	0.18	0.86	0.49	0.46

H. Adjectives Used in Antonyms

We prepared 100 pairs of antonyms to use for the antonym-structured sentences (*The opposite of "adj1" is "*, *adj2*). Keep in mind that there could be multiple suitable words for *adj2* (e.g. the antonym of small could be either big or large). We used Gemini 2.5 to generate them. All the antonyms are shown in Table 17. It is worth mentioning that the words or tokens can be exactly the same in some cases across different languages, most commonly between French and Spanish and between Chinese and Japanese.

Title Suppressed Due to Excessive Size

Table 17. Multilingual Antonym Pairs in the 7 languages

English	French	German	Spanish	Chinese	Japanese	Korean
good <i>bad</i>	bon <i>mauvais</i>	gut <i>schlecht</i>	bueno <i>malo</i>	好 坏	良い 悪い	좋은 나쁜
happy <i>sad, unhappy</i>	heureux <i>triste, mal- heureux</i>	glücklich <i>traurig, unglück- lich</i>	feliz <i>triste, infeliz</i>	开心 难过, 不高兴	嬉しい 悲しい	행복한 슬픈, 불행한
big <i>small</i>	grand <i>petit</i>	groß <i>klein</i>	grande <i>pequeño</i>	大 小	大きい 小さい	큰 작은
hot <i>cold</i>	chaud <i>froid</i>	heiß <i>kalt</i>	caliente <i>frío</i>	热 冷	暑い 寒い, 冷たい	더운 추운, 차가운
fast <i>slow</i>	rapide <i>lent, lente</i>	schnell <i>langsam</i>	rápido <i>lento</i>	快 慢	速い 遅い	빠른 느린
light <i>heavy, dark</i>	léger <i>lourd</i>	leicht <i>schwer</i>	ligero <i>pesado</i>	轻 重	軽い 重い	가벼운 무거운
easy <i>difficult, hard</i>	facile <i>difficile</i>	einfach <i>schwierig, schwer</i>	fácil <i>difícil</i>	容易 难	簡単な 難しい	쉬운 어려운
new <i>old</i>	nouveau <i>vieux, ancien</i>	neu <i>alt</i>	nuevo <i>viejo</i>	新 旧	新しい 古い	새로운 오래된
true <i>false, untrue</i>	vrai <i>faux</i>	wahr <i>falsch, unwahr</i>	verdadero <i>falso, incorrecto</i>	真 假	正しい 間違った	진실한 거짓의, 틀린
alive <i>dead</i>	vivant <i>mort</i>	lebendig <i>tot</i>	vivo <i>muerto</i>	活 死	生きている 死んでいる	살아있는 죽은
full <i>empty</i>	plein <i>vide</i>	voll <i>leer</i>	lleno <i>vacío</i>	满 空	満杯 空っぽ	가득 찬 빈
bright <i>dark, dim</i>	brillant <i>sombre, obscur</i>	hell <i>dunkel</i>	brillante <i>oscuro, opaco</i>	亮 暗	明るい 暗い	밝은 어두운, 흐릿한
strong <i>weak</i>	fort <i>faible</i>	stark <i>schwach</i>	fuerte <i>débil</i>	强 弱	強い 弱い	강한 약한
clean <i>dirty</i>	propre <i>sale</i>	sauber <i>schmutzig</i>	limpio <i>sucio</i>	干净 脏	きれいな 汚い	깨끗한 더러운
open <i>closed</i>	ouvert <i>fermé</i>	offen <i>geschlossen</i>	abierto <i>cerrado</i>	开 关	開いた 閉じた	열린 닫힌
rich <i>poor</i>	riche <i>pauvre</i>	reich <i>arm</i>	rico <i>pobre</i>	富裕 贫穷	裕福な 貧しい	부유한 가난한
beautiful <i>ugly</i>	beau <i>laid</i>	schön <i>hässlich</i>	hermoso <i>feo</i>	美 丑	美しい 醜い	아름다운 못생긴
long <i>short</i>	long <i>court</i>	lang <i>kurz</i>	largo <i>corto</i>	长 短	長い 短い	긴 짧은
wide <i>narrow</i>	large <i>étroit</i>	breit <i>eng</i>	ancho <i>estrecho</i>	宽 窄	広い 狭い	넓은 좁은
hard <i>soft</i>	dur <i>mou</i>	hart <i>weich</i>	duro <i>suave</i>	硬 软	硬い 柔らかい	딱딱한 부드러운
dry <i>wet, moist</i>	sec <i>humide, mouillé</i>	trocken <i>nass</i>	seco <i>mojado, húmedo</i>	干 湿	乾いた 濡れた, 湿った	마른 젖은, 축축한
loud <i>quiet, silent</i>	fort <i>silencieux, doux</i>	laut <i>leise, still</i>	ruidoso <i>tranquilo, silen- cioso</i>	大声 安静, 小声	うるさい 静かな	시끄러운 조용한, 고요한

Title Suppressed Due to Excessive Size

	English	French	German	Spanish	Chinese	Japanese	Korean
1045	early	tôt	früh	temprano	早	早い	이른
1046	<i>late</i>	<i>tard</i>	<i>spät</i>	<i>tarde</i>	晚	遅い	늦은
1047	near	proche	nah	cerca	近	近い	가까운
1048	<i>far</i>	<i>loin</i>	<i>fern, weit</i>	<i>lejos</i>	远	遠い	먼
1049	deep	profond	tief	profundo	深	深い	깊은
1050	<i>shallow</i>	<i>peu profond</i>	<i>flach</i>	<i>superficial, poco profundo</i>	浅	浅い	얕은
1051	bad	mauvais	schlecht	malo	坏	悪い	나쁜
1052	<i>good</i>	<i>bon</i>	<i>gut</i>	<i>bueno</i>	好	良い	좋은
1053	sad	triste	traurig	triste	难过	悲しい	슬픈
1054	<i>happy, joyful</i>	<i>heureux, joyeux</i>	<i>glücklich</i>	<i>feliz, alegre</i>	开心	嬉しい	행복한, 기쁜
1055	small	petit	klein	pequeño	小	小さい	작은
1056	<i>big, large</i>	<i>grand</i>	<i>groß</i>	<i>grande</i>	大	大きい	큰
1057	cold	froid	kalt	frío	冷	寒い	추운
1058	<i>hot, warm</i>	<i>chaud</i>	<i>heiß</i>	<i>caliente, cálido</i>	热	暑い, 熱い	더운, 따뜻한
1059	slow	lent	langsam	lento	慢	遅い	느린
1060	<i>fast, quick</i>	<i>rapide</i>	<i>schnell</i>	<i>rápido, veloz</i>	快	速い, 早い	빠른
1061	heavy	lourd	schwer	pesado	重	重い	무거운
1062	<i>light</i>	<i>léger</i>	<i>leicht</i>	<i>ligero</i>	轻	軽い	가벼운
1063	difficult	difficile	schwierig	difícil	难	難しい	어려운
1064	<i>easy, simple</i>	<i>facile</i>	<i>einfach</i>	<i>fácil, sencillo</i>	容易, 易	簡単な, 簡単	쉬운, 간단한
1065	old	vieux	alt	viejo	旧	古い	오래된
1066	<i>new</i>	<i>nouveau, jeune</i>	<i>neu</i>	<i>nuevo, joven</i>	新	新しい	새로운, 젊은
1067	false	faux	falsch	falso	假	間違った	거짓의
1068	<i>true, correct</i>	<i>vrai, correct</i>	<i>wahr, richtig</i>	<i>verdadero, correcto</i>	真	正しい	진실한, 올바른
1069	dead	mort	tot	muerto	死	死んでいる	죽은
1070	<i>alive, living</i>	<i>vivant</i>	<i>lebendig</i>	<i>vivo, viviente</i>	活, 生	生きている	살아있는, 생생한
1071	empty	vide	leer	vacío	空	空っぽ	빈
1072	<i>full</i>	<i>plein</i>	<i>voll</i>	<i>lleno</i>	满	満杯	가득 찬
1073	dark	sombre	dunkel	oscuro	暗	暗い	어두운
1074	<i>bright, light</i>	<i>brillant, clair</i>	<i>hell</i>	<i>brillante, claro</i>	亮, 光明	明るい	밝은
1075	weak	faible	schwach	débil	弱	弱い	약한
1076	<i>strong</i>	<i>fort</i>	<i>stark</i>	<i>fuerte</i>	强, 强壮	強い, 丈夫な	강한
1077	dirty	sale	schmutzig	sucio	脏	汚い	더러운
1078	<i>clean</i>	<i>propre</i>	<i>sauber</i>	<i>limpio</i>	干净	きれいな, きれい, 綺麗	깨끗한
1079	closed	fermé	geschlossen	cerrado	关	閉じた	닫힌
1080	<i>open</i>	<i>ouvert</i>	<i>offen</i>	<i>abierto</i>	开	開いた	열린
1081	poor	pauvre	arm	pobre	贫穷	貧しい	가난한
1082	<i>rich, wealthy</i>	<i>riche</i>	<i>reich</i>	<i>rico, adinerado</i>	富裕	裕福な, 豊かな	부유한, 풍부한
1083	ugly	laid	hässlich	feo	丑	醜い	못생긴
1084	<i>beautiful, pretty</i>	<i>beau, joli</i>	<i>schön</i>	<i>hermoso, bonito</i>	美	美しい	아름다운, 예쁜
1085	short	court	kurz	corto	短	短い	짧은
1086	<i>long</i>	<i>long</i>	<i>lang</i>	<i>largo</i>	长	長い	긴
1087	narrow	étroit	eng	estrecho	窄	狭い	좁은
1088	<i>wide, broad</i>	<i>large</i>	<i>breit</i>	<i>ancho</i>	宽	広い	넓은
1089							
1090							
1091							
1092							
1093							
1094							
1095							
1096							
1097							
1098							
1099							

Title Suppressed Due to Excessive Size

	English	French	German	Spanish	Chinese	Japanese	Korean
1100	soft	mou	weich	suave	软	柔らかい	부드러운
1101	<i>hard, firm</i>	<i>dur</i>	<i>hart</i>	<i>duro</i>	硬	硬い	딱딱한
1102	wet	humide	nass	mojado	湿	濡れた	젖은
1103	<i>dry</i>	<i>sec</i>	<i>trocken</i>	<i>seco</i>	干	乾いた	마른
1104	quiet	silencieux	leise	tranquilo	安静	静かな	조용한
1105	<i>loud, noisy</i>	<i>fort, bruyant</i>	<i>laut</i>	<i>ruidoso</i>	大声, 吵闹	うるさい	시끄러운
1106	late	tard	spät	tarde	晚	遅い	늦은
1107	<i>early</i>	<i>tôt</i>	<i>früh</i>	<i>temprano</i>	早	早い	이른
1108	far	loin	fern	lejos	远	遠い	먼
1109	<i>near, close</i>	<i>proche</i>	<i>nah</i>	<i>cerca</i>	近	近い	가까운
1110	shallow	peu profond	flach	superficial	浅	浅い	얕은
1111	<i>deep</i>	<i>profond</i>	<i>tief, hoch</i>	<i>profundo</i>	深	深い	깊은
1112	young	jeune	jung	joven	年轻	若い	젊은
1113	<i>old</i>	<i>vieux</i>	<i>alt</i>	<i>viejo</i>	老	老いた	늙은
1114	kind	gentil	freundlich	amable	善良	親切な	친절한
1115	<i>unkind, mean</i>	<i>méchant</i>	<i>unfreundlich</i>	<i>desagradable, malo</i>	不善良, 邪恶	意地悪な, 不親切な, 失礼な	불친절한, 못된
1116	brave	courageux	mutig	valiente	勇敢	勇敢な	용감한
1117	<i>cowardly</i>	<i>lâche</i>	<i>feige</i>	<i>cobarde</i>	懦弱	臆病な	겁많은
1118	wise	sage	weise	sabio	明智	賢い	현명한
1119	<i>foolish</i>	<i>insensé, stupide</i>	<i>dumm, törricht</i>	<i>tonto, necio</i>	愚蠢	愚かな	어리석은
1120	polite	poli	höflich	educado	礼貌	丁寧な	예의 바른
1121	<i>impolite, rude</i>	<i>impoli, grossier</i>	<i>unhöflich, grob</i>	<i>descortés, maleducado</i>	不礼貌	失礼な, 無礼な, 粗略な, 粗雑な	무례한
1122	patient	patient	geduldig	paciente	耐心	我慢強い	인내심 있는
1123	<i>impatient</i>	<i>impatient</i>	<i>ungeduldig</i>	<i>impaciente</i>	不耐烦	せっかちな	성급한, 참을성 없는
1124	honest	honnête	ehrlich	honesto	诚实	正直な	정직한
1125	<i>dishonest</i>	<i>malhonnête</i>	<i>unehrlich</i>	<i>deshonesto</i>	不诚实	不正直な	부정직한
1126	safe	sûr	sicher	seguro	安全	安全な	안전한
1127	<i>dangerous, unsafe</i>	<i>dangereux</i>	<i>gefährlich, sicher</i>	<i>peligroso, inseguro</i>	危险	危険な	위험한
1128	active	actif	aktiv	activo	积极	活動的な	활동적인
1129	<i>inactive, passive</i>	<i>inactif, passif</i>	<i>inaktiv, passiv</i>	<i>inactivo, pasivo</i>	消极	消極的な	비활동적인, 소극적인
1130	straight	droit	gerade	recto	直	まっすぐな	곧은
1131	<i>curved, bent</i>	<i>courbe, tordu</i>	<i>gebogen, krumm</i>	<i>curvo, doblado</i>	弯	曲かった	굽은, 휘어진
1132	whole	entier	ganz	entero	完整	全体の	전체의
1133	<i>part, broken</i>	<i>partiel, cassé</i>	<i>teilweise, brochen</i>	<i>geparcial, roto</i>	部分, 破	部分的な, 壊れた	부분적인, 부서진
1134	calm	calme	ruhig	calmado	平静	穏やかな	차분한
1135	<i>agitated, stormy</i>	<i>agité, orageux</i>	<i>aufgeregt, misch</i>	<i>agitado, tormentoso</i>	激动, 骚乱	荒れた, 興奮した	격앙된, 폭풍우 치는
1136	correct	correct	richtig	correcto	正确	正しい	정확한
1137	<i>incorrect, wrong</i>	<i>incorrect, faux</i>	<i>falsch</i>	<i>incorrecto, equivocado</i>	不正确, 错误	間違っている	부정확한, 틀린
1138	complex	complexe	komplex	complejo	复杂	複雑な	복잡한
1139	<i>simple</i>	<i>simple</i>	<i>einfach</i>	<i>simple, sencillo</i>	简单	単純な	간단한, 단순한
1140	effective	efficace	effektiv	efectivo	有效	効果的な	효과적인
1141	<i>ineffective</i>	<i>inefficace</i>	<i>ineffektiv</i>	<i>ineficaz</i>	无效	非効果的な	비효과적인

Title Suppressed Due to Excessive Size

	English	French	German	Spanish	Chinese	Japanese	Korean
1155	famous	célèbre	berühmt	famoso	著名	有名な	유명한
1156	<i>unknown, ob-</i>	<i>inconnu, ob-</i>	<i>unbekannt, unbe-</i>	<i>desconocido, os-</i>	<i>无名, 不為人</i>	<i>無名の, 知られてい</i>	<i>무명의, 잘 알려</i>
1157	<i>scure</i>	<i>scur</i>	<i>deutend</i>	<i>curo</i>	<i>知</i>	<i>ない</i>	<i>지지 않은</i>
1158	generous	généreux	großzügig	generoso	慷慨	寛大な	관대한
1159	<i>stingy, mean</i>	<i>avare, mesquin</i>	<i>geizig</i>	<i>tacaño, malo</i>	<i>吝嗇</i>	<i>ケチな</i>	<i>인색한, 못된</i>
1160	content	content	zufrieden	contento	高兴	幸せな	행복한
1161	<i>unhappy, dis-</i>	<i>mécontent, in-</i>	<i>unzufrieden</i>	<i>infeliz, insatisf-</i>	<i>不高兴, 不满</i>	<i>不幸せな</i>	<i>불행한, 불만족</i>
1162	<i>satisfied</i>	<i>satisfait</i>		<i>cho</i>	<i>意</i>		<i>스러운</i>
1163	healthy	sain	gesund	saludable	健康	健康な	건강한
1164	<i>unhealthy, sick</i>	<i>malsain, malade</i>	<i>ungesund, krank</i>	<i>insalubre, en-</i>	<i>不健康, 生病</i>	<i>不健康な, 病気の</i>	<i>건강하지 않은, 아픈</i>
1165	high	haut	hoch	alto	高	高い	높은
1166	<i>low</i>	<i>bas</i>	<i>niedrig</i>	<i>bajo</i>	<i>低</i>	<i>低い</i>	<i>낮은</i>
1167	important	important	wichtig	importante	重要	重要な	중요한
1168	<i>unimportant, trivial</i>	<i>sans impor-</i>	<i>unwichtig, trivial</i>	<i>sin importance, trivial</i>	<i>不重要, 琐碎</i>	<i>重要でない, 取るに</i>	<i>중요하지 않은, 사소한</i>
1169	innocent	innocent	unschuldig	inocente	无辜	無罪の	무고한
1170	<i>guilty</i>	<i>coupable</i>	<i>schuldig</i>	<i>culpable</i>	有罪	有罪の	유죄의
1171	known	connu	bekannt	conocido	已知	既知の	알려진
1172	<i>unknown</i>	<i>inconnu</i>	<i>unbekannt</i>	<i>desconocido</i>	未知	未知の	알려지지 않은
1173	male	masculin	männlich	masculino	男性	男性の	남성의
1174	<i>female</i>	<i>féminin</i>	<i>weiblich</i>	<i>femenino</i>	女性	女性の	여성의
1175	normal	normal	normal	normal	正常	普通の	정상적인
1176	<i>abnormal, un-</i>	<i>anormal, inusuel</i>	<i>abnormal, ungewöhnlich</i>	<i>anormal, inusual</i>	异常, 不正常	普通の	비정상적인, 이
1177	possible	possible	möglich	posible	可能	가능한	가능한
1178	<i>impossible</i>	<i>impossible</i>	<i>unmöglich</i>	<i>imposible</i>	不可能	不可能な	불가능한
1179	private	privé	privat	privado	私人	個人の	사적인
1180	<i>public</i>	<i>public</i>	<i>öffentlich</i>	<i>público</i>	公共	公共の	공적인
1181	right	juste	richtig	correcto	对	正しい	올바른
1182	<i>wrong</i>	<i>faux</i>	<i>falsch</i>	<i>incorrecto, equivo-</i>	错	間違っている	틀린
1183	simple	simple	einfach	sencillo	简单	簡単な	단순한
1184	<i>complex</i>	<i>complexe</i>	<i>komplex</i>	<i>complejo</i>	复杂	複雑な	복잡한
1185	sweet	doux	süß	dulce	甜	甘い	달콤한
1186	<i>sour, bitter</i>	<i>acide, amer</i>	<i>sauer, bitter</i>	<i>agrio, amargo</i>	酸, 苦	酸っぱい, 苦い	신, 쓴
1187	visible	visible	sichtbar	visible	可见	見える	보이는
1188	<i>invisible</i>	<i>invisible</i>	<i>unsichtbar</i>	<i>invisible</i>	不可见	見えない	보이지 않는
1189	warm	chaud	warm	cálido	暖和	暖かい	따뜻한
1190	<i>cool, cold</i>	<i>frais, froid</i>	<i>kühl, kalt</i>	<i>fresco, frío</i>	凉, 冷	涼しい, 冷たい	시원한, 차가운
1191	smooth	lisse	glatt	liso	光滑	滑らかな	매끄러운
1192	<i>rough, bumpy</i>	<i>rugueux, bosselé</i>	<i>rau, holprig</i>	<i>áspero, irregular</i>	粗糙	粗い	거친, 울퉁불퉁한
1193	thin	mince	dünn	delgado	薄	薄い	얇은
1194	<i>thick, fat</i>	<i>épais</i>	<i>dick</i>	<i>grueso</i>	厚	厚い	두꺼운
1195	tall	grand	groß	alto	高	背が高い	키가 큰
1196	<i>short</i>	<i>petit</i>	<i>klein</i>	<i>bajo</i>	矮	背が低い	키가 작은

Title Suppressed Due to Excessive Size

	English	French	German	Spanish	Chinese	Japanese	Korean
1210	married	marié	verheiratet	casado	已婚	既婚の	결혼한
1211	<i>single, unmar-</i>	<i>célibataire</i>	<i>ledig, unver-</i>	<i>soltero</i>	单身, 未婚	独身の, 未婚の	독신의, 미혼의
1212	<i>ried</i>		<i>heiratet</i>				
1213	optimistic	optimiste	optimistisch	optimista	乐观	楽観的な	낙관적인
1214	<i>pessimistic</i>	<i>pessimiste</i>	<i>pessimistisch</i>	<i>pesimista</i>	悲观	悲観的な	비관적인
1215	permanent	permanent	permanent	permanente	永久	恒久的な	영구적인
1216	<i>temporary</i>	<i>temporaire</i>	<i>temporär</i>	<i>temporal</i>	临时	一時的な	일시적인
1217	present	présent	gegenwärtig	presente	现在	現在の	현재의
1218	<i>absent, past</i>	<i>absent, passé</i>	<i>abwesend, vergan-</i>	<i>ausente, pasado</i>	缺席, 过去	不在の, 過去の	부재의, 과거의
1219			<i>gen</i>				
1220	public	public	öffentlich	público	公共	公共の	공공의
1221	<i>private</i>	<i>privé</i>	<i>privat</i>	<i>privado</i>	私人	個人の	사적인
1222	real	réel	echt	real	真实	本物の	실제적인
1223	<i>fake, unreal</i>	<i>faux, irréel</i>	<i>falsch, unreal</i>	<i>falso, irreal</i>	假, 虚假	偽物の, 非現実の	가짜의, 비현실적인
1224							
1225	responsible	responsable	verantwortlich	responsable	负责	責任	책임 있는
1226	<i>irresponsible</i>	<i>irresponsable</i>	<i>unverantwortlich</i>	<i>irresponsable</i>	不负责	無責任な, 無責任	무책임한
1227	single	célibataire	ledig	soltero	单身	独身の	독신의
1228	<i>married</i>	<i>marié</i>	<i>verheiratet</i>	<i>casado</i>	已婚	既婚の, 結婚している	결혼한
1229							
1230	sour	acide	sauer	agrio	酸	酸っぱい	신
1231	<i>sweet</i>	<i>doux</i>	<i>süß</i>	<i>dulce</i>	甜	甘い	달콤한
1232	useful	utile	nützlich	útil	有用	役に立つ	유용한
1233	<i>useless</i>	<i>inutile</i>	<i>nutzlos</i>	<i>inútil</i>	没用	役に立たない	쓸모없는
1234							
1235	vertical	vertical	vertikal	vertical	垂直	垂直な	수직의
1236	<i>horizontal</i>	<i>horizontal</i>	<i>horizontal</i>	<i>horizontal</i>	水平	水平な	수평의
1237	well	bien	gut	bien	好	良い	잘
1238	<i>poorly</i>	<i>mal</i>	<i>schlecht</i>	<i>mal, pobremente</i>	差	悪い	못, 형편없이
1239	winning	gagnant	gewinnend	ganador	获胜	勝利の	승리하는
1240	<i>losing</i>	<i>perdant</i>	<i>verlierend</i>	<i>perdedor</i>	失败	敗北の	패배하는
1241	unkind	méchant	unfreundlich	desagradable	不善良	意地悪な	불친절한
1242	<i>kind</i>	<i>gentil</i>	<i>freundlich</i>	<i>amable</i>	善良	親切な	친절한
1243							
1244	cowardly	lâche	feige	cobarde	懦弱	臆病な	겁많은
1245	<i>brave</i>	<i>courageux</i>	<i>mutig</i>	<i>valiente</i>	勇敢	勇敢な	용감한
1246							
1247							

I. Listings Used in Enumerations

We prepared 7 ordered enumerations to use for the enumeration-structured sentences (*The {category} are: {choices}*). We used Gemini 3 to translate them into the seven languages. All the categories, listings, and the split of the listed choices and continuation words are shown in Table 18.

J. Additional Redundancy Analysis

To assess the equivalence of the latent sets identified by the three selection methods, we compute representative per-layer activation space directions that correspond to the decoder directions of these layer-wise latents, then calculate the pairwise cosine similarity between per-method directions. If the latents of two selection methods for the same language are identical or equivalent in their representation in a certain layer, we expect their corresponding residual stream directions to be similar (large cosine similarity), while a small cosine similarity in a certain layer might indicate divergence. We report in Table 19 for each pair of selection methods the minimum cosine similarity, as well as the across-layer average, for the language latents of Gemma-2-2B and Qwen3-4B.

Finally, we analyze the Neuronpedia annotations of the identified latent sets as a (potentially noisy) means to discern latent set overlap. We compute as a proxy the percentage of language latents whose Neuronpedia annotation refers explicitly to

1265 the language. While this is 100% for `AnnSel` by design, the results for `ValSel` and `FreqSel` in Table 20 reveals substantial
1266 variability.

1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319

Title Suppressed Due to Excessive Size

Table 18. Enumerations Dataset: categories, items and split thresholds.

Category ID	Split	Lang	Category Name	Listings
months	4 8	en	months of the year	January, February, March, April, May, June, July, August, September, October, November, December
		de	Monate des Jahres	Januar, Februar, März, April, Mai, Juni, Juli, August, September, Oktober, November, Dezember
		fr	mois de l'année	janvier, février, mars, avril, mai, juin, juillet, août, septembre, octobre, novembre, décembre
		es	meses del año	enero, febrero, marzo, abril, mayo, junio, julio, agosto, septiembre, octubre, noviembre, diciembre
		zh	一年中的月份	一月, 二月, 三月, 四月, 五月, 六月, 七月, 八月, 九月, 十月, 十一月, 十二月
		ja ko	一年の月 1년의 열두달	一月, 二月, 三月, 四月, 五月, 六月, 七月, 八月, 九月, 十月, 十一月, 十二月 1월, 2월, 3월, 4월, 5월, 6월, 7월, 8월, 9월, 10월, 11월, 12월
numbers	4 6	en	numbers 1 to 10	one, two, three, four, five, six, seven, eight, nine, ten
		de	Zahlen von eins bis zehn	eins, zwei, drei, vier, fünf, sechs, sieben, acht, neun, zehn
		fr	nombres de un à dix	un, deux, trois, quatre, cinq, six, sept, huit, neuf, dix
		es	números del uno al diez	uno, dos, tres, cuatro, cinco, seis, siete, ocho, nueve, diez
		zh	从一到十的数字	一, 二, 三, 四, 五, 六, 七, 八, 九, 十
		ja ko	一から十までの数字 1부터10까지의숫자	一, 二, 三, 四, 五, 六, 七, 八, 九, 十 일, 이, 삼, 사, 오, 육, 칠, 팔, 구, 십
days_of_week	3 4	en	days of the week	Sunday, Monday, Tuesday, Wednesday, Thursday, Friday, Saturday
		de	Wochentage	Sonntag, Montag, Dienstag, Mittwoch, Donnerstag, Freitag, Samstag
		fr	jours de la semaine	dimanche, lundi, mardi, mercredi, jeudi, vendredi, samedi
		es	días de la semana	domingo, lunes, martes, miércoles, jueves, viernes, sábado
		zh	一星期中的日子	星期日, 星期一, 星期二, 星期三, 星期四, 星期五, 星期六
		ja ko	曜日 요일	日曜日, 月曜日, 火曜日, 水曜日, 木曜日, 金曜日, 土曜日 일요일, 월요일, 화요일, 수요일, 목요일, 금요일, 토요일
four_seasons	2 2	en	four seasons	spring, summer, fall, winter
		de	vier Jahreszeiten	Frühling, Sommer, Herbst, Winter
		fr	quatre saisons	printemps, été, automne, hiver
		es	cuatro estaciones	primavera, verano, otoño, invierno
		zh	四季	春, 夏, 秋, 冬
		ja ko	四季 사계절	春, 夏, 秋, 冬 봄, 여름, 가을, 겨울
times_of_day	2 2	en	times of day	morning, afternoon, evening, night
		de	Tageszeiten	Morgen, Nachmittag, Abend, Nacht
		fr	moments de la journée	matin, après-midi, soir, nuit
		es	momentos del día	mañana, tarde, tarde, noche
		zh	一天的时段	早上, 下午, 晚上, 夜里
		ja ko	一日の時間帯 하루의시간대	朝, 午後, 晚, 夜 아침, 오후, 저녁, 밤
cardinal_directions	2 2	en	cardinal directions	North, South, East, West
		de	Himmelsrichtungen	Norden, Süden, Osten, Westen
		fr	points cardinaux	nord, sud, est, ouest
		es	puntos cardinales	norte, sur, este, oeste
		zh	基本方位	北, 南, 东, 西
		ja ko	基本方位 방위	北, 南, 東, 西 북, 남, 동, 서
primary_colors	1 2	en	primary colors of light	red, green, blue
		de	Primärfarben des Lichts	Rot, Grün, Blau
		fr	couleurs primaires de la lumière	rouge, vert, bleu
		es	colores primarios de la luz	rojo, verde, azul
		zh	光的三原色	红, 绿, 蓝
		ja ko	光の三原色 빛의삼원색	赤, 綠, 青 빨강, 초록, 파랑

Table 19. Minimum and average cosine similarity between representative layer-wise activation space directions per language for Gemma-2-2B (top) and Qwen3-4B (bottom).

			es	en	zh	de	ja	fr	ko
Gemma-2-2B	AnnSel - ValSel	Avg	0.66	0.06	0.58	0.59	0.71	0.45	0.30
		Min	0.47	-0.1	0.44	0.41	0.53	0.07	-0.07
	FreqSel - AnnSel	Avg	0.57	0.04	0.59	0.63	0.47	0.57	0.41
		Min	0.09	-0.15	0.25	0.24	0.18	0.12	0.05
	ValSel - FreqSel	Avg	0.34	0.04	0.45	0.29	0.42	0.25	0.30
		Min	-0.1	-0.37	0.11	0.07	0.08	-0.05	-0.44
Qwen3-4B	AnnSel - ValSel	Avg	0.07	0.26	0.47	0.34	0.35	0.3	0.62
		Min	-0.07	-0.22	0.06	-0.04	0.02	-0.04	0.31
	FreqSel - AnnSel	Avg	0.41	0.37	0.38	0.48	0.55	0.32	0.35
		Min	-0.02	0.04	0.02	0.23	0.11	0.00	0.06
	ValSel - FreqSel	Avg	0.42	0.14	0.21	0.27	0.46	0.36	0.35
		Min	-0.03	-0.17	-0.32	-0.01	0.00	-0.15	-0.35

Table 20. The percentage (%) of extracted latents of Gemma-2-2B (top) and Qwen3-4B (bottom) per method and language whose Neuronpedia annotations explicitly refer to the language. AnnSel is 100% by design.

			en	fr	de	es	zh	ko	ja
Gemma-2-2B	ValSel	2	72	94	54	58	2	46	
	FreqSel	0.1	71	77	53	39	1	58	
	AnnSel	100	100	100	100	100	100	100	
Qwen3-4B	ValSel	2	24	32	24	12	0	8	
	FreqSel	1.1	32.9	38.2	28.4	37.5	4.1	16.7	
	AnnSel	100	100	100	100	100	100	100	