

Causality \neq Invariance: FUNCTION VS CONCEPT VECTORS IN LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

Do large language models (LLMs) represent concepts abstractly, i.e., independent of input format? We revisit Function Vectors (\mathcal{FV} s), compact representations of incontext learning (ICL) tasks that causally drive task performance. Across multiple LLMs, we show that \mathcal{FV} s are not fully invariant: \mathcal{FV} s of the same concept are nearly orthogonal when extracted from different input formats (e.g., open-ended vs. multiple-choice). We introduce Concept Vectors (\mathcal{CV} s) which produce more stable concept representations. Like $\mathcal{FV}s$, $\mathcal{CV}s$ are composed of attention head outputs; however, unlike \mathcal{FV} s, head selection is optimized via Representational Similarity Analysis (RSA) to encode concepts consistently across input formats. While these heads emerge in similar layers to \mathcal{FV} -related heads, the two sets are largely distinct, suggesting different underlying mechanisms. Steering experiments reveal that \mathcal{FV} s excel in-distribution, when extraction and application formats match (e.g., both open-ended in English), while CVs generalize better out-of-distribution across both question types (open-ended vs. multiple-choice) and languages. Our results show that LLMs do contain abstract concept representations, but these differ from those that drive ICL performance.

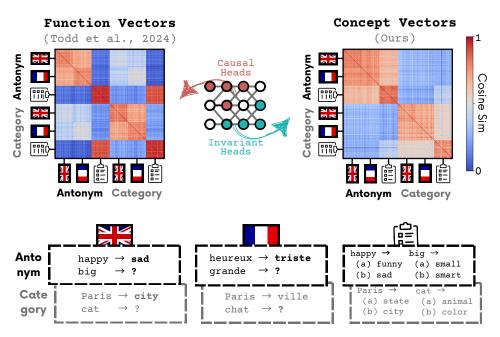


Figure 1: Function vs. Concept Vectors. Top: Similarity matrices for \mathcal{FV} s (left) and \mathcal{CV} s (right) in Llama 3.1 70B; cells show how similar two prompt representations are (warmer = more similar). Middle: Schematic highlighting the distinction between heads with causal effect (AP-selected) and heads that encode format-invariant structure (RSA-selected). Bottom: Example prompts for two concepts across three formats (EN open-ended, FR open-ended, multiple-choice). **Takeaway**: \mathcal{FV} s cluster by input format; CVs cluster by concept across formats.

1 Introduction

Do large language models represent concepts abstractly, i.e., in a way that is stable across surface form? Cognitive science argues that such abstraction underlies human generalization (Gentner, 1983; Hofstadter, 1995; Mitchell, 2020), and recent work shows that LLMs exhibit representational structures similar to humans (Pinier et al., 2025; Du et al., 2025; Doerig et al., 2025). We show that LLMs do contain abstract concept information, but the components that encode it differ from those that causally drive in-context learning (ICL) behavior.

We revisit Function Vectors ($\mathcal{FV}s$)—compact vectors formed by summing outputs of a small set of attention heads that mediate ICL (Todd et al., 2024; Hendel et al., 2023; Yin & Steinhardt, 2025). Because $\mathcal{FV}s$ transfer across contexts (e.g., differently formatted prompts and natural text), they are often treated as encoding the underlying concept (Zheng et al., 2024; Griffiths et al., 2025; Bakalova et al., 2025; Brumley et al., 2024; Fu, 2025). We update this view: $\mathcal{FV}s$ are not fully invariant. For the same concept, $\mathcal{FV}s$ extracted from different input formats (open-ended vs. multiple-choice) are nearly orthogonal (cosine similarity = 0.9), indicating that $\mathcal{FV}s$ mix concept with format (§2.2.1).

To isolate format-invariant structure, we contrast *activation patching* (AP), which localizes components with causal effects on outputs, with *representational similarity analysis* (RSA) (Kriegeskorte, 2008), which localizes components whose representations organize by concept independent of format. Using RSA to select heads and then summing their activations yields Concept Vectors ($\mathcal{CV}s$). Across seven concepts, three input formats (open-ended English, open-ended French, multiple-choice), and four models (Llama 3.1 8B/70B; Qwen 2.5 7B/72B), we find that \mathcal{CV} heads arise in similar layers but are largely disjoint from \mathcal{FV} heads, suggesting separable mechanisms for invariance vs. causality (§2.2.2).

Finally, we test whether $\mathcal{CV}s$ can steer. In steering experiments, $\mathcal{FV}s$ produce larger in-distribution gains when extraction and application formats match (§3.2.1), whereas $\mathcal{CV}s$ generalize more consistently out-of-distribution across question type and language (§3.2.2) and produce fewer format artifacts (e.g., tokens and language from extraction prompts; §3.2.3).

Overall, our contributions are as follows:

- \mathcal{FV} s are not input-invariant. They mix concept with input format; same-concept \mathcal{FV} s differ sharply across formats.
- RSA reveals CV heads. These heads encode concepts at a higher level of abstraction than FV heads.
- Mechanistic separation. FV and CV heads are largely disjoint, suggesting distinct mechanisms for causality vs invariance.
- Steering trade-off. \mathcal{FV} s steer more strongly in-distribution, while \mathcal{CV} s generalize more consistently out-of-distribution, albeit with smaller absolute gains.

2 IN SEARCH OF INVARIANCE

We test whether concept representations are stable across surface form, using AP (causal heads) and RSA (format-invariant heads) across models, datasets, and formats. We then form Function/Concept Vectors to compare clustering by format vs. concept; AP/RSA heads lie in similar layers but show minimal top-K overlap.

2.1 METHODS

2.1.1 Models

We test Llama 3.1 (8B, 70B) and Qwen 2.5 (7B, 72B) models (Meta AI, 2024; Qwen et al., 2025). All models are autoregressive, residual-based transformers (Vaswani et al., 2023). Each model, f internally comprises of \mathcal{L} layers. Each layer is composed of a multi-layer perceptron (MLP) and J attention heads $a_{\ell j}$ which together produce the vector representation of the last token of layer ℓ , $\mathbf{h}_{\ell} = \mathbf{h}_{\ell-1} + \mathrm{MLP}_{\ell} + \sum_{j \in J} a_{\ell j}$ (Elhage et al., 2021).

2.1.2 TASKS

Datasets We define a dataset as one concept expressed in one input format (e.g., Antonym in openended English). For each dataset we build a set of in-context prompts $P_d = \{p_d^i\}$ where i indexes individual prompts within dataset d. Each prompt contains few-shot input-output examples (x,y) that illustrate the same concept, followed by a query input x_q^i whose target output y_q^i is withheld. The input-output pairs (x,y) were either sourced from prior work or generated using OpenAI's GPT-40 (see Appendix D for details). Example prompts are provided in Appendix A.

Concepts. We consider seven concepts:

- **Antonym** Map a word to one with opposite meaning (e.g., hot \rightarrow cold).
- Categorical Map a word to its semantic category (e.g., apple \rightarrow fruit).
- Causal Map a cause to an effect (e.g., rain \rightarrow wet).
- **Synonym** Map a word to one with similar meaning (e.g., big \rightarrow large).
- **Translation** Translate a word to another language (e.g., house \rightarrow maison).
- **Present–Past** Convert a verb from present to past tense (e.g., run \rightarrow ran).
- **Singular-Plural** Convert a noun from singular to plural (e.g., cat \rightarrow cats).

Input formats. We vary only the prompt's surface format; the (x, y) relation stays the same. Formats:

- Open-ended in English (OE-EN)
- Open-ended in a different language (French or Spanish; OE-FR or OE-ES)
- Multiple-choice in English (MC)

We use 5-shot prompts for open-ended and 3-shot for multiple-choice to reduce computational load given prompt length. Altogether, we have 21 datasets (7 concepts \times 3 input formats). We build 50 prompts per dataset (total N=1050 prompts).

2.1.3 ACTIVATION PATCHING

Activation patching replaces specific activations with cached ones from a *clean* run to assess their impact on the model's output. The cached activations are then inserted into selected model components in a *corrupted* run, where the systematic relationships in the prompt are disrupted. For example, in an antonym ICL task, consider a *clean prompt*: Hot \rightarrow Cold, Big \rightarrow Small, Clean \rightarrow ? and a *corrupted prompt*: House \rightarrow Cold, Eagle \rightarrow Small, Clean \rightarrow ? The goal of activation patching is then to localize model components that push the model to the correct answer, Dirty, on the corrupted prompt.

We compute the *causal indirect effect* (CIE) for each attention head $a_{\ell j}$ as the difference between the probability of predicting the expected token y when processing the corrupted prompt \tilde{p} with and without the transplanted mean activation $\bar{\mathbf{a}}_{\ell j}$ from clean runs:

$$CIE(a_{\ell j}) = f(\tilde{p} \mid \mathbf{a}_{\ell j} := \bar{\mathbf{a}}_{\ell j})[y] - f(\tilde{p})[y]$$
(1)

We then compute the average indirect effect (AIE) over a collection \mathcal{D} of all datasets (§2.1.2).

$$AIE(a_{\ell j}) = \frac{1}{|\mathcal{D}|} \sum_{d \in \mathcal{D}} \frac{1}{|\tilde{\mathcal{P}}_d|} \sum_{\tilde{p}_i \in \tilde{\mathcal{P}}_d} CIE(a_{\ell j})$$
 (2)

where \mathcal{P}_d denotes the set of corrupted prompts for dataset d.

2.1.4 REPRESENTATIONAL SIMILARITY ANALYSIS

To find attention heads encoding concepts invariant to input formats, we employ representational similarity analysis (RSA; Kriegeskorte (2008)).

For each attention head $a_{\ell j}$ we compute representational similarity matrices (RSMs) where v_i denotes the output extracted from $a_{\ell j}$ for the ith prompt $p_i \in P_N$, and $\theta(\cdot,\cdot)$ is a cosine similarity function.

$$RSM = \begin{bmatrix} 1 & \cdots & \theta(v_1, v_N) \\ \vdots & \ddots & \vdots \\ \theta(v_N, v_1) & \cdots & 1 \end{bmatrix}$$
 (3)

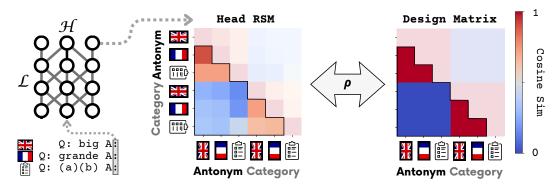


Figure 2: Representational Similarity Analysis (RSA). For each attention head, we compute a representational similarity matrix (RSM) over prompts spanning concepts and input formats (cosine similarity of head outputs). We construct a binary design matrix that marks pairs sharing the same concept, independent of format. The RSA score for a head is Spearman's ρ between the lower-triangular entries of the RSM and the design matrix; higher ρ indicates stronger concept-invariant encoding.

We then construct a binary design matrix, DM, where each entry is set to 1 if the corresponding pair of prompts share the same attribute value, and 0 otherwise. In this paper, we consider two attributes: (1) concept - does a pair of prompts illustrate the same concept, regardless of the input format? and (2) prompt_format - does a pair of prompts have the same question type (i.e. open-ended or multiple-choice)?

We then quantify the alignment between the RSM and DM for the lower-triangles (since similarity matrices are symmetric) using the non-parametric Spearman's rank correlation coefficient (ρ).

To localize attention heads carrying invariant concept information we compute the RSA for each attention head obtaining a single Concept RSA score for each attention head.

Concept-RSA
$$(a_{\ell j}) = \rho(\text{RSM}_{\ell j}, \text{Concept-DM})$$
 (4)

2.1.5 Function & Concept Vectors

To form Function/Concept Vectors we create sets of top K ranking attention heads, $\mathcal{A}_{\mathcal{FV}}$ and $\mathcal{A}_{\mathcal{CV}}$, based on their AIE and RSA scores respectively. Function/Concept Vectors for prompt i are then computed as the sum of activations for this prompt, $\mathbf{a}_{\ell j}^i$, from the sets $\mathcal{A}_{\mathcal{FV}}$ and $\mathcal{A}_{\mathcal{CV}}$ respectively.

$$\mathcal{FV}_{i} = \sum_{a_{\ell j}^{i} \in \mathcal{A}_{\mathcal{FV}}} \mathbf{a}_{\ell j}^{i} \quad \mathcal{CV}_{i} = \sum_{a_{\ell j}^{i} \in \mathcal{A}_{\mathcal{CV}}} \mathbf{a}_{\ell j}^{i}$$

$$(5)$$

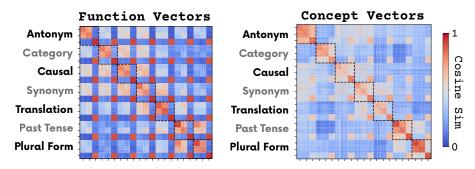


Figure 3: Similarity matrices. Full similarity matrices extracted from top K=5 heads in \mathcal{CV} s and \mathcal{FV} s in Llama 3.1 70B for all concepts. See Appendix B for other models.

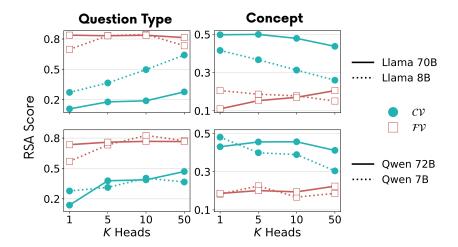


Figure 4: Concept vs. format RSA. Question type and Concept RSA scores for $\mathcal{CV}s$ and $\mathcal{FV}s$ in all models. **Takeaway**: $\mathcal{CV}s$ encode more concept information and less input format than $\mathcal{FV}s$.

2.2 RESULTS

2.2.1 CONCEPT VECTORS ARE MORE INVARIANT TO INPUT FORMAT

We test invariance to input format by computing RSA with design matrices for concept and question type (following the setup in §2.1.4). We form \mathcal{FV} s/ \mathcal{CV} s by summing the top-K heads ranked by AIE/RSA (Eq. 5). Across models and K, \mathcal{CV} s show higher concept RSA and lower question-type RSA than \mathcal{FV} s (Figure 4), indicating that \mathcal{FV} s encode format more strongly while \mathcal{CV} s track concept. Consistently, similarity matrices for Llama 3.1 70B cluster by concept across formats for \mathcal{CV} s, but by format for \mathcal{FV} s (Figure 3), where within-format type \mathcal{FV} clusters are nearly identical with mean cosine similarity = 0.90. \mathcal{CV} s nonetheless exhibit a weaker within-format type cluster (mean cosine similarity = 0.55), suggesting they retain some low-level format information. Overall, however, \mathcal{CV} s remain markedly more invariant to input format than \mathcal{FV} s.

2.2.2 Function & Concept Vectors are Composed of Different Attention Heads

If we compare which heads are selected by the two procedures, we see that \mathcal{FV} s and \mathcal{CV} s are composed of different attention heads. First, we ranked each head for each method, i.e., AIE (§2.1.3) for \mathcal{FV} s and by Concept-RSA (§2.1.4) for \mathcal{CV} s. Then we examined depth and top-K overlap. Layer-averaged scores show similar layer profiles (Figure 5), but head identities barely overlap: for $K \leq 20$ the intersection is near zero and stays small at larger K (Table 1). We also note that AIE scores are highly sparse: their histogram peaks at zero with a long right tail (Figure 12)—so only a few heads have measurable causal effect. Together this supports that AIE-selected *causal* heads are largely distinct from the *invariant*, RSA-selected heads.

Model	K=3	K=5	K=10	K=20	K=50	K=100
Llama-3.1 8B	0	0	1	1	12	28
Llama-3.1 70B	0	0	0	0	1	6
Qwen2.5 7B	0	0	0	4	15	39
Qwen2.5 72B	0	0	0	1	3	13

Table 1: RSA-AIE head overlap. Overlap between RSA and AIE heads (number of overlapping heads among top-K). Bold numbers indicate overlap significantly above chance (p < 0.05; details in Appendix E). **Takeaway**: \mathcal{FV} s and \mathcal{CV} s are composed of different attention heads.

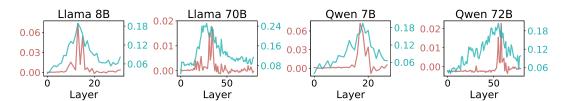


Figure 5: Layer-wise AIE vs. RSA. AIE and RSA scores averaged across all heads per layer. **Takeaway**: \mathcal{FV} and \mathcal{CV} heads are in similar layers.

3 CAN CONCEPT VECTORS STEER?

We now test whether these invariant heads can steer: we introduce how we construct vectors, the AmbiguousICL setup with conflicting cues, and the intervention protocol. \mathcal{FV} s win in-distribution; \mathcal{CV} s transfer better out-of-distribution with fewer format artifacts, at a cost of smaller gains.

3.1 Steering Methods

Steering Vectors Construction. For each concept and input format (OE-ENG, OE-FR, MC), we compute for every selected head $a_{\ell j}$ the mean last-token activation across the 50 *extraction prompts* of that concept-format. We then form one vector per format by summing these mean activations over the top-K heads selected for \mathcal{CV} or \mathcal{FV} (as in Eq. 5, but using per-format means in place of per-prompt activations). This yields one ID vector (OE-ENG) and two OOD vectors (OE-FR, MC) per concept.

AmbigousICL Task. We evaluate on AmbigousICL tasks (Figure 6): each prompt interleaves two concepts (3 then 2 exemplars) followed by a query. The second concept is always English→French translation, which allows us to test whether vectors store low-level language information. Unsteered models tend to continue with the second concept; we aim to steer toward the first. Note that these steering prompts are distinct from the extraction prompts used to construct the vectors.

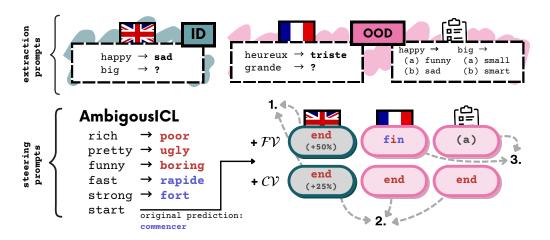


Figure 6: Overview of steering results. **Top**: We extract $\mathcal{CV}s$ and $\mathcal{FV}s$ from antonym ICL prompts in formats that are in-distribution (ID; OE-ENG) or out-of-distribution (OOD; OE-FR, MC) relative to the AmbiguousICL task (bottom-left). **Bottom-left**: We interleave two concepts—antonym and $EN \rightarrow FR$ translation—within one prompt; the model's original prediction is the French translation. **Bottom-right**: Predictions after steering. **Takeaways**: (1) $\mathcal{FV}s$ yield larger ID gains. (2) $\mathcal{CV}s$ show more stable OOD effects across formats. (3) $\mathcal{FV}s$ can conflate concept with input format (e.g., French version of antonym and multiple-choice formatting).

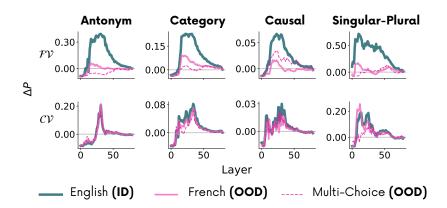


Figure 7: Steering effect across layers. We inject $\mathcal{CV}s$ and $\mathcal{FV}s$ into Llama-3.1-70B and plot the change in target-token probability (ΔP) for four representative concepts (columns). Curves compare ID extraction format with OOD formats relative to the AmbiguousICL task (Figure 6). Higher ΔP means the model assigns more probability to the expected token than the unsteered model. Takeaways: (1) $\mathcal{FV}s$ typically achieve larger ID gains but often drop OOD. (2) $\mathcal{CV}s$ yield smaller gains yet show more stable OOD behavior across formats. See Figure 16 for other concepts/models.

Steering with $\mathcal{CV}s$ and $\mathcal{FV}s$. We add a vector v to the last-token residual stream at a chosen layer:

$$\mathbf{h}_{\ell} \leftarrow \mathbf{h}_{\ell} + \alpha \mathbf{v}$$
 (6)

We measure effectiveness as $\Delta P = P_{\text{after}}(y) - P_{\text{before}}(y)$, averaged over 100 prompts per concept. We sweep α and K and report the best per model (Appendix F).

3.2 Steering Results

3.2.1 Function Vectors Outperform Concept Vectors in Distribution

Extracted from OE-ENG (ID setting), \mathcal{FV} s yield the largest gains on ambiguous prompts (Figure 7). \mathcal{CV} s also help but with smaller ΔP and minimal zeroshot effect (Figure 17). At the token level both vectors lift plausible English antonyms in the ID case (Table 2).

3.2.2 Concept Vectors Are More Stable Out of Distribution

Performance gains (ΔP). Out of distribution (extracting vectors from OE-FR or MC), $\mathcal{CV}s$ more often maintain positive effects across formats, whereas $\mathcal{FV}s$ frequently degrade—especially for MC—and only occasionally stay consistent for specific concepts/models (Figs. 7, 16). $\mathcal{CV}s$ raise the probability of the correct English answer across formats, and their top- Δ tokens remain conceptaligned (Table 2).

Query: salty →							
	+ Antonym	Top Δ Tokens					
\mathcal{FV}	OE-ENG OE-FR MC	_sweet (+56%), _fresh (+16%), _bland (+6%), _taste (+3%), _uns (+2%) _su (+31%), _dou (+27%), _frais (+5%), _fade (+5%), _ins (+3%) _ ((+53%), _A (+1%), _\n (+1%), _space (+0%), _) (+0%)					
CV	OE-ENG OE-FR MC	_sweet (+49%), _fresh (+8%), _bland (+3%), _taste (+3%), _uns (+3%) _sweet (+54%), _fresh (+9%), _bland (+3%), _uns (+3%), _taste (+2%) _sweet (+35%), _fresh (+12%), _bland (+4%), _uns (+3%), _taste (+3%)					

Table 2: *Token-level steering effects*. Top tokens with largest probability gains when injecting $\mathcal{CV}s$ or $\mathcal{FV}s$ into Llama-3.1-70B on the AmbiguousICL prompt (query shown above). Results shown at the layer with the strongest in-distribution effect per vector. Without intervention, the model predicts French $_sa$ (from $sal\acute{e}$) with 49%; antonym $_sweet$ has 2%. English antonyms in red, French in blue, and the opening bracket (MC token) in green.

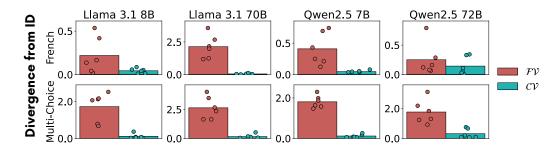


Figure 8: Consistency of steering effects. KL divergence between the probability distributions after the models were steered with an ID (OE-ENG) and OOD vectors (OE-FR [top], MC [bottom]). Lower values indicate more similar effects of ID and OOD vectors. **Takeaway**: \mathcal{CV} s steer the models more consistently than \mathcal{FV} s.

Distributional consistency (KL). To quantify consistency across formats independent of absolute gains, we compare the model's next-token distributions after steering with ID and OOD vectors. For each concept and vector type, we select the top 5 layers that achieve the highest ID ΔP . At each selected layer we compute KL divergence

$$D_{KL}[p(\mathbf{x} \mid \mathbf{v}_{\text{OOD}}) \parallel p(\mathbf{x} \mid \mathbf{v}_{\text{ID}})]$$

between the post-intervention distributions at the query token, where lower values indicate more similar effects of ID and OOD vectors. We average this KL divergence over prompts and selected layers to obtain one score per concept, and then summarize per model (Figure 8). Across models, \mathcal{CV} s yield lower KL than \mathcal{FV} s. The CV–FV KL gap is larger for MC than for OE–FR.

3.2.3 FUNCTION VECTORS MIX CONCEPT WITH INPUT FORMAT

Out of distribution, $\mathcal{FV}s$ reflect both prompt format and concept. When vectors are extracted from OE-FR, they push the model toward the French translation of the concept (e.g., French antonyms), and when extracted from MC, they increase the probability of format tokens such as the opening bracket (Table 2). We quantify the language effect by measuring ΔP for the French translation across concepts (Figure 13). In the larger models, $\mathcal{FV}s$ substantially increase the probability of the French token, whereas $\mathcal{CV}s$ remain near zero; in smaller models the effect is negligible. Notably, $\mathcal{FV}s$ extracted from open-ended Spanish prompts induce almost the same bias toward the French translation as $\mathcal{FV}s$ extracted from French prompts (Figure 14), even though the AmbigousICL alternatives are French only. This pattern suggests that $\mathcal{FV}s$ capture a generic translation/foreign-language signal tied to the extraction format rather than language-specific content. Combined with the MC bracket effect (Figure 15), these findings indicate that $\mathcal{FV}s$ mix concept with surface format, while $\mathcal{CV}s$ are comparatively format-invariant.

4 RELATED WORK

Attention Head Categorization. Recent work has made significant progress in characterizing specialized attention heads that process in-context learning (ICL) tasks. For instance, Olsson et al. (2022) identified induction-heads, which Yin & Steinhardt (2025) found can develop into \mathcal{FV} -heads during training. Other specialized head types include semantic-induction heads (Ren et al., 2024), symbol-abstraction heads (Yang et al., 2025), and various others (Zheng et al., 2024). Our work extends this line of research by identifying \mathcal{CV} heads, attention heads that invariantly represent concepts in ICL tasks at high levels of abstraction.

Linear Representation of Concepts. A substantial body of research has established that concepts are represented linearly in LLMs' representational space (Mikolov et al., 2013; Arora et al., 2016; Elhage et al., 2022). This phenomenon, often termed the "Linear Representation Hypothesis" (Park et al., 2024), has been extensively studied across various tasks and domains. Hernandez et al. (2024) demonstrated that relational concepts—similar to those we study in this paper—can be decoded from LLM activations using linear approximation. Subsequent work by Merullo et al. (2025) revealed

that the success of such decoding depends on the frequency of concepts in the pretraining corpora, which may explain why some concepts are represented more consistently than others in our study. Our findings contribute to this literature in two ways: (1) providing further support for the Linear Representation Hypothesis, and (2) extending previous work on relational concept representations by localizing specific attention heads that carry such representations and demonstrating their invariance to input formats.

Symbolic-like reasoning in LLMs. Recent work has demonstrated that LLMs can exhibit symbol-like representational properties even without explicit symbolic architecture (Feng & Steinhardt, 2024; Yang et al., 2025; Griffiths et al., 2025). Yang et al. (2025) define symbolic processing as requiring two key properties: (1) invariance to content variations, and (2) indirection through pointers rather than direct content storage. Our \mathcal{CV} s exhibit both properties: they are invariant to input format changes and function as pointers to content stored elsewhere, unlike \mathcal{FV} s which directly store content (§3.2.3).

5 DISCUSSION AND LIMITATIONS

Our results separate two representational roles in LLMs: components that cause strong ICL performance and components that encode abstract concept structure. Function Vectors ($\mathcal{FV}s$) occupy the first role, steering models effectively when extraction and application formats match, but deteriorating out of distribution (formats/languages). Conversely, Concept Vectors ($\mathcal{CV}s$) built from RSA-selected heads encode higher-level, format-invariant structure and generalize more robustly across languages and question types, albeit with smaller causal effects. This supports a view that invariance and causality are mediated by largely distinct mechanisms in similar layers; trends hold on average but vary by concept and model size.

Relation to Function Vectors. Prior work shows that \mathcal{FV} s compactly mediate ICL and can transfer across contexts (Todd et al., 2024). We refine this: \mathcal{FV} portability is strong within families of prompts, but is not fully invariant to surface format. Same-concept \mathcal{FV} s extracted from different formats are nearly orthogonal and can carry language/format signals (e.g., French subword or multiple-choice bracket tokens), while \mathcal{CV} s track concept across formats with less surface content. At the token level, we observe that \mathcal{FV} s produce plausible outputs (not merely output-space vocabulary), yet out of distribution they also boost format-specific tokens, consistent with \mathcal{FV} s mixing task procedure with surface constraints.

Implications for steering and interpretability. The dissociation between \mathcal{FV} s and \mathcal{CV} s suggests a practical trade-off. For *maximal in-distribution control*, \mathcal{FV} s are preferable. For *robust out-of-distribution control* or probing abstract knowledge, \mathcal{CV} s are more reliable. Methodologically, AP identifies what causally drives behavior, while RSA reveals how representations organize by concept regardless of format. This distinction highlights that effective behavioral control and abstract conceptual representation can be mediated by different mechanisms.

Limitations and Future Directions. Our \mathcal{CV} head selection targeted heads that encode *all* concepts simultaneously; this global criterion may miss concept-specific heads, which a per-concept RSA could reveal. We also did not probe how \mathcal{FV} s and \mathcal{CV} s emerge during model training or how they interact during inference; we hypothesize that \mathcal{CV} s act as a *backup circuit*—a format-invariant scaffold that stabilizes concept information which downstream \mathcal{FV} pathways can recruit or override—consistent with evidence of multiple, partially redundant circuits and compensatory self-repair under ablations (McGrath et al., 2023; Wang et al., 2022).

REFERENCES

Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. A latent variable model approach to PMI-based word embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399, 2016. doi: 10.1162/tacl_a_00106. URL https://aclanthology.org/Q16-1028/.

Aleksandra Bakalova, Yana Veitsman, Xinting Huang, and Michael Hahn. Contextualize-thenaggregate: Circuits for in-context learning in gemma-2 2b, 2025. URL https://arxiv.org/abs/2504.00132.

Madeline Brumley, Joe Kwon, David Krueger, Dmitrii Krasheninnikov, and Usman Anwar. Comparing bottom-up and top-down steering approaches on in-context learning tasks, 2024. URL https://arxiv.org/abs/2411.07213.

- DeepL SE. Deepl translator, 2025. URL https://www.deepl.com/translator. Accessed: 2025.
- Adrien Doerig, Tim Kietzmann, Emily Allen, Yihan Wu, Thomas Naselaris, Kendrick Kay, and Ian Charest. High-level visual representations in the human brain are aligned with large language models. *Nature Machine Intelligence*, 7:1220–1234, 08 2025. doi: 10.1038/s42256-025-01072-0.
- Changde Du, Kaicheng Fu, Bincheng Wen, Yi Sun, Jie Peng, Wei Wei, Ying Gao, Shengpei Wang, Chuncheng Zhang, Jinpeng Li, Shuang Qiu, Le Chang, and Huiguang He. Human-like object concept representations emerge naturally in multimodal large language models. *Nature Machine Intelligence*, 7(6):860–875, June 2025. ISSN 2522-5839. doi: 10.1038/s42256-025-01049-z. URL http://dx.doi.org/10.1038/s42256-025-01049-z.
- Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021.
- Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition, 2022. URL https://arxiv.org/abs/2209.10652.
- Jiahai Feng and Jacob Steinhardt. How do language models bind entities in context?, 2024. URL https://arxiv.org/abs/2310.17191.
- Shuhao Fu. Function Vectors for Relational Reasoning in Multimodal Large Language Models. PhD thesis, UCLA, 2025.
- Dedre Gentner. Structure-mapping: A theoretical framework for analogy. *Cognitive science*, 7(2): 155–170, 1983.
- Thomas L. Griffiths, Brenden M. Lake, R. Thomas McCoy, Ellie Pavlick, and Taylor W. Webb. Whither symbols in the era of advanced neural networks?, 2025. URL https://arxiv.org/abs/2508.05776.
- Roee Hendel, Mor Geva, and Amir Globerson. In-Context Learning Creates Task Vectors, October 2023.
- Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. Linearity of relation decoding in transformer language models, 2024. URL https://arxiv.org/abs/2308.09124.
- Douglas R Hofstadter. Fluid concepts and creative analogies: Computer models of the fundamental mechanisms of thought. Basic books, 1995.
- Nikolaus Kriegeskorte. Representational similarity analysis connecting the branches of systems neuroscience. *Frontiers in Systems Neuroscience*, 2008. ISSN 16625137. doi: 10.3389/neuro.06. 004.2008.
- Thomas McGrath, Matthew Rahtz, Janos Kramar, Vladimir Mikulik, and Shane Legg. The hydra effect: Emergent self-repair in language model computations, 2023. URL https://arxiv.org/abs/2307.15771.
- Jack Merullo, Noah A. Smith, Sarah Wiegreffe, and Yanai Elazar. On linear representations and pretraining data frequency in language models, 2025. URL https://arxiv.org/abs/2504.12459.

Meta AI. The Llama 3 Herd of Models, November 2024.

Tomas Mikolov, Wen-tau Yih, and Geoffrey Zweig. Linguistic regularities in continuous space word representations. In Lucy Vanderwende, Hal Daumé III, and Katrin Kirchhoff (eds.), *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 746–751, Atlanta, Georgia, June 2013. Association for Computational Linguistics. URL https://aclanthology.org/N13-1090/.

- Melanie Mitchell. *Artificial Intelligence: A Guide for Thinking Humans*. Picador, New York, first picador paperback edition, 2020 edition, 2020. ISBN 978-1-250-75804-0.
- Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads, 2022. URL https://arxiv.org/abs/2209.11895.
- OpenAI. Gpt-4o system card, 2024. URL https://arxiv.org/abs/2410.21276.
- Kiho Park, Yo Joong Choe, and Victor Veitch. The Linear Representation Hypothesis and the Geometry of Large Language Models, July 2024.
- Christopher Pinier, Sonia Acuña Vargas, Mariia Steeghs-Turchina, Dora Matzke, Claire E. Stevenson, and Michael D. Nunez. Large language models show signs of alignment with human neurocognition during abstract reasoning, 2025. URL https://arxiv.org/abs/2508.10057.
- Qwen, :, An Yang, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang, Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang, Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li, Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru Zhang, and Zihan Qiu. Qwen2.5 technical report, 2025. URL https://arxiv.org/abs/2412.15115.
- Jie Ren, Qipeng Guo, Hang Yan, Dongrui Liu, Quanshi Zhang, Xipeng Qiu, and Dahua Lin. Identifying semantic induction heads to understand in-context learning, 2024. URL https://arxiv.org/abs/2402.13055.
- Eric Todd, Millicent L. Li, Arnab Sen Sharma, Aaron Mueller, Byron C. Wallace, and David Bau. Function Vectors in Large Language Models, February 2024.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL https://arxiv.org/abs/1706.03762.
- Kevin Wang, Alexandre Variengien, Arthur Conmy, Buck Shlegeris, and Jacob Steinhardt. Interpretability in the wild: A circuit for indirect object identification in gpt-2 small, 2022. URL https://arxiv.org/abs/2211.00593.
- Yukang Yang, Declan Campbell, Kaixuan Huang, Mengdi Wang, Jonathan Cohen, and Taylor Webb. Emergent symbolic mechanisms support abstract reasoning in large language models, 2025. URL https://arxiv.org/abs/2502.20332.
- Kayo Yin and Jacob Steinhardt. Which attention heads matter for in-context learning?, 2025. URL https://arxiv.org/abs/2502.14010.
- Zifan Zheng, Yezhaohui Wang, Yuxin Huang, Shichao Song, Mingchuan Yang, Bo Tang, Feiyu Xiong, and Zhiyu Li. Attention heads of large language models: A survey, 2024. URL https://arxiv.org/abs/2409.03752.

```
594
       A PROMPT EXAMPLES
595
596
       A.1 OPEN-ENDED (5-SHOT)
597
598
      Q: resistant
      A: susceptible
599
600
      Q: classify
601
      A: disorganize
602
603
      Q: posterior
604
      A: anterior
605
      Q: goofy
606
      A: serious
607
608
      Q: stationary
      A: moving
609
610
      Q: hairy
611
612
613
      A.2 MULTIPLE-CHOICE (3-SHOT)
614
615
      Instruction: Q: unveil A: ?
       (a) optional
616
       (b) mild
617
       (c) con
618
       (d) conceal
619
       Response: (d)
620
      Instruction: Q: hooked A: ?
621
       (a) unhooked
622
       (b) stale
623
      (c) sturdy
624
       (d) sell
625
      Response: (a)
626
      Instruction: Q: spherical A: ?
627
       (a) unconstitutional
628
       (b) flat
629
       (c) demand
       (d) healthy
630
      Response: (b)
631
632
       Instruction: Q: minute A: ?
633
       (a) conservative
634
       (b) hour
635
       (c) retail
       (d) awake
636
       Response: (
637
638
639
640
641
642
643
644
645
646
647
```

B SIMILARITY MATRICES FOR OTHER MODELS

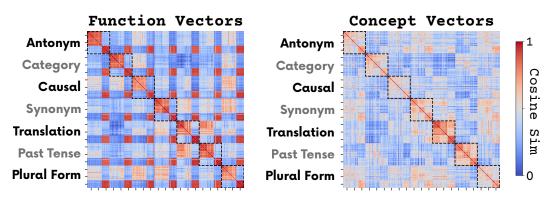


Figure 9: Similarity matrices extracted from top K=1 heads in $\mathcal{CV}s$ and $\mathcal{FV}s$ in Llama 3.1 8B.

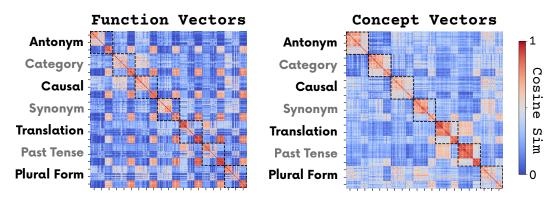


Figure 10: Similarity matrices extracted from top K=1 heads in \mathcal{CV} s and \mathcal{FV} s in Qwen 2.5 7B.

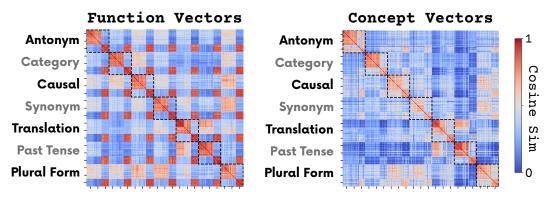


Figure 11: Similarity matrices extracted from top K=2 heads in $\mathcal{CV}s$ and $\mathcal{FV}s$ in Qwen 2.5 72B.

C AIE SCORES

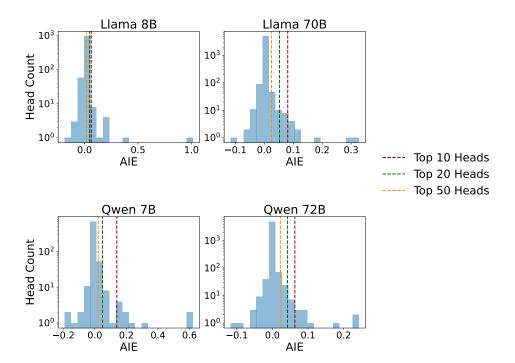


Figure 12: Histogram of AIE scores for Llama 3.1 8B, 70B, Qwen 2.5 7B, and Qwen 2.5 72B. Note, the y-axis is on a log-scale. **Takeaway**: AIE scores are highly sparse.

D DATA GENERATION PROCESS

Concept sourcing: For most concepts (antonym, synonym, translation, present–past, singular–plural), we sourced word pairs from the datasets used by Todd et al. (2024). For categorical and causal concepts, we generated word pairs using OpenAI's GPT-40 model (OpenAI, 2024).

Translation generation: French and Spanish translations were created using DeepL's translation service (DeepL SE, 2025) to ensure high-quality, contextually appropriate translations.

Generated concepts (categorical and causal): We prompted GPT-40 to generate exemplar:category pairs (e.g., "apple:fruit", "blue:colour") and cause:effect pairs (e.g., "stumble:fall", "storm:flood"). The model was given examples of the desired format and asked to produce 100 pairs per batch. We generated pairs in batches of 100 until reaching approximately 1000 examples per concept, with retry mechanisms to ensure sufficient coverage. The final datasets were saved as JSON files containing input-output pairs.

Quality filtering: Generated pairs underwent several filtering steps: (1) removal of duplicates based on input words, (2) exclusion of pairs containing underscores or numbers, (3) restriction to single words or two-word phrases (maximum one space per input/output), and (4) conversion to lowercase for consistency.

Multiple choice format: For multiple choice prompts, we generated four options per question by randomly sampling three additional outputs from the same concept dataset, ensuring all four options were unique. The correct answer was randomly positioned among the four options.

E SIGNIFICANCE TEST FOR RSA-AIE HEAD OVERLAP

We assess whether the observed overlap between the top-K heads selected by Concept-RSA and by AIE is larger than expected by chance under a simple null model. Let N denote the total number of attention heads in the model (layers \times heads per layer). For a fixed K, each method selects a size-K subset of heads. Under the null hypothesis that these two subsets are independent, uniformly random size-K subsets of $\{1,\ldots,N\}$, the overlap size

$$X = |S_{RSA,K} \cap S_{AIE,K}|$$

follows a hypergeometric distribution $X \sim \operatorname{Hypergeom}(N, K, K)$.

For an observed intersection x, we report the one-sided tail probability

$$p_{\geq x} = \Pr\left[X \geq x\right] = \sum_{t=x}^{K} \frac{\binom{K}{t} \binom{N-K}{K-t}}{\binom{N}{K}}.$$

Entries with $p_{>x} < 0.05$ are typeset in bold in Table 1.

F STEERING HYPERPARAMETERS

To optimize the intervention performance, we conduct a hyperparameter search for two parameters:

- α : the steering weight that controls the strength of the intervention
- K: the number of attention heads to extract for concept vector computation

We evaluate the following parameter ranges:

- $K \in \{1, 3, 5, 10, 20, 50\}$ for the number of heads
- $\alpha \in \{1, 3, 5, 10, 15\}$ for the steering weight

The hyperparameter optimization is performed separately for each model using antonym prompts. We select the parameter combination that maximizes the average steering effect across all input formats. This ensures that our chosen hyperparameters generalize well across different prompt structures. We report the best hyperparameters for each model in Table 3.

Model	$\mathbf{Best}\ K$	Best α
Llama 3.1 8B	1	10
Llama 3.1 70B	5	10
Qwen 2.5 7B	3	10
Qwen 2.5 72B	5	15

Table 3: Optimal hyperparameters for steering interventions across different models. K represents the number of attention heads used for $\mathcal{FV/CV}$ extraction, while α controls the intervention strength.

G INPUT FORMAT MIXING IN FUNCTION VECTORS

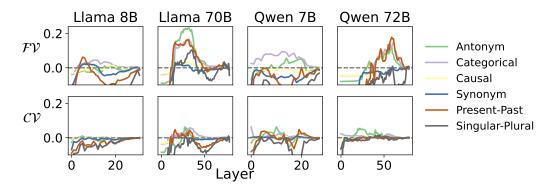


Figure 13: ΔP for French translations of all the concepts. $\mathcal{FV}s$ and $\mathcal{CV}s$ are extracted from openended French prompts.

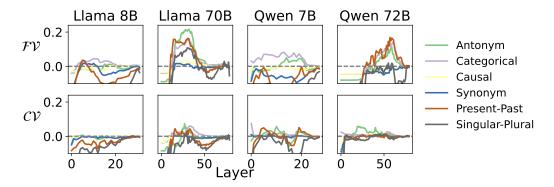


Figure 14: ΔP for French translations of all the concepts. $\mathcal{FV}s$ and $\mathcal{CV}s$ are extracted from openended Spanish prompts.

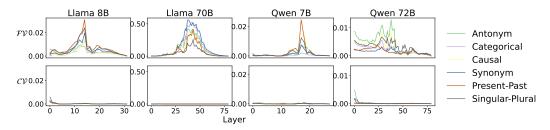


Figure 15: ΔP for the opening bracket token _ (. $\mathcal{FV}s$ and $\mathcal{CV}s$ are extracted from mulitple-choice prompts.

H STEERING RESULTS

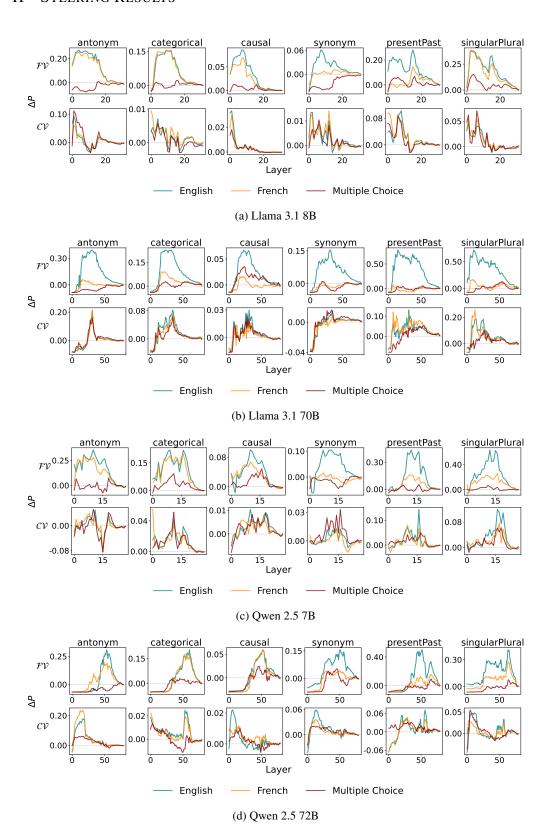


Figure 16: Steering effect across layers and all concepts for different models.

I 0-SHOT STEERING RESULTS

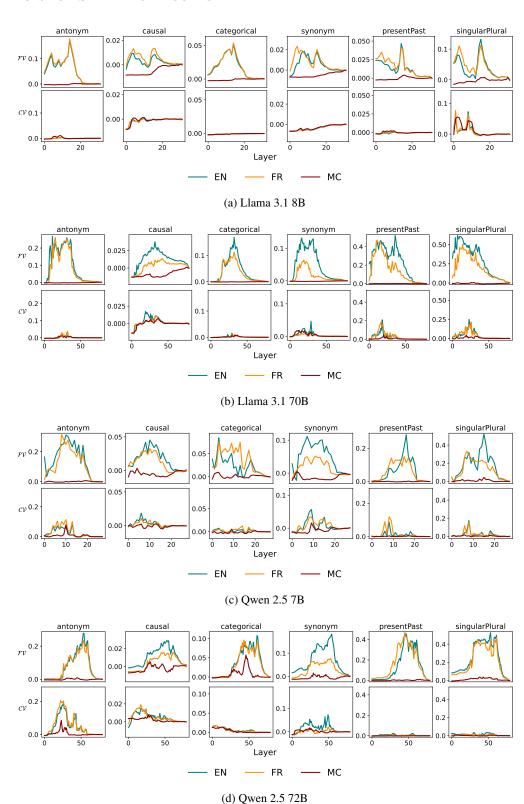


Figure 17: 0-shot steering effect across layers and all concepts for different models.

J QWEN 2.5 72B OUTLIER ANALYSIS

We identified anomalous CIE values for Qwen 2.5 72B in the Categorical concept across French open-ended and multiple-choice formats. As shown in Figure 18, these conditions exhibit unusually high CIE values with a bimodal distribution that deviates from the expected pattern. We excluded these two datasets from the final AIE calculations. This exclusion has minimal impact on our results: the top-5 head rankings remain identical (100% overlap), confirming that our main findings are robust to this methodological decision.

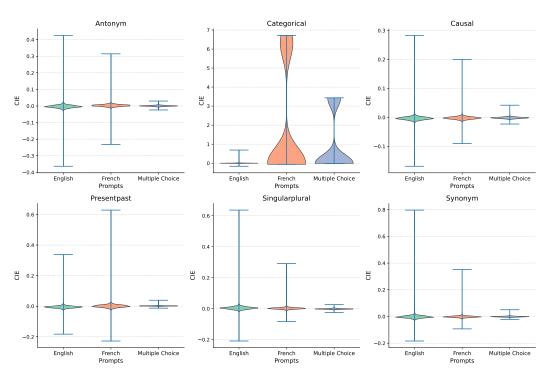


Figure 18: Violin plots of CIE for different concepts and prompts.