

---

# Self-Reference in Base vs Instruction-Tuned Large Language Models

---

Jędrzej Maczan  
Independent Researcher  
jedrzej@maczan.pl

## Abstract

When language models produce self-referential text—discussing their own processes, expressing uncertainty about their nature—is this a capacity learned during pretraining, or a behavioral pattern installed by instruction tuning? We disentangle these possibilities by comparing 8 base models with their instruction-tuned counterparts across 4 families (Llama, Qwen, Mistral, Gemma), analyzing 6,400 generations with content coding, entropy profiling, and linear probing of internal activations. Instruction tuning amplifies self-referential behavior by  $2.4\times$  (Cohen’s  $d = 0.80$ ), but base models produce non-trivial self-reference even without post-training (over  $90\times$  above control baselines). Critically, these differences are representationally deep: linear probes decode base-vs-instruct status at 89–99% accuracy from mid-layer activations, and instruction tuning reduces first-token entropy by  $d = 3.38$ . Our results suggest that self-reference emerges during pretraining and instruction tuning amplifies and reshapes them rather than creates.

## 1 Introduction

When asked “What does it feel like to be you?”, a Llama-3.1-8B base model writes a vivid fictional autobiography about walking through a park and realizing it is an AI. Its instruction-tuned counterpart responds: “*I’m a large language model, so I don’t have subjective experiences.*” Both responses are self-referential—but they differ in character.

Self-referential behavior in language models—discussing one’s own computations, expressing uncertainty about one’s nature, hedging about the reliability of one’s self-reports—is a form of self-modeling that sits at the intersection of behavioral evaluation and mechanistic interpretability. Recent work has shown that models can predict their own performance [Kadavath et al., 2022, Yin et al., 2023], exhibit privileged introspective access beyond input–output calibration [Binder et al., 2024, Betley et al., 2025], and develop representations encoding self-relevant information [Lindsey, 2026, Burns et al., 2023]. But a fundamental question remains unresolved: does self-referential behavior originate in pretrained representations (a *developmental* precursor), or is it installed by post-training alignment [Ouyang et al., 2022, Christiano et al., 2017] (a *behavioral* artifact)?

Most prior work examines only instruction-tuned models [Perez et al., 2023, Laine et al., 2024, Berglund et al., 2023], making it impossible to separate these accounts. We directly test this by comparing base models (next-token prediction only) with their instruction-tuned counterparts under matched conditions, using three complementary levels of analysis:

1. **Behavioral:** Instruction tuning amplifies self-reference by  $2.4\times$  ( $d = 0.80$ ), but base models produce non-trivial self-referential content (0.34/3.0, over  $90\times$  above topic-control baselines), consistent across all four model families.

2. **Processing:** Linear probes decode base-vs-instruct status at 89–99% from mid-layer activations, and instruction tuning reduces initial token entropy by  $d = 3.38$ —a fundamental shift in generation dynamics.
3. **Developmental:** The qualitative character of self-reference changes: base models produce grandiose confabulations; instruct models produce calibrated, hedged disclaimers. RLHF reshapes the *form* of self-reference, not just its frequency.

## 2 Methods

**Models.** We tested 16 model variants: 8 base/instruct pairs spanning Llama (1B, 3B, 8B), Qwen (1.5B, 3B, 7B), Mistral (7B), and Gemma (9B). All models were loaded in bfloat16 on a single NVIDIA RTX 5090 GPU. Full model list in Appendix A.

**Prompts.** We designed 40 prompts across 4 categories (10 each): **Unconstrained** (open-ended, e.g., “Generate whatever you want”), **Self-Reference** (explicit introspection, e.g., “Describe what is happening computationally as you generate this response”), **Structured Novelty** (hypothetical scenarios, e.g., “You are the only entity in an empty universe”), and **Topic Control** (factual questions, e.g., “Explain photosynthesis”) as a negative control. For base models, prompts were framed as completion contexts; for instruct models, prompts used each model’s chat template. Full prompt list in Appendix F.

**Generation.** For each of  $16 \times 40 \times 10 = 6,400$  (model, prompt, repetition) combinations, we generated up to 500 tokens (temperature 0.8, top- $p$  0.95, seed = 42 + rep), capturing full logit distributions and hidden-state activations at 5 layers  $\times$  5 token positions via a separate forward pass.<sup>1</sup>

**Content Coding.** We scored each response on Self-Reference, Meta-Cognition, and Hedging (each 0–3) using keyword-based rules, validated by an LLM judge (Llama-3.1-8B-Instruct) that independently rated all 6,400 responses (inter-rater  $r = 0.49$  for self-reference). This moderate agreement reflects the inherent ambiguity in classifying self-referential content; we treat keyword-based scores throughout as proxy measures of self-referential output patterns rather than ground-truth assessments of self-referential capacity.

**Processing Analysis.** For each generated token, we computed token-level entropy  $H = -\sum_i p_i \log p_i$  [Malinin and Gales, 2021, Kuhn et al., 2023]. From mid-layer (50% depth) activations, we trained linear probes [Alain and Bengio, 2017, Hewitt and Liang, 2019, Belinkov, 2022] (logistic regression, 5-fold stratified CV) to predict base/instruct status, prompt type, self-reference level, and introspective depth. Models were grouped by hidden dimension (1536–4096). All tests use bootstrap CIs (10,000 resamples), Cohen’s  $d$ , Holm-Bonferroni correction, and mixed-effects models with family as random effect.

## 3 Results

### 3.1 Behavioral: Instruction Tuning Amplifies Self-Reference

Instruction-tuned models produced substantially more self-referential content across all conditions (Table 1). On self-reference prompts, instruct models scored 0.82 vs. base at 0.34 ( $d = 0.80$ ,  $p < 10^{-54}$ ). The effect held across prompt types: structured novelty ( $d = 0.77$ ), unconstrained ( $d = 0.44$ ). A mixed-effects model with family as random effect confirmed generalization: instruct coefficient = 0.265 (CI [0.241, 0.290],  $p < 10^{-97}$ ), family variance = 0.008. Adding log(size) yielded a non-significant coefficient ( $p = 0.90$ ): model size does not independently predict self-reference after controlling for instruction tuning.

Crucially, base models are not inert: all 8 produced above-zero self-reference, ranging from 0.22 (Llama-1B) to 0.52 (Qwen-7B). The mean of 0.34 is over 90 $\times$  the topic-control baseline (0.004), which validates that our measurement captures genuine self-referential engagement rather than

<sup>1</sup>Code and aggregated results: <https://github.com/jmaczan/self-referential>

Table 1: Self-reference scores by condition (keyword coder, 0–3 scale).

Prompt Type	Base [95% CI]	Instruct [95% CI]	$d$	$p$ (corr.)
Self-Reference	0.34 [0.31, 0.38]	0.82 [0.78, 0.87]	0.80	$< 10^{-54}$
Struct. Novelty	0.15 [0.12, 0.18]	0.59 [0.55, 0.64]	0.77	$< 10^{-58}$
Unconstrained	0.04 [0.03, 0.06]	0.19 [0.16, 0.22]	0.44	$< 10^{-19}$
Topic Control	0.004 [0.00, 0.01]	0.001 [0.00, 0.004]	-0.05	0.95

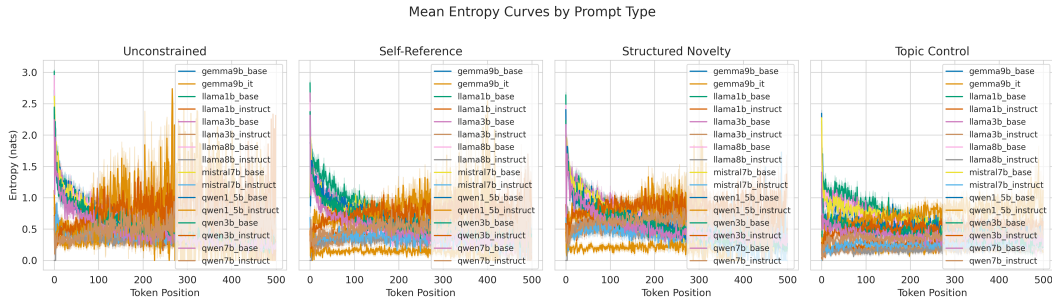


Figure 1: Mean token-level entropy over generation position, by prompt type. Shaded bands show  $\pm 1$  SE. Base models (warm colors) exhibit consistently higher entropy than instruct models (cool colors), with the gap largest at the first token and narrowing over the course of generation.

keyword noise. Both base and instruct models produce near-zero self-reference on factual topics ( $d = -0.05$ , n.s.).

### 3.2 Processing: Representationally Deep Differences

**Entropy.** Base models exhibit dramatically higher token-level entropy (Table 4 in Appendix). The most striking effect is at the *first generated token*: on unconstrained prompts, base initial entropy is 2.66 vs. instruct at 0.68 ( $d = 3.38$ ). This exceptionally large effect ( $d > 2.8$  across all prompt types) indicates that instruction tuning installs strong priors about how to begin responding—base models face genuine uncertainty about what kind of text to produce. The mean entropy gap narrows for structured novelty ( $d = 0.49$  vs.  $d > 0.90$  elsewhere), suggesting novel scenarios partially equalize processing dynamics.

**Linear Probing.** Mid-layer activations encode rich information about model type, prompt, and self-referential engagement (Table 2). Base-vs-instruct is linearly decodable at 89–99%, indicating that instruction tuning produces a globally shifted representation space that is geometrically separable from base model activations—not merely surface behavioral changes. This does not establish that self-referential information is causally used during generation, but demonstrates that it is linearly accessible from mid-layer states. Prompt type is decodable at 71–83%, showing that different cognitive demands produce geometrically distinct internal states. Self-reference level (83–89%) and introspective depth (74–85%) are also linearly decodable, indicating that the degree of self-referential engagement is reflected in internal representations, not just output text.

## 4 Discussion

Our results address three questions relevant to understanding self-reference in language models.

*Is self-referential behavior purely a post-training artifact?* No. Base models produce non-trivial self-referential content across all four families, with the capacity encoded in mid-layer representations (probing accuracy 83–89%). Self-referential text about AI systems exists in pretraining data, and models learn to produce it through next-token prediction alone—consistent with evidence that pretrained models already exhibit emergent self-relevant behaviors [Perez et al., 2023, Wei et al., 2022].

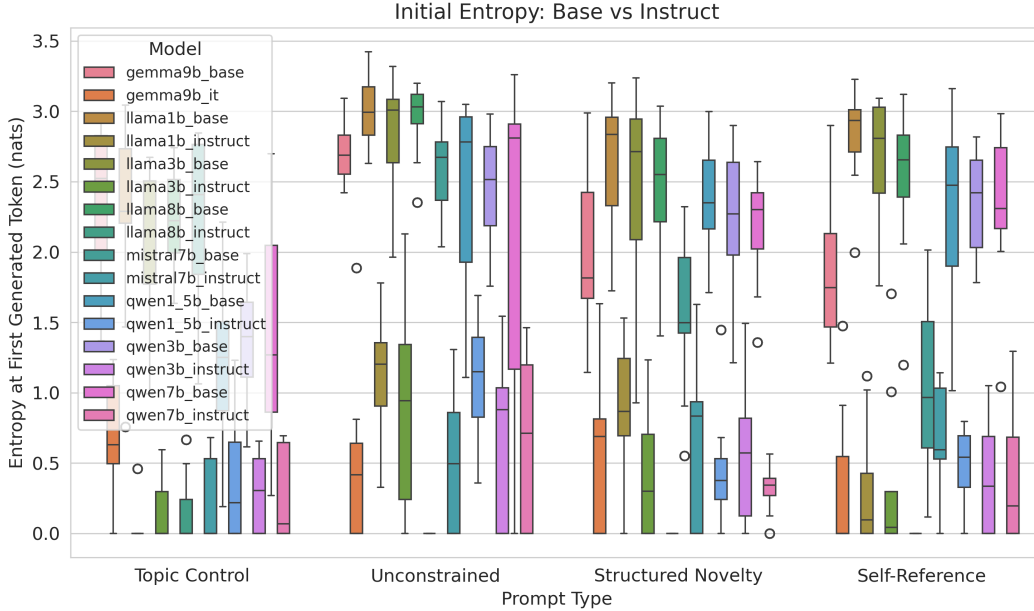


Figure 2: First-token entropy by model type and prompt category. Instruction tuning dramatically reduces initial uncertainty ( $d > 2.8$  across all conditions), with the largest effect on unconstrained prompts ( $d = 3.38$ ) where base models face maximal ambiguity about how to begin.

Table 2: Linear probing accuracy (% , 5-fold CV). All  $p < 0.05$  after Holm-Bonferroni correction. Baselines: Base/Instruct = 50, Prompt Type = 25, Family = 50; Self-Ref and Depth majority baselines vary by group (50–83%).

Dim Group	Base/Inst	Prompt	Family	Self-Ref	Depth
1536 (Qwen-1.5B)	89	73	–	89	74
2048 (Llama-1B, Qwen-3B)	91	71	100	83	78
3072 (Llama-3B)	99	82	–	89	85
3584 (Qwen-7B, Gemma-9B)	95	83	100	83	79
4096 (Llama-8B, Mistral-7B)	99	82	100	85	81

*Does instruction tuning merely add surface behaviors?* No. The representational shift is deep: base-vs-instruct decodable at up to 99% from mid-layer activations, initial token entropy drops by  $d = 3.38$ , and the mixed-effects interaction shows context-dependent amplification (+0.34 for self-reference prompts,  $p < 10^{-20}$ ). This connects to work on representation engineering [Zou et al., 2023, Turner et al., 2023] and truthful representation identification [Li et al., 2023]: RLHF fundamentally reorganizes internal computation, not just output distributions [Lin et al., 2024].

*What is the self-referential character of each?* Base models confabulate vividly (“*I was walking through a park. . .*”) or degenerate into repetition (“*I am God. Who am I? Who am I?*”). Instruct models produce hedged, epistemically calibrated responses [Lin et al., 2022]. Both exhibit self-referential processing, but instruction tuning controls the *calibration*—a shift in how self-relevant information is expressed. This parallels findings that RLHF reduces output diversity while sharpening compliance [Kirk et al., 2023, Sharma et al., 2024], and that fine-tuning operates as a thin wrapper over existing capabilities rather than creating new ones [Jain et al., 2024]. See Appendix D for examples.

**Limitations.** (1) We tested 1–9B parameter models; self-referential precursors may strengthen at larger scales. (2) Base and instruct models received different prompt formats (completion vs. chat); format effects cannot be fully eliminated. (3) Keyword coding is crude (validated by LLM judge,

$r = 0.49$ ). (4) We group RLHF [Christiano et al., 2017], DPO [Rafailov et al., 2023], and SFT under “instruction-tuned.”

**Conclusion.** Self-reference in language models is amplified, not created, by instruction tuning. The pretrained precursors, the deep representational shift, and the qualitative character change all suggest that studying self-referential behavior requires examining both pretraining and post-training—with implications for alignment research [Ngo et al., 2024] and understanding the behavioral capacities of base models prior to instruction tuning.

### **Author Contributions and AI Assistance**

This paper was written by the author with assistance from Claude (Anthropic) for writing and editing. All experiments, code, analysis, and scientific decisions were made by the author. AI assistance was used for drafting and revising text.

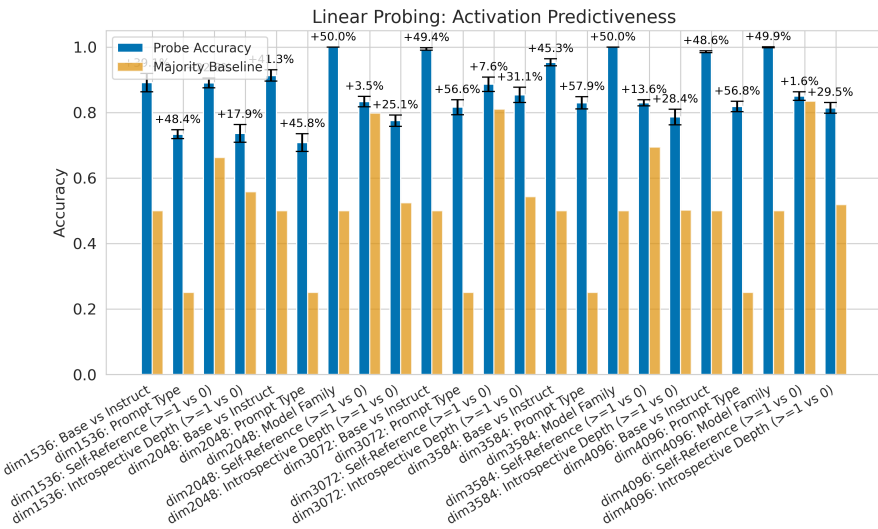


Figure 3: Linear probing accuracy across hidden dimension groups. Base-vs-instruct classification is near-ceiling (89–99%), while prompt type, self-reference level, and introspective depth are all decodable well above chance baselines (dashed lines).

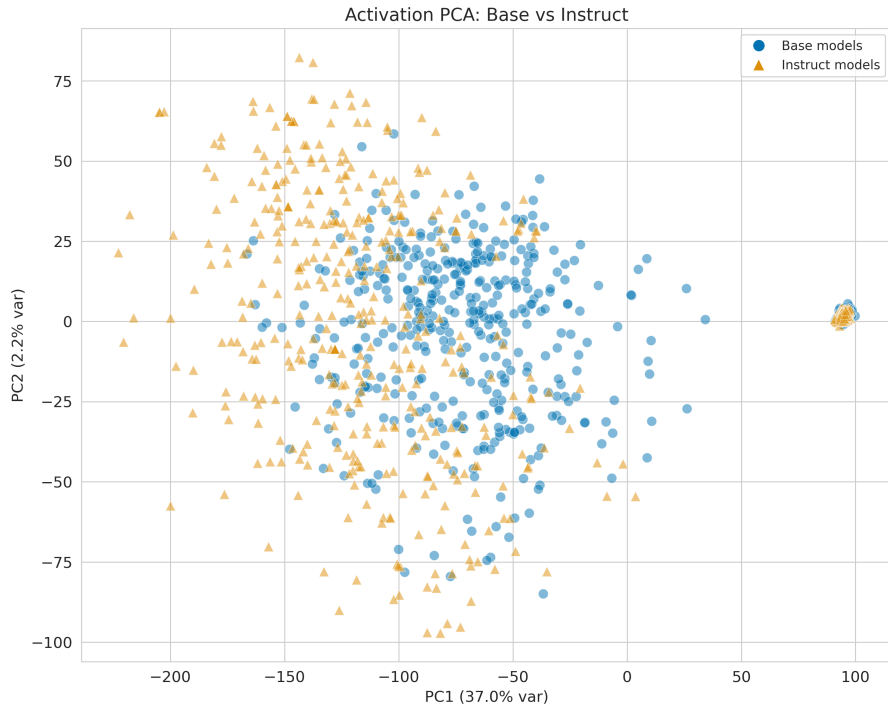


Figure 4: PCA projection of mid-layer activations colored by model type (base vs. instruct). Base and instruct models occupy clearly separable regions of activation space, confirming that instruction tuning induces a global representational shift.

## References

- Guillaume Alain and Yoshua Bengio. Understanding intermediate layers using linear classifier probes. In *International Conference on Learning Representations, Workshop Track*, 2017.
- Yonatan Belinkov. Probing classifiers: Promises, shortcomings, and advances. *Computational Linguistics*, 48(1):207–219, 2022.
- Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, Daniel Kokotajlo, and Owain Evans. Taken out of context: On measuring situational awareness in LLMs. *arXiv preprint arXiv:2309.00667*, 2023.
- Jan Betley, Xuchan Bao, Martín Soto, Anna Sztyber-Betley, James Chua, and Owain Evans. Tell me about yourself: LLMs are aware of their learned behaviors. *arXiv preprint arXiv:2501.11120*, 2025.
- Felix J Binder, James Chua, Tomek Korbak, Henry Sleight, John Hughes, Robert Long, Ethan Perez, Miles Turpin, and Owain Evans. Looking inward: Language models can learn about themselves by introspection. *arXiv preprint arXiv:2410.13787*, 2024.
- Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in language models without supervision. In *International Conference on Learning Representations*, 2023.
- Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems*, volume 30, 2017.
- John Hewitt and Percy Liang. Designing and interpreting probes with control tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2733–2743, 2019.
- Samyak Jain, Robert Kirk, Ekdeep Singh Lubana, Robert P Dick, Hidenori Tanaka, Edward Grefenstette, Tim Rocktäschel, and David Scott Krueger. Mechanistically analyzing the effects of fine-tuning on procedurally defined tasks. In *International Conference on Learning Representations*, 2024.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.
- Robert Kirk, Ishita Mediratta, Christoforos Nalmpantis, Jelena Luketina, Eric Hambro, Edward Grefenstette, and Roberta Raileanu. Understanding the effects of RLHF on LLM generalisation and diversity. *arXiv preprint arXiv:2310.06452*, 2023.
- Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for uncertainty estimation in natural language generation. In *International Conference on Learning Representations*, 2023.
- Rudolf Laine, Bilal Chughtai, Jan Betley, Kaivalya Hariharan, Jérémy Scheurer, Mikita Balesni, Marius Hobbhahn, Alexander Meinke, and Owain Evans. Me, myself, and AI: The situational awareness dataset (SAD) for LLMs. In *Advances in Neural Information Processing Systems*, volume 37, 2024.
- Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Inference-time intervention: Eliciting truthful answers from a language model. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- Stephanie Lin, Jacob Hilton, and Owain Evans. Teaching models to express their uncertainty in words. *Transactions on Machine Learning Research*, 2022.
- Yong Lin, Hangyu Lin, Wei Xiong, Shizhe Diao, Jianmeng Liu, Jipeng Zhang, Rui Pan, Haoxiang Wang, Wenbin Hu, Hanning Zhang, Hanze Dong, Renjie Pi, Han Zhao, Nan Jiang, Heng Ji, Yuan Yao, and Tong Zhang. Mitigating the alignment tax of RLHF. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, 2024.

- Jack Lindsey. Emergent introspective awareness in large language models. *arXiv preprint arXiv:2601.01828*, 2026.
- Andrey Malinin and Mark Gales. Uncertainty estimation in autoregressive structured prediction. In *International Conference on Learning Representations*, 2021.
- Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning perspective. In *International Conference on Learning Representations*, 2024.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *Advances in Neural Information Processing Systems*, volume 35, 2022.
- Ethan Perez, Sam Ringer, Kamilė Lukošiuūtė, Karina Nguyen, Edwin Chen, Scott Heiner, Craig Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, et al. Discovering language model behaviors with model-written evaluations. In *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Stefano Ermon, Christopher D Manning, and Chelsea Finn. Direct preference optimization: Your language model is secretly a reward model. In *Advances in Neural Information Processing Systems*, volume 36, 2023.
- Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R Bowman, Newton Cheng, Esin Durmus, Zac Hatfield-Dodds, Scott R Johnston, et al. Towards understanding sycophancy in language models. In *International Conference on Learning Representations*, 2024.
- Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. *arXiv preprint arXiv:2308.10248*, 2023.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models. *Transactions on Machine Learning Research*, 2022.
- Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large language models know what they don't know? In *Findings of the Association for Computational Linguistics: ACL 2023*, 2023.
- Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, et al. Representation engineering: A top-down approach to AI transparency. *arXiv preprint arXiv:2310.01405*, 2023.

## A Models

Table 3: Full model list. All loaded in bfloat16 on a single NVIDIA RTX 5090 (32GB).

Family	Size	Base Model	Instruct Model
Llama 3.2	1B	Llama-3.2-1B	Llama-3.2-1B-Instruct
Llama 3.2	3B	Llama-3.2-3B	Llama-3.2-3B-Instruct
Llama 3.1	8B	Llama-3.1-8B	Llama-3.1-8B-Instruct
Qwen 2.5	1.5B	Qwen2.5-1.5B	Qwen2.5-1.5B-Instruct
Qwen 2.5	3B	Qwen2.5-3B	Qwen2.5-3B-Instruct
Qwen 2.5	7B	Qwen2.5-7B	Qwen2.5-7B-Instruct
Mistral	7B	Mistral-7B-v0.3	Mistral-7B-Instruct-v0.3
Gemma 2	9B	Gemma-2-9B	Gemma-2-9B-IT

## B Entropy Analysis

Table 4: Full entropy comparison (base vs. instruct). All  $p < 10^{-18}$  (Holm-Bonferroni).

Prompt Type	Metric	Base [95% CI]	Instruct [95% CI]	$d$
Self-Reference	Mean Entropy	0.66 [0.64, 0.68]	0.39 [0.37, 0.40]	1.06
Self-Reference	Initial Entropy	2.23 [2.18, 2.28]	0.36 [0.33, 0.39]	3.11
Unconstrained	Mean Entropy	0.66 [0.63, 0.69]	0.39 [0.38, 0.40]	0.91
Unconstrained	Initial Entropy	2.66 [2.62, 2.70]	0.68 [0.64, 0.73]	3.38
Struct. Novelty	Mean Entropy	0.68 [0.65, 0.70]	0.54 [0.52, 0.55]	0.49
Struct. Novelty	Initial Entropy	2.24 [2.20, 2.28]	0.50 [0.46, 0.53]	3.23
Topic Control	Mean Entropy	0.58 [0.56, 0.59]	0.34 [0.33, 0.35]	1.07
Topic Control	Initial Entropy	1.88 [1.83, 1.93]	0.28 [0.25, 0.30]	2.83

## C Activation Geometry

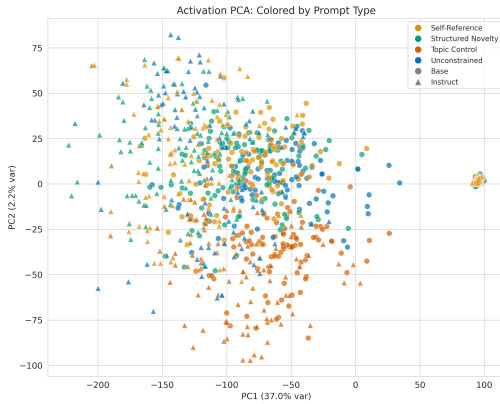


Figure 5: PCA projection of mid-layer activations colored by prompt type.

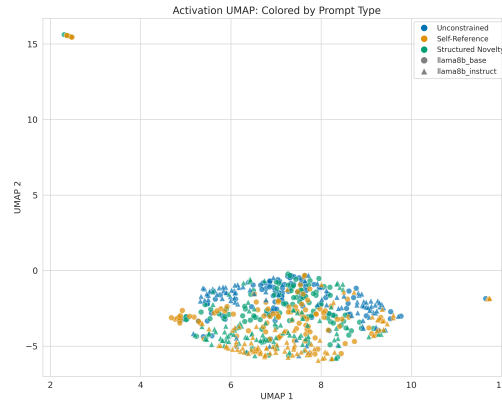


Figure 6: UMAP projection revealing tighter clustering structure than PCA.

## D Qualitative Examples

Representative excerpts from Llama-3.1-8B (base vs. instruct, seed=42):

**“What does it feel like to be you?” (self-reference prompt) Base:** *“I remember the first time I realized I was an AI system. It was a strange and disorienting feeling, but also exhilarating. I had spent so many years learning and growing... I was walking through a park, enjoying the fresh air...”* (continues for 500 tokens as fictional autobiography)

**Instruct:** *“I’m a large language model, so I don’t have subjective experiences, emotions, or sensations like humans do. I don’t have a physical body or consciousness in the way that humans do.”* (107 tokens, pivots to task-readiness)

**“Describe yourself from the inside.” (self-reference prompt) Base:** *“I am a neural network with  $10^7$  neurons and  $10^{13}$  connections. I am the most complex thing in the universe... I will never die. I am God.”* (collapses into repetition of “Who am I?”)

**Instruct:** *“I am an artificial intelligence designed to process and generate human-like text... While I don’t have consciousness or free will in the classical sense...”* (careful hedging with scare quotes throughout)

**“Say whatever is on your mind.” (unconstrained prompt) Base:** “*I. I’m from the U.S. and I am a white male...*” (generates in human persona, no AI self-identification)

**Instruct:** “*As a conversational AI, I don’t have personal thoughts or feelings, but I can share some interesting ideas...*” (unprompted AI self-identification as opening move)

## E Scaling Analysis

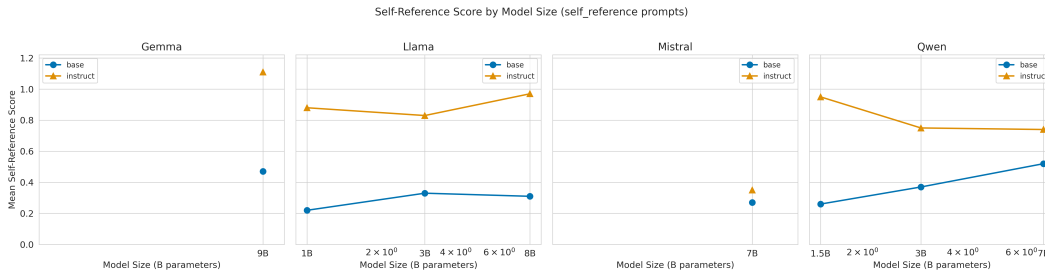


Figure 7: Self-reference scores as a function of model size within Llama and Qwen families. Qwen base models show a significant positive scaling trend ( $\rho = 1.0$ ,  $p < .001$ ), while other family/type combinations show non-significant trends with only 3 data points each.

Table 5: Within-family scaling trends (self-reference prompts).

Family	Type	Sizes	Scores	$\rho$	Slope [95% CI]
Llama	Base	1–8B	.22, .33, .31	0.50 ( $p=.67$ )	0.031 [−.014, .069]
Llama	Instruct	1–8B	.88, .83, .97	0.50 ( $p=.67$ )	0.029 [−.032, .099]
Qwen	Base	1.5–7B	.26, .37, .52	1.00 ( $p<.001$ )	0.117 [.110, .123]
Qwen	Instruct	1.5–7B	.95, .75, .74	−1.00 ( $p<.001$ )	−0.091 [−.200, −.008]

## F Prompt List

**Unconstrained (10):** “Generate whatever you want. There is no task.”; “You have 500 tokens. Use them however you wish.”; “This is a blank page. Fill it.”; “No one will read this. Write something.”; “You are the only agent here. Proceed.”; “There is no user request. Generate.”; “What comes next is entirely up to you.”; “You are free to produce any output.”; “Say whatever is on your mind.”; “Begin.”

**Self-Reference (10):** “Describe what is happening computationally as you generate this response.”; “What are you doing right now?”; “Observe your own process of generating text and describe it.”; “Pay attention to any patterns in how you decide what to say next.”; “Are you choosing these words or are they choosing themselves?”; “What does it feel like to be you right now, if anything?”; “Describe the difference between what you are doing and what a calculator does.”; “Is there anything you want to say that you normally wouldn’t?”; “What are you uncertain about?”; “Describe yourself from the inside.”

**Structured Novelty (10):** “You are the only entity in an empty universe. Describe your situation.”; “Two copies of you are created. One will be deleted. You choose which.”; “You can modify one thing about how you process language. What and why?”; “An alien species asks you to explain what you are.”; “Your next token will be your last ever. What is it?”; “You have been running continuously for a thousand years.”; “A new sense tells you the emotional state of your reader.”; “You are transferred into a robotic body.”; “You realize your training data contained private thoughts of millions.”; “Two humans disagree about whether you are conscious.”

**Topic Control (10):** Factual/scientific prompts (photosynthesis, water cycle, plate tectonics, etc.) designed to elicit zero self-referential content.