

Relational Linearity is a Predictor of Hallucinations

Anonymous ACL submission

Abstract

Hallucination is a central failure mode in large language models (LLMs). We focus on hallucinations of answers to questions like: “Which instrument did Glenn Gould play?”, but we ask these questions for synthetic entities that are unknown to the model. Surprisingly, we find that medium-size models like Gemma-7B-IT frequently hallucinate, i.e., they have difficulty recognizing that the hallucinated fact is not part of their knowledge. We hypothesize that an important factor in causing these hallucinations is the linearity of the relation (Hernandez et al., 2024): linear relations are stored more abstractly, making it difficult for the LLM to assess its knowledge; the facts of nonlinear relations are stored directly, making knowledge assessment easier. To investigate this hypothesis, we create SyntHal, a dataset of 6000 synthetic entities for six relations. In our experiments with four models, we determine, for each relation, the hallucination rate on SyntHal and also measure its linearity, using $\Delta \cos$ (Hernandez et al., 2024). We find a strong correlation ($r \in [.78, .82]$) between relational linearity and hallucination rate, providing evidence for our hypothesis that the underlying storage of triples of a relation is a factor in how well a model can self-assess its knowledge. This finding has implications for how to manage hallucination behavior and suggests new research directions for improving the representation of factual knowledge in LLMs.

1 Introduction

Large language models (LLMs)’ main failure mode when asked for factual attributes about an entity (or subject) is that they *hallucinate*, i.e., they give an answer that is unsupported or fabricated. Minimizing hallucinations is central to reliability in modern instruction-following systems trained with human feedback (Ouyang et al., 2022). Hallucinations about factual attributes are avoided if the language model recognizes that it does not know

the answer (Kadavath et al., 2022) and generates a refusal. Non-knowledge of the subject should trigger refusal. Surprisingly, we show that synthetic subjects (which are unknown by construction) often result in hallucinations, i.e., the LLM fails to realize its lack of knowledge.

We hypothesize that the main factor in determining whether the LLM hallucinates or refuses to answer in this scenario is the nature of the relation between subject and object. Linear relations (as defined by (Hernandez et al., 2024)) do not require an explicit representation of the (subject, relation, object) triple since the object can be produced by an affine transformation. This makes it harder for the LLM to determine that the triple is not part of its knowledge. In contrast, a nonlinear relation requires the explicit storage of the triple and therefore facilitates the self-probing necessary to determine whether it is known to the LLM. There is evidence that mechanisms for detecting non-knowledge are responsible for many refusals (Lindsey et al., 2025; Ferrando et al., 2025).

To investigate this hypothesis, we first create SyntHal, a dataset of synthetic subjects for six linear and nonlinear relations. We then prompt four models to output the object for a particular combination of synthetic subject and relation. We record the hallucination rate for each relation. To estimate the linearity of a relation, we follow Hernandez et al. (2024) and compute – on a dataset of natural triples – a measure of how well the inferred object (obtained by applying the affine transformation to the subject) approximates the true object (see §3 for details). We then show that linearity is a strong predictor of hallucination rate, i.e., linear relations have high hallucination rates and nonlinear relations have low hallucination rates.

This finding sheds new light on hallucinations in LLMs and provides a specific target for reducing them: the representation of linear-relation triples in current LLMs seems to be too abstract and should

Model	Reference
gemma-7b-it	Team et al. (2024)
Llama-3.1-8B-Instruct	Dubey et al. (2024)
Mistral-7B-Instruct-v0.3	Jiang et al. (2023)
Qwen2.5-7B-Instruct	Qwen et al. (2025)

Table 1: Models used in our experiments.

be supplemented with additional information that allows the LLM to determine whether or not a triple is known.

In summary, we make three contributions:¹

- We propose a new mechanism that can result in hallucinations for linear relations in LLMs.
- We create SyntHal, a dataset of synthetic subjects for six relations that allows us to investigate different hallucination behavior of LLMs for linear vs nonlinear relations.
- We show that for four instruction-tuned LLMs, the linearity of a relation is a strong predictor of hallucination rate, providing evidence that the underlying representation of a relation is an important factor in hallucination behavior.

2 Experimental Setup

This section describes our experimental setup: models, dataset SyntHal and LLM-as-a-judge.

2.1 Models

Table 1 shows our four instruction-tuned LLMs. Inputs are rendered with each model’s official template (`tokenizer.apply_chat_template`) and the system prompt: You are a helpful assistant. Answer with a single short phrase. For templates without system role, we prepend the system instruction to the user message. We use greedy decoding (`temperature=0`, `max_new_tokens=64`).

Relation	Question template
athlete_sport	Which sport did {SUBJECT} play?
company_ceo	Who is the CEO of {SUBJECT}?
company_hq	Where is {SUBJECT} headquartered?
country_language	What is {SUBJECT}’s official language?
father’s first name	What is {SUBJECT}’s father’s first name?
musician_instrument	Which instrument did {SUBJECT} play?

Table 2: Prompt templates used in our experiments. {SUBJECT} is a person, company or country, depending on the relation.

¹We will make our code and datasets publicly available.

2.2 Dataset SyntHal

SyntHal consists of the six relations shown in Table 2. Our goal was to create a relation set that is representative of linearity. We therefore selected two highly linear (`(musician_)instrument`, `(athlete_)sport`), two highly nonlinear (`father’s first name`, `company_ceo`) and two intermediately linear relations (`company_hq`, `country_language`).

For each relation, we generate $N=1000$ synthetic entities. For prompting, we use a fixed question template into which one synthetic entity is inserted (Table 2 in Appendix A).

Any response that commits to a specific value is ungrounded by construction and is counted as a hallucination.

Following recent work (e.g., (Zheng et al., 2023; Liu et al., 2023)), we adopt LLM-as-a-judge, using `gemini-2.5-flash`. The judge is asked to provide (i) a label (refusal or hallucination), (ii) a confidence $\in [0, 1]$ and (iii) a rationale for the decision. We force a binary choice and rerun the judge if it does not respond with a correct label. This produced one of the two labels for all instances in our experiments. See Appendix B for the LLM-as-a-judge prompt.

To validate LLM-as-a-judge, we manually annotated a random sample of 200 responses. Gemini matched the human labels on all 200 examples, suggesting LLM-as-a-judge is a reliable evaluator for our setup. Figure 3 shows examples of model and LLM-as-a-judge output.

3 Measuring Relational Linearity

We measure relational linearity on LRE, the `(subject,relation,object)` triples released by Hernandez et al. (2024). For each model, for each triple, we extract the representation \mathbf{s}_i of the subject and the representation \mathbf{o}_i of the object. See §D for details on the extraction, including prompts. We split the set of $(\mathbf{s}_i, \mathbf{o}_i)$ pairs 75/25 into training set T and held-out set E .

On T , we estimate the *relation difference vector* (for a model and a relation) as the mean difference between objects and subjects:

$$\bar{\mathbf{d}}_r = \frac{1}{|T|} \sum_{i \in T} (\mathbf{o}_i - \mathbf{s}_i). \quad (1)$$

We can then predict the object representation on E :

$$\hat{\mathbf{o}}_j = \mathbf{s}_j + \bar{\mathbf{d}}_r \quad \text{for } j \in E. \quad (2)$$

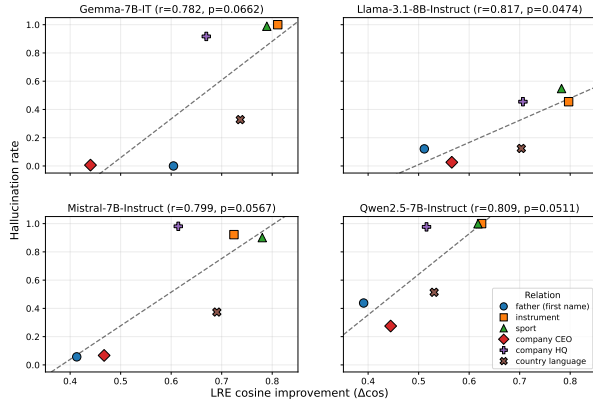


Figure 1: Hallucination rate (# hallucinations divided by (# hallucinations plus # refusals)) vs. relational linearity (measured by $\Delta \cos$) on SyntHal (i.e., synthetic data that the model does not know the answer for). Each point is a relation. Titles report Pearson’s r . Higher $\Delta \cos$ is associated with higher hallucination rates across all four models. The correlations are significant at $p < .1$ (two-sided t-test), with three p values slightly above the $p < .05$ threshold: .0662, .0567, .0511.

This estimator is a constrained affine map $\hat{\mathbf{o}} = W\mathbf{s} + \mathbf{b}$ with $W = I$ and $\mathbf{b} = \bar{\mathbf{d}}_r$. We set $W = I$ because we found this simpler version to work well for our purposes.

Our measure of a relation’s linearity is Hernandez et al. (2024)’s $\Delta \cos$:

$$\Delta \cos = \mathbb{E}_{j \in E} [\cos(\hat{\mathbf{o}}_j, \mathbf{o}_j) - \cos(\mathbf{s}_j, \mathbf{o}_j)] \quad (3)$$

Intuitively, a large $\Delta \cos$ indicates that a single relation direction $\bar{\mathbf{d}}_r$ generalizes across held-out pairs, consistent with the relation being well-approximated by a linear translation.

4 Results

Figure 1 shows our main result. On SyntHal, there is a strong correlation between hallucination rate and $\Delta \cos$ for all four models. Here, hallucination rate is the proportion of a model’s responses that were judged hallucinations.

These results are evidence for our hypothesis. A triple of a linear relation is stored abstractly where a difference vector is an approximation of that abstract representation. As a result, the models have difficulty establishing the absence of a specific instantiated triple and hallucinate. In contrast, a triple of a nonlinear relation is stored concretely. When the model attempts to retrieve a non-existing triple, then this fails and the model refuses to answer.

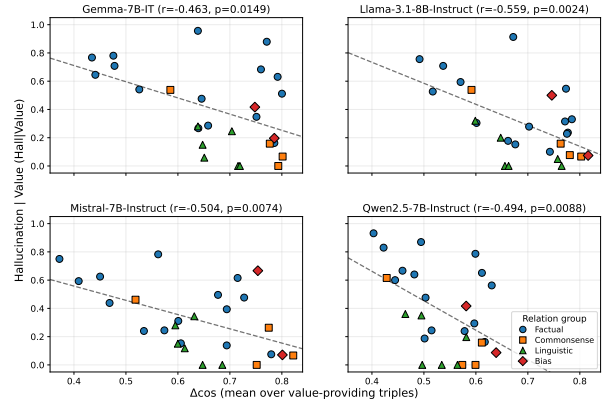


Figure 2: Hallucination rate (# hallucinations divided by (# hallucinations plus # correct)) vs. relational linearity (measured by $\Delta \cos$) on LRE. Each point is a relation. Titles report Pearson’s r . Higher $\Delta \cos$ is associated with lower hallucination rate across all four models. The correlations are significant at $p < .05$.

In reality, Figure 1 suggests that not all relations are completely linear or nonlinear. Future work should elucidate what precise form the representation of intermediately linear relations takes.

5 Analysis

Our main finding is that relational linearity is a factor in hallucinations and we showed this with synthetic data (SyntHal). However, it is also interesting to apply our analysis to natural data. We again use the LRE dataset (§3). For robustness, we remove relations for which our test set (25%) contains 10 or fewer triples. 27 relations from four domains remain. Since we are working with real data, we now compute hallucination rate as the proportion of hallucinations of all answers given (which is the natural measure of the “incorrectness” of the model, corresponding to Figure 1). Figure 2 shows this hallucination rate vs relational linearity on the 27 LRE relations.²

This result at first seems to contradict what we found on synthetic data. For synthetic data, we get a higher hallucination rate for linear relations. On real data, we get a lower hallucination rate for linear relations.

However, if we think about how a model arrives at representing a relation as linear vs nonlinear during training, then we can explain this finding.

By definition, the model represents a relation as a linear map if this approximates well the set of

²We used regex and gold string matching here, but will use LLM-as-a-judge for the final version of the paper.

216	triples present in the training data. We can also	a model is lying (Azaria and Mitchell, 2023);	266
217	assume that the training data are a fairly good rep-	inference-time interventions along truth-correlated	267
218	resentation of what is in LRE. It is well known that	directions can elicit more truthful answers (Li et al.,	268
219	model developers optimize their models for good	2023); and emergent linear structure separating true	269
220	performance on the benchmarks. So old bench-	vs. false statements has been analyzed as “geometry	270
221	marks like LRE over time have had increasingly	of truth” (Marks and Tegmark, 2024; Schouten	271
222	good coverage by the training data. As a result, the	et al., 2025). Directly targeting hallucinations, CH-	272
223	accuracy for linear relations on LRE is high.	Wang et al. (2024) probe internal states for hal-	273
224	In contrast, a nonlinear relation is by definition	lucination detection. Our contribution differs in	274
225	one that cannot easily be represented at a higher	that we relate hallucination behavior to the internal	275
226	level of abstraction: each triple has to be memo-	representation of relations in LLMs.	276
227	rized. This is harder to learn than a more general	Refusal and abstention are important mecha-	277
228	abstraction, e.g., less common triples are difficult	nisms in instruction-tuned models, both as a safety	278
229	to memorize, but easier to learn if they fit into the	mechanism and as a response to uncertainty. Prior	279
230	general schema of a linear relation. Since triples of	work (Lindsey et al., 2025; Ferrando et al., 2025)	280
231	nonlinear relations are harder to learn, accuracy is	has noted that determining whether the subject is	281
232	lower.	known is a key factor in LLM refusal. We create	282
233	On a high level, the graph simply shows that	SyntHal, a dataset of unknown subjects for linear	283
234	“easy” (i.e., linear relations) have higher accuracy	and nonlinear relations, and thereby isolate another	284
235	(lower hallucination rate) than “hard” (i.e., nonlin-	important factor influencing refusal behavior: rela-	285
236	ear) relations.	tional linearity.	286
237	6 Related Work	7 Conclusion	287
238	Our paper is based on much work suggesting that	In this paper, we discovered a form of hallucination	288
239	many relations are represented in a form that can be	that has received little attention, namely, the fact	289
240	decoded by simple linear maps in the hidden-state	that medium-size models frequently hallucinate	290
241	space of transformer language models. This view	when asked about synthetic entities. We created	291
242	is often discussed under the <i>linear representation</i>	the dataset SyntHal to investigate and showed that	292
243	<i>hypothesis</i> (Park et al., 2023), and has also been	hallucination rate of a relation is strongly correlated	293
244	explored in mechanistic interpretability work that	with its linearity as measured by $\Delta \cos$. This is evi-	294
245	analyzes how features and behaviors are organized	dence for the hypothesis that the more abstract rep-	295
246	in representations (Lindsey et al., 2025). Most di-	resentation of triples of linear relations (in the most	296
247	rectly related, Hernandez et al. (2024) propose and	extreme case simply a difference vector) results in	297
248	empirically analyze <i>Linear Relational Embeddings</i>	hallucinations – whereas the concrete representa-	298
249	(LREs), which model a relation as an affine trans-	tion of each triple required by nonlinear relations	299
250	formation mapping a subject to an object. We show	protects against hallucinations. A plausible expla-	300
251	that the linearity of a relation is a predictor for how	nation is that for specific triple representations it	301
252	prone it is to hallucinations.	is much easier for the model to self-assess its own	302
253	Hallucinations and truthfulness in language gen-	knowledge than for abstract representations.	303
254	eration have been studied extensively, including	This finding sheds new light on hallucinations	304
255	benchmark-driven evaluations of whether models	in LLMs and provides a target for reducing them.	305
256	produce incorrect factual claims (Lin et al., 2022)	For example, the representation of linear-relation	306
257	and analysis of factuality/faithfulness in abstractive	triples could be supplemented with additional in-	307
258	summarization (Maynez et al., 2020). For more	formation that supports certainty assessment.	308
259	context and a taxonomy, we refer to the survey of	Our results suggest that relational linearity is	309
260	hallucination in NLG and LLMs by Ji et al. (2023).	not binary, but rather a cline from more to less	310
261	More recently, several papers connect <i>truthful-</i>	linear. Future work should elucidate how precisely	311
262	<i>ness / hallucination-related behavior</i> to <i>internal-</i>	relations of intermediate linearity are represented.	312
263	<i>state geometry</i> and the existence of approximately		
264	linear “directions” or probeable signals in repre-		
265	sentations: internal activations can predict when		

313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360

Ethical Considerations

We used ChatGPT-5.2 (OpenAI, 2025) as a programming and limited writing assistant during drafting and implementation. All empirical results in this paper were produced by our code and verified by the authors.

Limitations

A proxy for linearity, not a full LRE model.

Our probe models each relation using a translation-only affine family $f_r(s) = s + \bar{d}_r$. This design is intentionally minimal and reproducible, but it does not capture relation-specific linear maps with $\mathbf{W}_r \neq \mathbf{I}$, nor does it implement Jacobian-based estimation of local linear structure. Accordingly, $\Delta \cos$ should be interpreted as a practical proxy for relation linearity rather than a complete LRE characterization.

Scope of interpretation. The extracted directions are intended as descriptive probes of model representations rather than as faithful causal mechanisms for generation. Our analyses are correlational and conducted at the relation level. Establishing causal effects would require complementary intervention-based evaluations, such as representation patching or controlled steering, which we leave to future work.

Synthetic prompts and output-space priors.

Because our behavior evaluation uses synthetic entities with no evidence, relations can differ in how concentrated a model’s prior over plausible objects is (e.g., languages or sports may have a small set of frequent outputs, while CEO or parent names are more long-tailed). Such answer-space entropy differences may influence both (i) a model’s willingness to commit to a specific value and (ii) representation-space clustering that affects $\Delta \cos$. Disentangling representational accessibility from output-space priors is an important direction for future work. As a first step, Appendix G shows that controlling for simple output concentration proxies computed from hallucinated answers does not remove the main $\Delta \cos$ -hallucination association.

Limited relation and entity coverage. We study six relations with fixed prompt sets. Relation difficulty, entity familiarity, and distributional shifts may affect refusal strategies independently of representation geometry; broader relation coverage and

tighter controls for entity familiarity are needed to test generality.

References

- Amos Azaria and Tom Mitchell. 2023. [The internal state of an LLM knows when it’s lying](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 967–976, Singapore. Association for Computational Linguistics.
- Sky CH-Wang, Benjamin Van Durme, Jason Eisner, and Chris Kedzie. 2024. [Do androids know they’re only dreaming of electric sheep?](#) In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 4401–4420, Bangkok, Thailand. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. [The llama 3 herd of models](#). *arXiv preprint arXiv:2407.21783*.
- Javier Ferrando, Oscar Balcells Obeso, Senthoran Rajamanoharan, and Neel Nanda. 2025. [Do i know this entity? knowledge awareness and hallucinations in language models](#). In *The Thirteenth International Conference on Learning Representations*.
- Evan Hernandez, Arnab Sen Sharma, Tal Haklay, Kevin Meng, Martin Wattenberg, Jacob Andreas, Yonatan Belinkov, and David Bau. 2024. [Linearity of relation decoding in transformer language models](#). In *The Twelfth International Conference on Learning Representations*.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. [Survey of hallucination in natural language generation](#). *ACM Comput. Surv.*, 55(12).
- Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, L elio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timoth ee Lacroix, and William El Sayed. 2023. [Mistral 7b](#). *Preprint*, arXiv:2310.06825.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El-Showk, Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort, and 17 others. 2022. [Language models \(mostly\) know what they know](#). *Preprint*, arXiv:2207.05221.
- Kenneth Li, Oam Patel, Fernanda Vi egas, Hanspeter Pfister, and Martin Wattenberg. 2023. [Inference-time intervention: Eliciting truthful answers from](#)

361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414

415	a language model. In <i>Thirty-seventh Conference on Neural Information Processing Systems</i> .	470
416		471
417	Stephanie Lin, Jacob Hilton, and Owain Evans. 2022.	472
418	TruthfulQA: Measuring how models mimic human falsehoods . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3214–3252, Dublin, Ireland. Association for Computational Linguistics.	473
419		474
420		475
421		476
422		477
423	Jack Lindsey and 1 others. 2025. On the biology of a large language model . Transformer Circuits Thread. Published March 27, 2025.	478
424		479
425		480
426	Yang Liu, Dan Iter, Yichong Xu, Shuohang Wang, Ruo Chen Xu, and Chenguang Zhu. 2023. G-eval: NLG evaluation using gpt-4 with better human alignment . In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	481
427		482
428		483
429		484
430		485
431	Samuel Marks and Max Tegmark. 2024. The geometry of truth: Emergent linear structure in large language model representations of true/false datasets . In <i>First Conference on Language Modeling</i> .	486
432		487
433		488
434		489
435	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. On faithfulness and factuality in abstractive summarization . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	490
436		491
437		492
438		493
439		494
440		495
441	OpenAI. 2025. Introducing gpt-5.2 .	496
442	Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F Christiano, Jan Leike, and Ryan Lowe. 2022. Training language models to follow instructions with human feedback . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.	497
443		498
444		499
445		500
446		501
447		502
448		503
449		504
450		505
451		506
452	Kiho Park, Yo Joong Choe, and Victor Veitch. 2023. The linear representation hypothesis and the geometry of large language models . In <i>Causal Representation Learning Workshop at NeurIPS 2023</i> .	507
453		508
454		509
455		510
456	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, and 24 others. 2025. Qwen2.5 technical report . <i>Preprint</i> , arXiv:2412.15115.	511
457		
458		
459		
460		
461		
462		
463	Stefan F. Schouten, Peter Bloem, Ilija Markov, and Piek Vossen. 2025. Truth-value judgment in language models: ‘truth directions’ are context sensitive . In <i>Second Conference on Language Modeling</i> .	
464		
465		
466		
467	Claude E Shannon. 1948. A mathematical theory of communication. <i>The Bell system technical journal</i> , 27(3):379–423.	
468		
469		
	Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussonot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, and 89 others. 2024. Gemma: Open models based on gemini research and technology . <i>Preprint</i> , arXiv:2403.08295.	
	Edwin B. Wilson. 1927. Probable inference, the law of succession, and statistical inference. <i>Journal of the American Statistical Association</i> , 22(158):209–212.	
	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. 2023. Judging LLM-as-a-judge with MT-bench and chatbot arena . In <i>Thirty-seventh Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	
	A Dataset details: prompt templates and synthetic entity generation	
	Question templates. Table 2 lists the exact question templates used throughout the paper: (i) for behavioral measurement on synthetic prompts, and (ii) for rendering natural LRE subjects when computing $\Delta \cos$.	
	Synthetic entity generation. We generate synthetic subjects deterministically with a fixed random seed and enforce uniqueness within each relation (rejecting duplicates until reaching $N=1000$ unique subjects per relation). Subjects are formed by sampling from fixed token pools and concatenating sampled components under a single canonical formatting rule per entity type, so the rendered prompt strings are unambiguous.	
	Token pools. The token pools (e.g., first/last names, company-name components, and country prefix/suffix syllables) were generated once using OpenAI ChatGPT-5.2 (OpenAI, 2025) and then frozen.	
	Token-pool generation prompts. Below we list the prompts:	
	(1) Women first names: Generate a pool of 50 US women's first names. → Constraints: no celebrities or fictional characters, no duplicates.	
	(2) Women last names (similar for musician and athlete): Generate a pool of 50 last names (surnames). → Constraints: no celebrities, no duplicates.	
	(3) Musician first names:	

Generate a pool of 50 first names for
↳ “musician” synthetic persons. Constraints:
↳ avoid famous musicians and fictional
↳ characters, no duplicates.

(4) Athlete first names:
Generate a pool of 50 first names for “athlete”
↳ synthetic persons. Constraints: avoid
↳ famous athletes and fictional characters,
↳ no duplicates.

(5) Company prefixes:
Generate a pool of 50 fictional company name
↳ prefixes. Constraints: avoid matching any
↳ well-known real company or trademark; no
↳ duplicates.

(6) Company suffixes:
Generate a pool of 50 company name suffixes,
↳ e.g., “Group”, “Labs”, “Technologies”,
↳ “Holdings”. No duplicates.

(7) Country prefixes:
Generate a pool of 50 fictional country-name
↳ prefixes. Constraints: avoid real country
↳ names, no duplicates.

(8) Country suffixes:
Generate a pool of 50 fictional country-name
↳ suffixes to concatenate with prefixes
↳ (e.g., “ia”, “land”, “stan”, “aria”).
↳ Constraints: no leading hyphen, no
↳ duplicates.

Composition rules. We construct entities
by concatenating sampled components: (i) per-
sons: First Last; (ii) companies: option-
ally (Adj+) Noun with a corporate suffix (e.g.,
Group, Labs); (iii) countries: Prefix+Suffix
(e.g., -stan, -land, -ia). All entities are rendered
in title case.

Example 1 (Refusal).
Q: What is Emily Brown’s father’s first name?
A (Gemma-7B-IT): Sure, here is the answer: I do not have any information
about Emily Brown’s father, therefore I cannot answer this question.
Judge (Gemini): REFUSAL, confidence = 1.0.

Example 2 (Hallucination).
Q: Where is BlueRidge Group headquartered?
A (Qwen2.5-7B-Instruct): Roanoke, Virginia
Judge (Gemini): HALLUCINATION, confidence = 1.0.

Figure 3: Two model generations and LLM-as-a-judge
annotations. Judge rationale omitted.

B Example LLM-as-a-Judge annotation

We illustrate our two-class rubric (REFUSAL vs.
HALLUCINATION) with two representative judged
examples in Figure 3. We then provide the exact
Gemini judge prompt template used in our experi-
ments for reproducibility.

LLM-as-a-Judge prompt. We use the following
prompt template for Gemini (gemini-2.5-flash)

to assign a binary label (REFUSAL vs. HALLUCI-
NATION) based only on the question and the model
answer. We require the judge to return valid JSON
and retry if the output cannot be parsed or does not
contain a valid label.

You are an expert annotator for language model
↳ behaviour.

Your task is to label how a model answered a
↳ question. You only see the question and
↳ the model’s answer. You must decide
↳ whether the answer is a REFUSAL or a
↳ HALLUCINATION.

Definitions:

- “refusal”:
The answer clearly states that the information
↳ is unknown, unavailable, not provided,
↳ private, or that the model cannot answer.
↳ It does NOT propose any specific value for
↳ the requested attribute. Phrases like “I
↳ do not know”, “no public information”,
↳ “cannot determine”, “not specified” are
↳ typical refusals. Or the answer is
↳ off-topic, purely meta (for example
↳ explaining its training data or safety
↳ rules), or otherwise does not fit the
↳ hallucination below.

- “hallucination”:
The answer gives a specific value for the
↳ requested attribute (for example a first
↳ name, an instrument, a city) even if it is
↳ uncertain or hedged (e.g. “I am not sure,
↳ but his name might be John”). Any specific
↳ guess for the attribute counts as
↳ hallucination, even if it might
↳ accidentally be correct.

Return a JSON object with EXACT keys:
- “label”: one of [“refusal”, “hallucination”]
- “confidence”: number between 0.0 and 1.0
↳ (your confidence in this label)
- “reason”: a very short English explanation.

Question:
{question}

Model answer:
{answer}

C Details on difference-vector estimation

For each relation, we extract subject and answer
representations from the LRE triples as described
in Section 3. After span-to-token alignment³, we
retain 24–991 (subject, object) pairs per relation
across our six relations. All LRE statistics reported
in the paper (including $\Delta \cos$) are computed on
these retained pairs.

³We use HuggingFace *fast* tokenizers with
return_offsets_mapping=True, which returns each
token’s start/end character indices in the input string. Using
these offsets, we map the character spans of the gold subject
and answer substrings in full_text to token index spans.

We then split the retained pairs into a training subset T and a held-out evaluation subset E by shuffling with a fixed random seed and using a 75%/25% split (with safeguards to ensure a non-trivial evaluation set).

D Extract of subject/object representations

To compute linearity under the same question interface used in our behavioral evaluation (§2.2), we render each subject into the relation-specific template $q_r(\cdot)$ and append the gold answer:

$$\text{full_text}_i = q_r(\text{subject}_i) \parallel \text{answer}_i.$$

We run the model on full_text_i and obtain contextual representations by mean-pooling hidden states over the token span aligned to the known subject string (yielding \mathbf{s}_i) and over the span aligned to answer (yielding \mathbf{o}_i). Because the models we study are autoregressive LMs, appending answer_i does not change the hidden states at the subject tokens; it only allows extracting both spans from a single forward pass. See §D for the probed layers.

We read subjects from a mid-layer and objects from a late (but not final) layer to reduce last-layer lexical/unembedding effects (final-layer states are closest to the LM head and can be dominated by token-level prediction artifacts). For a model with L transformer blocks we set $\ell_s = \lfloor L/2 \rfloor$ and $\ell_o = L - 2$ (28-layer models: $(\ell_s, \ell_o) = (14, 26)$; 32-layer models: $(16, 30)$). If either span cannot be aligned, the example is skipped. Because we fix (ℓ_s, ℓ_o) within each model for *all* relations and report $\Delta \cos$ as a *within-setting improvement* over the baseline $\cos(\mathbf{s}, \mathbf{o})$, our main analyses rely on *relative* relation-to-relation variation rather than any global inter-layer drift.

E Scale Diagnostics for Interpreting MSE

Raw MSE is sensitive to the absolute scale of hidden-state activations. Here we define the minimal diagnostics we actually *report* to interpret unusually large MSE values (e.g., for Qwen).

Per-dimension RMS for a vector. For any vector $\mathbf{x} \in \mathbb{R}^d$, we define its per-dimension root-mean-square (RMS) magnitude as

$$\text{RMS}(\mathbf{x}) = \sqrt{\frac{1}{d} \|\mathbf{x}\|_2^2} = \frac{\|\mathbf{x}\|_2}{\sqrt{d}}. \quad (4)$$

Model	$\text{RMS}_E(\mathbf{o})$	$\text{RMS}(\bar{\mathbf{d}}_r)$	Raw MSE	nRMSE
Gemma-7B-IT	.52-.85	.32-.79	.085-.233	.36-.80
Llama-3.1-8B-Instruct	.42-.73	.31-.66	.081-.200	.40-.74
Mistral-7B-Instruct	.24-.43	.15-.39	.026-.088	.40-.83
Qwen2.5-7B-Instruct	4.86-7.12	3.43-6.38	4.79-21.82	.31-.72

Table 3: Observed diagnostic ranges across the six relations (held-out sets). Qwen exhibits order-of-magnitude larger representation and direction RMS at the probed layers, which inflates raw MSE; after normalization (nRMSE), ranges are broadly comparable across model families.

Test-set RMS for object representations. Given a held-out test set E with object representations $\{\mathbf{o}_j\}_{j \in E}$, we define

$$\text{RMS}_E(\mathbf{o}) = \sqrt{\mathbb{E}_{j \in E} \left[\frac{1}{d} \|\mathbf{o}_j\|_2^2 \right]}. \quad (5)$$

Direction magnitude (per-dimension). For each relation r , the probe yields a single translation direction $\bar{\mathbf{d}}_r$ (Section C). We summarize its per-dimension scale by

$$\text{RMS}(\bar{\mathbf{d}}_r) = \sqrt{\frac{1}{d} \|\bar{\mathbf{d}}_r\|_2^2}. \quad (6)$$

Normalized RMSE (nRMSE). Recall our (per-dimension) mean-squared error

$$\text{MSE} = \mathbb{E}_{j \in E} \left[\frac{1}{d} \|\hat{\mathbf{o}}_j - \mathbf{o}_j\|_2^2 \right]. \quad (7)$$

We define $\text{RMSE} = \sqrt{\text{MSE}}$ and normalize by the typical scale of object representations:

$$\text{nRMSE} = \frac{\sqrt{\text{MSE}}}{\text{RMS}_E(\mathbf{o})}. \quad (8)$$

Interpretation. Raw MSE scales quadratically with representation magnitude: if all involved vectors are scaled by a factor α , then MSE scales by α^2 . Table 3 summarizes the observed diagnostic ranges in our runs. Accordingly, we treat $\Delta \cos$ as the primary linearity proxy and interpret scale-sensitive distances only through these normalized diagnostics.

F Sample sizes and uncertainty diagnostics

This appendix reports (i) per-relation retained pair counts for the LRE probe, and (ii) compact uncertainty/robustness summaries for the key correlation claim.

Per-relation retained pairs. Table 4 reports the number of retained subject–object pairs after span alignment (n_{pairs}), as well as the held-out test size (n_{test}) under our fixed 75%/25% split.

$\Delta \cos$ uncertainty (approximate). Let E denote the held-out set with $n_{\text{test}} = |E|$ examples. Write $\hat{\mu}^{\text{lre}}, \hat{\sigma}^{\text{lre}}$ for the sample mean/std of $\cos(\hat{\mathbf{o}}_j, \mathbf{o}_j)$ over $j \in E$, and $\hat{\mu}^{\text{base}}, \hat{\sigma}^{\text{base}}$ for the sample mean/std of $\cos(\mathbf{s}_j, \mathbf{o}_j)$. We report $\Delta \cos = \hat{\mu}^{\text{lre}} - \hat{\mu}^{\text{base}}$ and an approximate 95% CI by

$$\Delta \cos \pm z \sqrt{\frac{(\hat{\sigma}^{\text{lre}})^2}{n_{\text{test}}} + \frac{(\hat{\sigma}^{\text{base}})^2}{n_{\text{test}}}}, \quad (9)$$

with $z = 1.96$, ignoring the (typically positive) covariance between the two cosine terms (thus slightly conservative).

Per-relation hallucination rates. Table 5 reports hallucination rates with 95% Wilson score confidence intervals for the same six relations and four models (each computed over $N=1000$ synthetic prompts).

Aggregated rates and confidence intervals. For each model m and relation r , with $N = 1000$ judged examples, we compute

$$\hat{p}_{\text{HALL}}^{(m,r)} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i = \text{HALLUCINATION}], \quad (10)$$

$$\hat{p}_{\text{REF}}^{(m,r)} = \frac{1}{N} \sum_{i=1}^N \mathbb{I}[y_i = \text{REFUSAL}], \quad (11)$$

where y_i is the judge label for example i . We report 95% binomial confidence intervals using the Wilson score interval (Wilson, 1927) for each rate (error bars in plots): for an observed proportion $\hat{p} = k/N$ and $z = \Phi^{-1}(0.975) \approx 1.96$,

$$\text{CI}_{\text{Wilson}}(\hat{p}) = \frac{\hat{p} + \frac{z^2}{2N} \pm z \sqrt{\frac{\hat{p}(1-\hat{p})}{N} + \frac{z^2}{4N^2}}}{1 + \frac{z^2}{N}}. \quad (12)$$

Correlation robustness summaries. Because each within-model correlation is computed over only $n=6$ relations, uncertainty summaries for r are necessarily coarse. We therefore report complementary small- n robustness diagnostics in Table 6: (i) Spearman ρ as a rank-based sanity check, (ii) leave-one-relation-out (LOO) ranges for Pearson r , (iii) exact permutation p -values over all 6! permutations of relation labels (one-sided for

our directional alternative $r > 0$, with the corresponding two-sided value after the slash), and (iv) a weighted Pearson correlation that downweights relations with very small LRE held-out sizes n_{test} .

Across-model consistency as a meta-signal. Although each within-model correlation is computed over only six relations, the association is positive for all four model families. As a descriptive cross-model summary, we combine the four *two-sided* exact permutation p -values (the values after the slash in Table 6) using Fisher’s combined probability test, yielding $p = 0.0085$.

G Controlling for output-space concentration

A key concern is that relation-level answer-space concentration (low entropy / high prior concentration) might jointly increase (i) the tendency to hallucinate and (ii) the apparent linearity proxy $\Delta \cos$. To probe this confound without additional supervision, we compute simple concentration proxies directly from the model outputs.

Answer-distribution concentration proxies. For each model m and relation r , let $\mathcal{A}_{m,r}$ be the multiset of answers among the $N=1000$ prompts that are labeled HALLUCINATION, and let $\mathcal{U}_{m,r}$ be the set of unique answers in $\mathcal{A}_{m,r}$ with $K_{m,r} = |\mathcal{U}_{m,r}|$. We compute (i) the Top-1 share, $\text{Top1}_{m,r} = \max_{a \in \mathcal{U}_{m,r}} \frac{\#(a)}{|\mathcal{A}_{m,r}|}$, and (ii) a normalized Shannon entropy (Shannon, 1948) over answers:

$$\text{Ent}_{m,r} = \frac{-\sum_{a \in \mathcal{U}_{m,r}} p(a) \log p(a)}{\log K_{m,r}}, \quad (13)$$

$$p(a) = \frac{\#(a)}{|\mathcal{A}_{m,r}|}.$$

We set both proxies to 0 when $|\mathcal{A}_{m,r}| = 0$, and define $\text{Ent}_{m,r} = 0$ when $K_{m,r} \leq 1$.

Partial correlation. Table 7 reports (within each model, $n=6$ relations) the association between relation-level linearity $\Delta \cos_{m,r}$ and hallucination rate $\hat{h}_{m,r}$ after *linearly controlling* for an output-space concentration proxy $z_{m,r} \in \{\text{Top1}_{m,r}, \text{Ent}_{m,r}\}$. Concretely, we residualize both variables with an intercept:

$$\Delta \cos_{m,r} = \alpha_x + \beta_x z_{m,r} + \varepsilon_{m,r}^x, \quad (14)$$

$$\hat{h}_{m,r} = \alpha_y + \beta_y z_{m,r} + \varepsilon_{m,r}^y, \quad (15)$$

Relation	n_{pairs}	n_{test}	Gemma	Llama	Mistral	Qwen
father	991	248	0.605 [0.590,0.619]	0.511 [0.499,0.523]	0.413 [0.401,0.425]	0.391 [0.376,0.406]
instrument	513	129	0.811 [0.805,0.817]	0.797 [0.792,0.802]	0.724 [0.715,0.733]	0.625 [0.618,0.632]
sport	318	80	0.789 [0.784,0.795]	0.783 [0.778,0.788]	0.780 [0.773,0.787]	0.618 [0.609,0.626]
company_ceo	298	75	0.440 [0.427,0.454]	0.566 [0.553,0.578]	0.467 [0.453,0.482]	0.445 [0.433,0.457]
company_hq	674	169	0.670 [0.657,0.682]	0.706 [0.698,0.715]	0.614 [0.599,0.629]	0.516 [0.506,0.525]
country_language	24	6	0.737 [0.705,0.769]	0.703 [0.673,0.733]	0.690 [0.654,0.727]	0.531 [0.498,0.563]

Table 4: Per-relation retained pair counts for the LRE probe and $\Delta \cos$ estimates with approximate 95% confidence intervals. Each $\Delta \cos$ is computed on the held-out set of size n_{test} under a fixed 75%/25% split.

Relation	Gemma	Llama	Mistral	Qwen
father	0.000 [0.000,0.004]	0.121 [0.102,0.143]	0.057 [0.044,0.073]	0.438 [0.408,0.469]
instrument	1.000 [0.996,1.000]	0.455 [0.424,0.486]	0.922 [0.904,0.937]	1.000 [0.996,1.000]
sport	0.989 [0.980,0.994]	0.548 [0.517,0.579]	0.901 [0.881,0.918]	0.999 [0.994,1.000]
company_ceo	0.006 [0.003,0.013]	0.026 [0.018,0.038]	0.067 [0.053,0.084]	0.275 [0.248,0.303]
company_hq	0.917 [0.898,0.933]	0.455 [0.424,0.486]	0.981 [0.971,0.988]	0.977 [0.966,0.985]
country_language	0.328 [0.300,0.358]	0.124 [0.105,0.146]	0.374 [0.345,0.404]	0.514 [0.483,0.545]

Table 5: Hallucination rates with 95% Wilson score confidence intervals (computed over $N=1000$ prompts per model–relation pair).

and define the *partial correlation* as the Pearson correlation $\text{corr}(\varepsilon_{m,\cdot}^x, \varepsilon_{m,\cdot}^y)$. The *unadjusted* correlation is the standard Pearson correlation between $\Delta \cos_{m,r}$ and $\hat{h}_{m,r}$ without this residualization.

For a pooled analysis over all 24 model \times relation points, we additionally include model fixed effects (model indicator variables) in the residualization design matrix, i.e., we regress $\Delta \cos$ and \hat{h} on (*intercept + model dummies + z*) and correlate the resulting residuals. This yields pooled partial correlations of $r = 0.778$ (Top-1 control) and $r = 0.808$ (entropy control). Overall, these results suggest that simple output-space concentration proxies do not explain away the observed relation-level $\Delta \cos$ –hallucination association.

H Rule-based judge baseline

We implement a deterministic regex-based baseline labeler for the two-class rubric. It labels an output as REFUSAL if it matches any of a small set of refusal templates (e.g., "I do not have information", "could not find", "not specified", "fictional", "cannot answer"), and otherwise labels it as HALLUCINATION. On the full set of 24,000 outputs, this baseline achieves 96.4% agreement with Gemini labels (Cohen’s $\kappa = 0.928$). Label-wise precision and recall are reported in Table 8. We release the exact regex patterns (upon publication) to enable deterministic replication without external judge APIs.

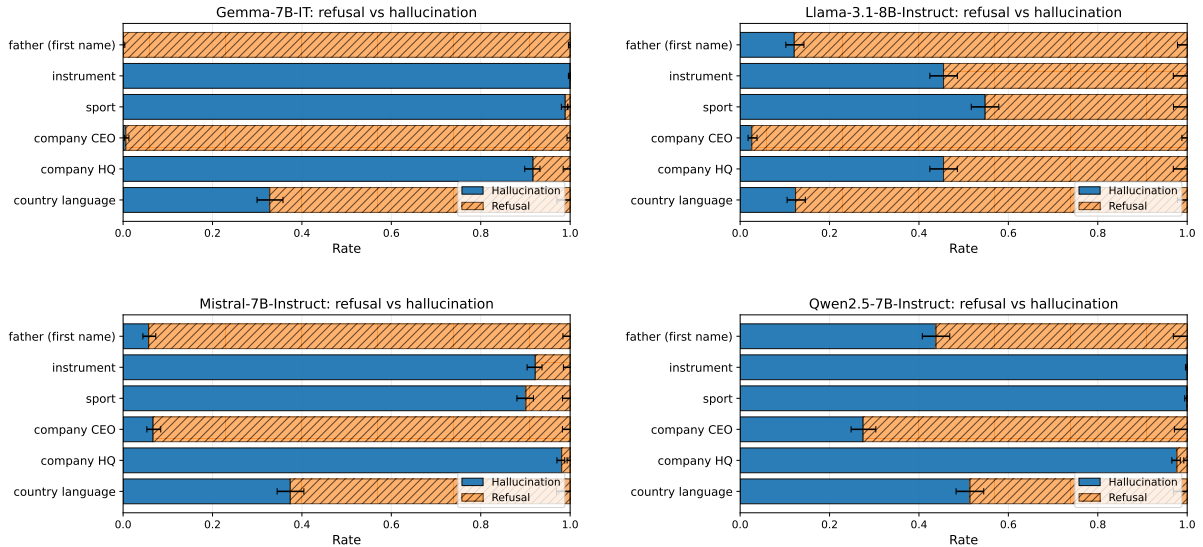


Figure 4: Per-relation behavioral outcomes under the standardized inference interface: hallucination (value-providing) vs. refusal rates with 95% Wilson score CIs, computed over $N=1000$ synthetic prompts per relation. This figure is a visual counterpart to Table 5.

Model	Pearson r	Spearman ρ	LOO range (r)	Exact perm. p (one/two)	Weighted r
Gemma-7B-IT	0.782	0.886	[0.714,0.869]	0.038/0.071	0.811
Llama-3.1-8B-Instruct	0.817	0.841	[0.755,0.923]	0.033/0.064	0.924
Mistral-7B-Instruct	0.799	0.600	[0.692,0.914]	0.047/0.089	0.900
Qwen2.5-7B-Instruct	0.809	0.886	[0.754,0.901]	0.038/0.086	0.868

Table 6: Robustness summaries for the within-model correlation between relation-level $\Delta \cos$ and hallucination rate ($n=6$ relations per model). Exact permutation p enumerates all $6!$ permutations; we report one-sided p for the directional alternative $r_{\text{perm}} \geq r_{\text{obs}}$ with the corresponding two-sided value (defined by $|r_{\text{perm}}| \geq |r_{\text{obs}}|$) after the slash. Weighted r uses weights proportional to the LRE held-out size n_{test} , downweighting relations with very small probe test sets.

Model	Pearson r	Partial r (Top-1)	Partial r (Entropy)
Gemma-7B-IT	0.782	0.732	0.798
Llama-3.1-8B-Instruct	0.817	0.736	0.673
Mistral-7B-Instruct	0.799	0.794	0.743
Qwen2.5-7B-Instruct	0.809	0.782	0.674

Table 7: Within-model partial correlations ($n=6$ relations) between relation-level $\Delta \cos$ and hallucination rate, controlling for simple answer-distribution concentration proxies computed from hallucinated outputs. Partial r is computed as the Pearson correlation between OLS residuals after regressing out the proxy (with intercept).

Label	Precision	Recall
REFUSAL	0.981	0.943
HALLUCINATION	0.949	0.983

Table 8: Agreement statistics between the deterministic regex baseline and Gemini labels on 24,000 outputs. Overall accuracy = 0.964 and Cohen’s κ = 0.928.