

How Context Shapes Truth: Geometric Transformations of Statement-level Truth Representations in LLMs

Anonymous ACL submission

Abstract

Large Language Models (LLMs) often encode whether a statement is true as a vector in their residual stream activations. These vectors, also known as *truth vectors*, have been studied in prior work, however how they change when context is introduced remains unexplored. We study this question by measuring (1) the directional change (θ) between the truth vectors with and without context and (2) the relative magnitude of the truth vectors upon adding context. Across four LLMs and four datasets, we find that (1) truth vectors are roughly orthogonal in early layers, converge in middle layers, and may stabilize or continue increasing in later layers; (2) adding context generally increases the truth vector magnitude, i.e., the separation between true and false representations in the activation space is amplified; (3) larger models distinguish relevant from irrelevant context mainly through directional change (θ), while smaller models show this distinction through magnitude differences. We also find that context conflicting with parametric knowledge produces larger geometric changes than parametrically aligned context. To the best of our knowledge, this is the first work that provides a geometric characterization of how context transforms the truth vector in the activation space of LLMs¹.

1 Introduction

As Large language models (LLMs) get increasingly adopted in high stakes applications, it becomes important to understand how they process and represent information internally. Prior work (Hollinsworth et al., 2024; Gurnee and Tegmark, 2023; Marks and Tegmark, 2024) studies how concepts are encoded in model activations, specifically using activations from residual stream (after the MLP layer)². They find that many high-level concepts, including whether a statement is true, are

represented as linear directions (i.e. vectors) in the activation space (termed “truth directions”). Prior work (Burns et al., 2023; Azaria and Mitchell, 2023; Marks and Tegmark, 2024; Li et al., 2023; Bao et al., 2025) shows that linear classifiers can reliably separate true from false statements in the LLM activation space, implying a geometric structure to how truth is represented. However, these studies do not study how this geometry changes when context is added. While in-context learning and retrieval-augmented-generation have proven effective at improving model outputs without re-training (Brown et al., 2020; Min et al., 2022; Wei et al., 2023; Lewis et al., 2020; Gao et al., 2023), how the geometric structure of statement-level truth changes when context is added remains underexplored. It is precisely these geometric changes in the direction and magnitude of residual stream activations when context is added that we study in this work. We contribute the first characterisation of how truth geometry transforms when context is added. Understanding this has theoretical implications for how LLMs process context, and practical implications for designing retrieval-augmented and in-context learning systems that more reliably integrate contextual knowledge.

We analyze the residual stream activations when an LLM processes a statement with and without context. For both conditions, we extract the vectors in activation space that separate true from false statements i.e., the “*truth vectors*”. We hypothesize that adding context should alter this geometric structure. To test this, we examine two geometric properties: the angle between the truth vectors with and without context (θ), which captures directional change, and the relative magnitude of the truth vectors, which captures whether context amplifies or compresses the separation between true and false representations in the activation space.

Experiments with four LLMs and four datasets, spanning diverse domains and context types, show

¹Our code is anonymized and available [here](#)

²In the rest of the paper, by “residual stream” we will mean after the MLP layer without explicitly clarifying it.

the following three findings: (1) **Three-phase pattern of directional change**: Comparing truth vectors with and without context, we find that truth vectors are approximately orthogonal in early layers, converge sharply in early to middle layers, and then either stabilize or continue increasing in later layers depending on the dataset. (2) **Increase in Relative Magnitudes**: Adding context generally increases the truth vector magnitude, i.e., the separation between true and false representations in the activation space increases. (3) **Sensitivity to relevant vs irrelevant context**: On comparing relevant context with randomly generated and irrelevant context, we find that relevant context generally produces a higher directional or magnitudinal change. These findings are statistically significant across models and datasets. Collectively, our results provide novel empirical evidence on how context reshapes the geometric structure of statement-level truth representations in the activation space of LLMs.

2 Related Work

Truth Representations in LLMs Understanding how LLMs represent truth has received attention. Burns et al. (2023) introduce Contrast-Consistent Search (CCS), an unsupervised methodology showing that truth directions can be extracted from model activations. This work shows that LLMs encode truth as a linear direction in their representation space. Marks and Tegmark (2024) extend this using mass-mean probes, which compute the mean difference between activations for true and false statements to identify truth directions. Li et al. (2023) introduce Inference-Time Intervention (ITI), showing that shifting model activations along truthful directions can significantly improve LLM truthfulness. This work distinguishes between generation accuracy (measured by model output) and probe accuracy (classifying statements using intermediate activations); similarly to this, our work also focuses on internal representations rather than output behavior. Bürger et al. (2024) address the failure of truth probes to generalize across negated statements by showing that truth is represented in a two-dimensional subspace rather than a single direction. Lastly, Bao et al. (2025) find that consistent truth directions emerge in more capable models and that probes trained on factual statements generalize to in-context settings, including question answering grounded in provided passages

and abstractive summarization. However, they test whether a single probe transfers across these settings, not whether the geometric structure of truth vectors change when context is introduced. Our work addresses this gap directly.

While the above work establishes that truth has a geometric structure in the activation space of LLMs and tests if truth probes generalize in different settings, it does not directly examine how truth vectors change when context is added. It is precisely this gap that our work addresses by measuring the geometric transformations, namely, the directional change θ and relative magnitude shift between truth vectors with and without context, showing that context induces consistent layer-dependent changes.

Activation Steering and Contrastive Vectors

Prior work has shown that LLM behavior can be steered by adding contrastive vectors to model activations (Turner et al., 2024; Rimsky et al., 2024; Zou et al., 2023; Subramani et al., 2022). These vectors are typically computed as the mean difference between the activations of two contrasting conditions, such as truth and false (Li et al., 2023), toxic and non-toxic (Liu et al., 2024), or positive and negative sentiment (Turner et al., 2024). During inference, the contrasting vectors are added back to shift the model’s behavior. Here, magnitude of the contrastive vector is a key hyperparameter, serving as the strength of intervention. We take inspiration from steering techniques for our method and instead of using vectors to modify behavior, we observe how truth vectors change when context is introduced.

Context Utilization Research on in-context learning has focused on how models use instructions and exemplars to recognize tasks and learn input-output mappings (Brown et al., 2020; Min et al., 2022; Wei et al., 2023), with evidence that task recognition occurs in the middle layers (Sia et al., 2024). While prior work has focused on how LLMs utilize context by analyzing the generated outputs (Du et al., 2024; Marjanovic et al., 2024; Hagström et al., 2025) or by probing the residual stream activations directly to detect knowledge conflict signals (Zhao et al., 2024), our work focuses on how context geometrically changes the direction of truth in the residual stream activations.

false representations when context is added as:

$$rm_{k,tc-fc}^{(l)} = \frac{\|v_{k,c}^{(l)}\|^2}{\|v_{k,nc}^{(l)}\|^2} \quad (6)$$

Eq. 6 corresponds to measuring $\frac{CD}{AB}$ in Figure 2. For a dataset D , we average across all statements N_k to get the relative magnitudes for the entire dataset in each layer l as:

$$rm_{D,tc-fc}^{(l)} = \frac{1}{|N_k|} \sum_k rm_{k,tc-fc}^{(l)} \quad (7)$$

where $|N_k|$ is the total number of statements. We also calculate the vectors $v_{k,tc-fnc}^{(l)}$ and $v_{k,tnc-fc}^{(l)}$ to measure the relative magnitudes in the case when context is added to generate either true or false completions while generating the other completion without any context (see Appendix A.2).

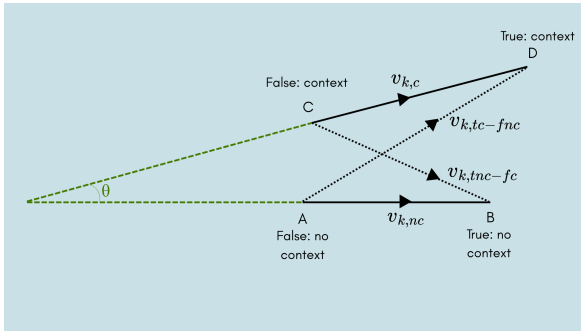


Figure 2: For a statement k , $v_{k,nc}$ (AB) is the truth vector without context and $v_{k,c}$ (CD) is the truth vector when context is added. θ is the angle between $v_{k,nc}$ and $v_{k,c}$ denoting the directional change (Eq. 4). To track relative magnitudes, we compute the ratio of L_2 distances: $\frac{CD}{AB}$, $\frac{AD}{AB}$ and $\frac{BC}{AB}$ as per Eq. 6, 10 & 11.

4 Experiments

We aim to study how the geometric structure of the truth vector (specifically its directional change θ (Eq. 5) and relative magnitude (Eq. 7)) changes when context is introduced. We select only statements where the LLM follows instructions across all four prompts (Figure 1b); see Appendix A.8.

LLMs We use four instruction-tuned models spanning different scales (3B–12B) and families: Llama-3.1-8B-Instruct (Grattafiori et al., 2024), Mistral-Nemo-12B-Instruct (Mistral AI and NVIDIA, 2024), Qwen3-4B-Instruct (Yang et al., 2025) and SmoLLM3-3B (Bakouch et al., 2025). This selection allows us to examine whether the observed geometric transformations generalize across

Dataset	Rows	Len.	Read.	Context Type
Borderlines	982	153.9	41.3	Geo. factcheck
Politifact	907	114.6	47.7	Pol. factcheck
ScienceFeedback	618	128.5	39.6	Sci. factcheck
MF2	1736	457.9	57.2	Movie synopsis
CL-Bill	500	185.3	6.9	Legal bills
CL-Company	500	657.5	10.7	Company descr.
ConflictQA-Counter	1244	82.4	46.0	Parametrically counter context
ConflictQA-Parametric	1244	50.8	53.9	Parametrically aligned context

Table 1: Dataset statistics. Len. is mean context length in words. Read. is the Flesch Reading Ease (0–100, the lower, the harder the text). Borderlines, Politifact and ScienceFeedback are subsets of DRUID. CL denotes Corporate Lobbying datasets from LegalBench.

model scale. As our task is text-generation following a specific set of instructions, we use off-the-shelf instruction fine-tuned models. We use Huggingface API for inference with greedy decoding sampling to ensure reproducibility. All experiments were conducted on NVIDIA A100 and H100 GPUs, requiring approximately 500 GPU hours.

Datasets We use datasets containing statements and relevant contexts: Druid (Hagström et al., 2025), MF2 (Zaranis et al., 2025), ConflictQA (Xie et al., 2024) and LegalBench (Guha et al., 2023). We select three subsets from Druid: Borderlines, Politifact and ScienceFeedback and analyze them separately as the context type varies across them. See Table 1 for a dataset summary and Appendix A.1 for details. While Druid, MF2 and LegalBench contain real world data, ConflictQA is a synthetic dataset. We use two subsets from ConflictQA: Parametric and Counter. ConflictQA-Parametric contains context which is aligned to the LLM’s parametric knowledge and ConflictQA-Counter contains context which goes against the parametric knowledge. In Table 1 we also show the Fleisch score of context, which approximates human difficulty in understanding text (Flesch, 1948).

5 Results and Discussion

5.1 Directional Change across Layers

To understand how context changes the truth representations in the residual stream, we begin with a layer-wise analysis of directional change θ . Figure 3 shows how θ changes across layers. A lower θ means higher similarity between truth vectors with and without context. All four LLMs show a consistent 3-phase pattern: θ remains high (near

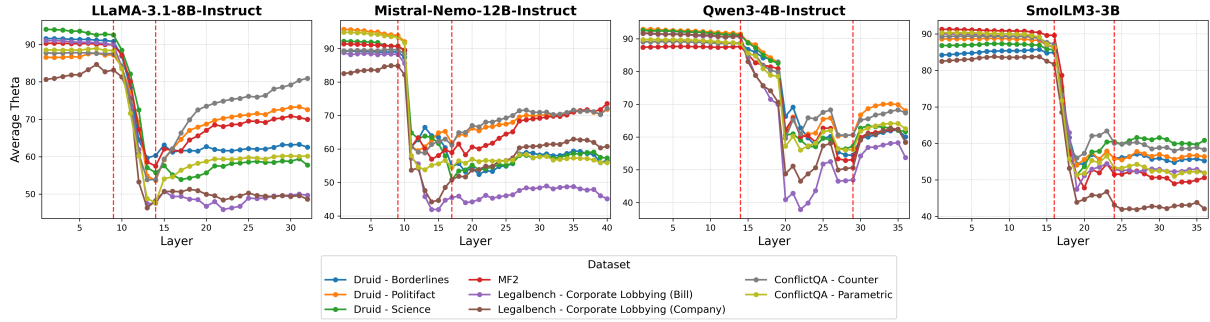


Figure 3: Layer wise plot of average θ in degrees across different models and datasets indicating the directional change in truth vectors when context is added. Vertical red lines indicate the beginning of a new phase. Across all the settings, we observe three phases: Phase-1, where the truth vectors are almost orthogonal, Phase-2, where the truth vectors become more similar and finally Phase-3, where truth vectors stabilize or continue increasing.

orthogonal) in early layers, drops sharply in middle layers to reach a minimum, and then either stabilizes or increases in later layers. LLaMA and Mistral begin decreasing around layer 9, reaching minima near layer 15, while smaller models (Qwen, SmoLLM) show prolonged early phases until layers 14–16 with later minima (layers 20–25). In later layers, behavior varies by model-dataset combination. This 3-phase pattern, and especially the convergence of the truth vectors in the middle layers, is consistent with prior findings that early layers handle low-level input processing, middle layers encode semantic information, and later layers shift toward next-token prediction (Ghandeharioun et al., 2024). Early LLM layers have been related to capturing syntactic meaning (Li and Subramani, 2025). As such, the direction of “truth” can potentially have less meaning in early layers - leading to orthogonality in phase-1. We also verify this using probes built to classify truth using residual stream activations. We observe that accuracies often peak in the middle layers and are usually low in the earlier layers (Appendix A.4). Further, we note that while larger models compress the initial stage into fewer layers (until layer 9), smaller models take longer (until layers 14-16).

The convergence in the middle layers indicates that truth vectors with and without context become more similar. This aligns with prior work showing that middle layers are the primary site for semantic encoding (Ghandeharioun et al., 2024; Geva et al., 2023), factual knowledge retrieval (Meng et al., 2022), and task-relevant representations (Hendel et al., 2023). Ghandeharioun et al. (2024) observe that steering vectors are most effective in middle layers, where input processing has concluded but next-token prediction has not yet dominated, and

Sia et al. (2024) find that task recognition in machine translation occurs in similar layers. Similarly to Azaria and Mitchell (2023), we also observe that accuracies of probes often peak in the middle layers (Appendix A.4). Notably, θ never reaches zero, suggesting that while truth vectors converge, the models maintain distinct representations for statements with and without context.

In phase 3, θ shows a flat trend for most datasets, suggesting that truth vectors have largely converged by the middle layers. However, for ConflictQA-Counter and Politifact, θ increases in later layers for LLaMA, Mistral, and Qwen, possibly reflecting continued processing of context that conflicts with parametric knowledge. Notably, θ values for ConflictQA-Counter consistently exceed those for ConflictQA-Parametric, indicating that contradictory context induces greater directional shift than aligned context. This is consistent with prior findings that LLMs exhibit confirmation bias towards memory-aligned information (Xie et al., 2024) and that knowledge conflicts arise from competing memory heads and context heads in later layers (Jin et al., 2024). Further, prior work suggests that deeper layers are often redundant and can be pruned with limited performance loss (Men et al., 2025), though the final layer remains important. Our results suggest that later-layer contributions may also be context-dependent.

5.2 Relative Magnitude

To understand how context affects the separation between true and false representations, we compute the relative magnitude of the truth vector when context is added (as described in Section 3). The results for the relative magnitudes in the final layer are shown in Table 2 and a layerwise analysis is pre-

Dataset	LLaMA	Mistral	Qwen	SmolLM
Borderlines	1.18*	1.08*	1.13*	1.11*
Politifact	1.01	0.85	1.07*	1.15*
ScienceFeedback	1.10*	0.87	1.00	1.19*
MF2	1.13*	1.00	1.06*	0.96
CL-Bill	1.06*	1.06*	1.00	0.95
CL-Company	1.15*	1.18*	1.06*	1.00
ConflictQA - Counter	1.20*	0.98	0.98	1.26*
ConflictQA - Param	1.34*	1.02	1.06*	1.16*

Table 2: Relative magnitude (Eq. 7) averaged over statements from the final LLM layer across datasets. Values above 1 mean that the truth vector magnitude increases when context is added. * marks stat. significance of $p < 0.05$ with the Wilcoxon signed-rank test.

sented in Figure 4. Relative magnitude values A above (resp. below) 1 mean that the separation between true and false representations increase (resp. decrease) when context is added.

Figure 4 shows a certain variability across LLM layers with respect to relative magnitude, however one common pattern is a peak in middle layers followed by a decline and eventual stabilization. LLaMA shows early-layer variability, an upward spike around layers 15–20, then stabilization. Mistral exhibits a spike around layers 10–15, a sharp decline through layers 15–20, then stabilization. Qwen shows a spike around layer 22, then declining until layer 27 before stabilizing. SmolLM displays a spike around layers 17–19, a slight decrease, and a secondary smaller spike around layers 25–27, though this later spike is absent for MF2 and Corporate Lobbying. Across models, middle-layer spikes almost always exceed 1, indicating that the separation between true and false representations are maximum in the middle layers. Notably, middle layers are often responsible for semantic encoding (Ghandeharioun et al., 2024). Although relative magnitudes decrease toward later layers, they generally remain above 1, even in the final layers (Table 2). In the final layer, LLaMA increases the average relative magnitude of the truth vector across 7 out of 8 datasets. However, the results are mixed for other models. Additional results are found in Appendix A.2 and Appendix A.7.

Note that we also examine whether θ and relative magnitude correlate with changes in output probability for “True” and “False” tokens when context is added (Appendix A.5). We find some correlations, but not consistently across datasets and models. This suggests that θ and relative magnitudes do not necessarily translate to probabilistic differences in the output generation.

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	2.84*	6.87*	0.71	6.32*	5.32*
	Politifact	11.81*	13.88*	11.92*	12.22*	13.47*
	ScienceFeedback	2.55*	4.79*	3.36*	7.07*	3.97*
	MF2	1.13*	1.71*	4.91*	7.12*	2.08*
	CL-Bill	-1.73	-0.15	-1.13	-1.39	-1.63
	CL-Company	-10.43	-10.9	-3.93	-3.81	-2.86
	ConflictQA-Counter	22.38*	22.16*	18.18*	18.10*	13.01*
	ConflictQA-Param	2.03	2.49	-7.51	-7.00	-10.29
Mistral	Borderlines	3.09*	-0.04	3.66*	0.21	1.07
	Politifact	18.05*	19.12*	20.97*	19.01*	18.53*
	ScienceFeedback	1.49	4.58*	8.05*	4.69*	5.49*
	MF2	7.61*	10.22*	9.96*	12.97*	5.41*
	CL-Bill	-6.05	-5.4	-2.2	-3.74	-1.43
	CL-Company	-4.95	-2.50	1.84*	5.13*	2.12*
	ConflictQA-Counter	14.97*	16.67*	15.94*	16.18*	12.46*
	ConflictQA-Param	0.54	3.19*	2.63*	4.89*	0.12
Qwen	Borderlines	0.73	1.48	0.96	-6.08	-12.78
	Politifact	1.04	0.09	-0.84	-5.15	-2.49
	ScienceFeedback	4.43	4.97*	3.74	2.99	2.02
	MF2	-2.51	-2.07	-4.85	-12.05	-5.81
	CL-Bill	0.28	-2.16	1.78	-2.17	-4.86
	CL-Company	-0.34	-0.26	-1.67	-3.87	-5.16
	ConflictQA-Counter	6.97*	8.51*	6.47*	-0.49	-0.79
	ConflictQA-Param	-4.11	-1.49	-7.53	-7.94	-9.11
SmolLM	Borderlines	-4.63	-2.35	-3.98	1.54	-3.44
	Politifact	-4.80	1.15	0.14	1.34	0.41
	ScienceFeedback	-1.42	2.58*	1.62	4.65*	0.62
	MF2	-8.54	-6.47	-1.85	-1.54	-0.65
	CL-Bill	-1.61	-1.46	-1.30	0.10	-0.78
	CL-Company	-6.55	-5.65	-1.70	-2.44	-2.45
	ConflictQA-Counter	2.29*	2.06*	2.17*	4.94*	2.04*
	ConflictQA-Param	-3.07	-2.64	-2.13	1.73	-1.83

Table 3: Random VS relevant context. Each value is the mean difference in θ between relevant and random context(s) in the final LLM layer. Char, Word, Salad, Wiki and Shuffle represent various random contexts. * marks stat. significance of $p < 0.05$ with the Wilcoxon signed-rank test, meaning that relevant context induces more directional change than random context.

5.3 Relevant versus Random Context

Motivated by prior work showing that adding unrelated context dramatically reduces model performance (Shi et al., 2023; Yoran et al., 2024), we compare the effect of adding relevant versus random context to study if relevant context produces different geometric changes than random context, we experiment with five different contexts varying in degree of randomness: (1) context of “random characters”, such that words have no linguistic meaning; (2) context of “random words”, randomly sampled from the NLTK english corpus and ordered randomly, such that the sentence has no meaning; (3) context of “random salad”, where the sentence is grammatical but incoherent (e.g. *colorless green ideas sleep furiously*); (4) “random wiki” context, where paragraphs are randomly sampled from wikipedia; and (5) “random shuffle” context, where we shuffle the contexts from the same dataset such the statement and contexts do not match. Except for (5), in (1)-(4) we control for the length of contexts so that the random context has the same number of words as the original context for that statement. See Appendix A.3 for examples details on the length distribution.

Tables 3 and 4 show the effect of random con-

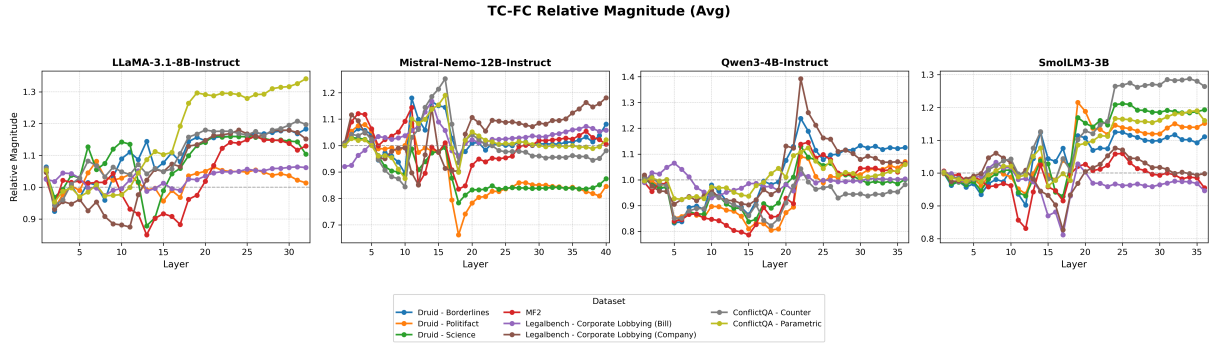


Figure 4: Layer wise plot of average relative magnitudes across different models and datasets indicating the increase in the magnitude of truth vector when context is added. Early layers show variability, followed by a peak in the middle layers. The values decrease and stabilize towards the final layers.

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	0.22*	0.26*	0.24*	-0.08	-0.19
	Politifact	0.10*	0.11*	0.00	0.01	-0.03
	ScienceFeedback	0.13*	0.16*	0.19*	0.14*	0.00
	MF2	0.12*	0.10*	0.05*	-0.18	-0.05
	CL-Bill	0.07*	0.05*	-0.11	-0.01	0.00
	CL-Company	0.07*	-0.05	0.04*	0.01	0.04*
	ConflictQA - Counter	0.39*	0.40*	0.42*	0.20*	0.21*
	ConflictQA - Param	0.60*	0.65*	0.63*	0.39*	0.38*
Mistral	Borderlines	0.08*	0.17*	0.15*	0.09*	-0.02
	Politifact	0.08*	0.08*	-0.03	0.00	0.01
	ScienceFeedback	0.02*	0.06*	-0.02	0.01	-0.01
	MF2	-0.11	-0.05	-0.03	0.01*	0.01*
	CL-Bill	0.01*	0.04*	0.07*	0.01*	0.00
	CL-Company	0.04*	0.05*	0.06*	-0.03	-0.01
	ConflictQA - Counter	0.12*	0.22*	0.10*	0.06*	-0.03
	ConflictQA - Param	0.21*	0.25*	0.15*	0.09*	0.06*
Qwen	Borderlines	0.07*	0.00	0.15*	0.19*	0.22*
	Politifact	0.08*	0.01	0.04*	0.03*	0.10*
	ScienceFeedback	-0.03	-0.08	-0.01	-0.02	0.04*
	MF2	0.08*	0.03*	0.03*	0.05*	0.02*
	CL-Bill	0.01*	-0.04	-0.01	0.00	0.01*
	CL-Company	0.08*	0.01*	-0.03	-0.01	0.01*
	ConflictQA - Counter	0.14*	0.06*	0.09*	0.06*	0.09*
	ConflictQA - Param	0.20*	0.13*	0.16*	0.14*	0.14*
SmoLLM	Borderlines	0.15*	0.20*	0.22*	0.09*	0.10*
	Politifact	0.18*	0.20*	0.23*	0.11*	0.14*
	ScienceFeedback	0.11*	0.14*	0.20*	0.14*	0.08*
	MF2	0.17*	0.17*	0.10*	0.05*	0.03*
	CL-Bill	0.05*	0.04*	0.04*	0.02*	0.02*
	CL-Company	0.02*	-0.01	0.03*	0.06*	0.01*
	ConflictQA - Counter	0.35*	0.34*	0.37*	0.25*	0.23*
	ConflictQA - Param	0.27*	0.25*	0.26*	0.13*	0.15*

Table 4: Random VS relevant context. Each value is the mean difference in relative magnitude between relevant and random context(s) in the final LLM layer. * marks stat. significance of $p < 0.05$ with the Wilcoxon signed-rank test, meaning that the true and false representations are more separated for relevant than random context. The remaining notation is as in Table 3.

text upon θ and relative magnitude. Specifically, we show the difference in θ and relative magnitude between the original relevant context and “random contexts”. We use the final layer of the model for comparison, since this is the closest layer responsible for text generation. We also show the Bonferroni corrected differences in Appendix A.6. We describe our findings next.

Larger Models show directional sensitivity: Each value in Table 3 is the difference between θ with the relevant context and θ with a random

context from the final LLM layer. A significant difference means that relevant context causes a greater directional shift in the residual stream than random context. We see that for larger models, LLaMA and Mistral, relevant context generally induces a significantly higher θ compared to random contexts, specifically for Borderlines, Politifact, ScienceFeedback, MF2 and ConflictQA-Counter. The primary exceptions are the Corporate Lobbying datasets from LegalBench, where random contexts sometimes result in a higher θ . However, for smaller models, we generally observe that θ values are smaller for relevant context when compared to random contexts, with the exception of ConflictQA-Counter dataset, where the contexts are designed to contradict the parametric knowledge of the model. We discuss this in Section 5.4.

Smaller Models show magnitudinal sensitivity: Each value in Table 4 is the difference in relative magnitude between relevant and random context from the final LLM layer. A significant difference means that the true and false representations are more separated for relevant context than random context. We see that for smaller models (Qwen and SmoLLM) the relative magnitudes are significantly higher for relevant context compared to non-relevant context across most datasets, even though the difference in θ is often negative or insignificant. This suggests that smaller models encode contextual relevance through magnitude scaling rather than directional changes. SmoLLM shows positive magnitudinal differences in almost all settings, despite showing negative differences in θ . Larger models, LLaMA and Mistral, also generally show higher relative magnitudes for relevant context, except for specific instances in the Corporate Lobbying dataset.

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	Both	Both	Mag	Theta	Theta
	Politifact	Both	Both	Theta	Theta	Theta
	ScienceFeedback	Both	Both	Both	Theta	Theta
	MF2	Both	Both	Both	Theta	Theta
	CL-Bill	Mag	Mag	None	None	None
	CL-Company	Mag	None	Mag	Mag	Mag
	ConflictQA - Counter	Both	Both	Both	Both	Both
	ConflictQA - Param	Mag	Mag	Mag	Mag	Mag
Mistral	Borderlines	Both	Mag	Both	Mag	None
	Politifact	Both	Both	Theta	Theta	Theta
	ScienceFeedback	Mag	Both	Theta	Theta	Theta
	MF2	Theta	Theta	Theta	Both	Both
	CL-Bill	Mag	Mag	Mag	Mag	None
	CL-Company	Mag	Mag	Both	Theta	Theta
	ConflictQA - Counter	Both	Both	Both	Both	Theta
	ConflictQA - Param	Mag	Both	Both	Both	Mag
Qwen	Borderlines	Mag	None	Mag	Mag	Mag
	Politifact	Mag	None	Mag	Mag	Mag
	ScienceFeedback	None	Theta	None	None	Mag
	MF2	Mag	Mag	Mag	Mag	Mag
	CL-Bill	Mag	None	None	None	Mag
	CL-Company	Mag	Mag	None	None	Mag
	ConflictQA - Counter	Both	Both	Both	Mag	Mag
	ConflictQA - Param	Mag	Mag	Mag	Mag	Mag
SmolLM	Borderlines	Mag	Mag	Mag	Mag	Mag
	Politifact	Mag	Mag	Mag	Mag	Mag
	ScienceFeedback	Mag	Both	Mag	Both	Mag
	MF2	Mag	Mag	Mag	Mag	Mag
	CL-Bill	Mag	Mag	Mag	Mag	Mag
	CL-Company	Mag	None	Mag	Mag	Mag
	ConflictQA - Counter	Both	Both	Both	Both	Both
	ConflictQA - Param	Mag	Mag	Mag	Mag	Mag

Table 5: Comparison between random and relevant context. **Both** means θ and relative magnitude is significantly greater for relevant than random context. **Theta** (resp. **Mag**) means that only θ (resp. relative magnitude) is significantly greater for relevant context. **None** means that neither θ or relative magnitude is greater than random context. The rest of notation is as in Table 3.

Lastly, Table 5 shows a joint overview of the results from θ and relative magnitude. Overall, we see that either θ or relative magnitude is significantly greater for relevant context than random context. This means that, in general, meaningful context tends to have a greater impact on the geometry of the representations of truth statement.

Collectively, our results show that θ and relative magnitude capture aspects of how context changes truth vectors. Across models, relevant context produces significantly higher θ (in larger models) or higher relative magnitude (in smaller models) compared to random context, indicating sensitivity to context relevance. However, larger representational changes do not imply beneficial utilization. ConflictQA-Counter yields the highest θ values yet has contradictory information processing, while LegalBench shows minimal differences, suggesting models struggle with complex legal text.

5.4 ConflictQA and LegalBench

We now discuss some idiosyncrasies of two particular datasets. One dataset shows consistent effects across all models: ConflictQA-Counter. Both θ and magnitude are significantly greater for relevant context compared to random context across LLaMA, Mistral, Qwen, and SmolLM (Table 5 shows “Both” for most random context types). This dataset contains contexts that explicitly contradict the model’s parametric knowledge, suggesting that counter-memory information produces a particularly strong directional and magnitudinal signal. We also observe that ConflictQA-Parametric has much lower, and often negative directional shift, even for larger models (Table 3). This could be a result of the confirmation bias towards parametrically aligned context (Xie et al., 2024).

LegalBench Corporate Lobbying datasets often fail to show significant differences between relevant and random context, particularly for θ . These datasets have notably low Flesch readability scores (6.9 and 10.7 compared to 40–57 for other datasets from Table 1), indicating highly technical legal language. This suggests that when context is sufficiently complex or domain-specific, models may struggle to extract a meaningful signal that distinguishes it from random text.

6 Conclusion

We investigate how context shapes truth representations in large language models by analyzing directional changes (θ) and separation (relative magnitude) between true and false statements across layers. To the best of our knowledge, this is the first work to analyze how truth vectors change when context is added. First, we observe a three-phase pattern: truth vectors are orthogonal in the early layers, converge in the middle layers and depending on the context, may stabilize or continue increasing in the later layers. Second, adding context generally increases the separation between true and false representations. Third, we observe that relevant context produces larger changes than random context in most cases. Our work provides a useful lens for understanding how models process context to shape the truth vectors.

7 Limitations

Our study has several limitations. First, we extract truth representations from only the first token position, though relevant information may be

563	distributed across multiple tokens. Second, our	and Jacob Steinhardt. 2023. Eliciting Latent Predictions from Transformers with the Tuned Lens . <i>arXiv preprint</i> . Version Number: 6.	614
564	comparisons between relevant and random context		615
565	rely on the final layer only. Third, the ConflictQA		616
566	dataset was constructed to evaluate parametric	Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie	617
567	knowledge models for comparatively larger models.	Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind	618
568	Models in our study may lack this knowledge or	Neelakantan, Pranav Shyam, Girish Sastry, Amanda	619
569	may have encountered the dataset during pretrain-	Askeell, Sandhini Agarwal, Ariel Herbert-Voss,	620
570	ing. Fourth, our analysis is correlational rather than	Gretchen Krueger, Tom Henighan, Rewon Child,	621
571	causal, and we leave interventional experiments for	Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu,	622
572	future work. Finally, we evaluate on a limited set	Clemens Winter, and 12 others. 2020. Language	623
573	of English-language datasets; extending to other	models are few-shot learners. In <i>Proceedings of the</i>	624
574	languages and domains is a direction for future	<i>34th International Conference on Neural Information</i>	625
575	work.	<i>Processing Systems, NIPS '20</i> , Red Hook, NY, USA.	626
		Curran Associates Inc. Event-place: Vancouver, BC,	627
		Canada.	628
576	8 Ethical Consideration	Collin Burns, Haotian Ye, Dan Klein, and Jacob Stein-	629
577	This work is primarily aimed at understanding how	hardt. 2023. Discovering Latent Knowledge in Language Models Without Supervision . In <i>The Eleventh International Conference on Learning Representations</i> .	630
578	LLMs represent truth internally when context is		631
579	added. We do not foresee direct negative societal		632
580	impacts from this interpretability study. However,	Lennart Bürger, Fred A. Hamprecht, and Boaz Nadler.	634
581	understanding truth vectors could potentially be	2024. Truth is Universal: Robust Detection of Lies in LLMs . <i>arXiv preprint</i> . ArXiv:2407.12831 [cs].	635
582	dual-use. While it may help improve factuality		636
583	in LLMs and detect misinformation, it could the-	Kevin Du, Vésteinn Snæbjarnarson, Niklas Stoehr, Jen-	637
584	oretically inform adversarial attacks that manipu-	nifer White, Aaron Schein, and Ryan Cotterell. 2024.	638
585	late model outputs. All datasets used are publicly	Context versus Prior Knowledge in Language Models . In <i>Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 13211–13235, Bangkok, Thailand. Association for Computational Linguistics.	639
586	available and do not contain personally identifiable		640
587	information.		641
			642
			643
588	References	Rudolph Flesch. 1948. A new readability yardstick . <i>Journal of Applied Psychology</i> , 32(3):221–233.	644
589	Amos Azaria and Tom Mitchell. 2023. The Internal State of an LLM Knows When It’s Lying . In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 967–976, Singapore. Association for Computational Linguistics.		645
590		Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	646
591		Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang,	647
592		and Haofen Wang. 2023. Retrieval-Augmented Generation for Large Language Models: A Survey . <i>arXiv preprint</i> . Version Number: 5.	648
593			649
594	Elie Bakouch, Loubna Ben Allal, Anton Lozhkov, Nou-		650
595	mane Tazi, Lewis Tunstall, Carlos Miguel Patiño,	Mor Geva, Jasmijn Bastings, Katja Filippova, and Amir	651
596	Edward Beeching, Aymeric Roucher, Aksel Joonas	Globerson. 2023. Dissecting Recall of Factual Associations in Auto-Regressive Language Models . In <i>The 2023 Conference on Empirical Methods in Natural Language Processing</i> .	652
597	Reedi, Quentin Gallouédec, Kashif Rasul, Nathan		653
598	Habib, Clémentine Fourrier, Hynek Kydlicek, Guil-		654
599	herme Penedo, Hugo Larcher, Mathieu Morlon, Vaib-	Asma Ghandeharioun, Ann Yuan, Marius Guerard,	656
600	hav Srivastav, Joshua Lochner, and 4 others. 2025.	Emily Reif, Michael A. Lepori, and Lucas Dixon.	657
601	SmolLM3: smol, multilingual, long-context reasoner .	2024. Who’s asking? User personas and the mechanics of latent misalignment . In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	658
602			659
603	Yuntai Bao, Xuhong Zhang, Tianyu Du, Xinkui Zhao,		660
604	Zhengwen Feng, Hao Peng, and Jianwei Yin. 2025.	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	662
605	Probing the Geometry of Truth: Consistency and Generalization of Truth Directions in LLMs Across Logical Transformations and Question Answering Tasks . In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 682–700, Vienna, Austria. Association for Computational Linguistics.	Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh Goyal, Anthony Hartshorn, Aobo Yang, Archi Mitra, Archie Sravankumar, Artem Korenev, Arthur Hinsvark, and 540 others. 2024. The Llama 3 Herd of Models . <i>arXiv preprint</i> . Version Number: 3.	663
606			664
607			665
608			666
609			667
610			668
611			669
612	Nora Belrose, Igor Ostrovsky, Lev McKinney, Zach Fur-		670
613	man, Logan Smith, Danny Halawi, Stella Biderman,		671

670	Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher	Systems, NIPS '23, Red Hook, NY, USA. Curran As-	729
671	Re, Adam Chilton, Aditya Narayana, Alex Chohlas-	sociates Inc. Event-place: New Orleans, LA, USA.	730
672	Wood, Austin Peters, Brandon Waldon, Daniel Rock-		
673	more, Diego Zambrano, Dmitry Talisman, Enam	Michael Li and Nishant Subramani. 2025. Echoes of	731
674	Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gre-	BERT: Do Modern Language Models Rediscover the	732
675	gory M. Dickinson, Haggai Porat, Jason Hegland,	Classical NLP Pipeline? <i>arXiv preprint</i> . Version	733
676	and 21 others. 2023. LegalBench: A Collaboratively	Number: 4.	734
677	Built Benchmark for Measuring Legal Reasoning		
678	in Large Language Models . In <i>Thirty-seventh Con-</i>	Sheng Liu, Haotian Ye, Lei Xing, and James Zou. 2024.	735
679	<i>ference on Neural Information Processing Systems</i>	In-context vectors: making in context learning more	736
680	<i>Datasets and Benchmarks Track</i> .	effective and controllable through latent space steer-	737
681		ing. In <i>Proceedings of the 41st International Confer-</i>	738
682	Wes Gurnee and Max Tegmark. 2023. Language Mod-	<i>ence on Machine Learning, ICML'24, Vienna, Aus-</i>	739
683	els Represent Space and Time . <i>arXiv preprint</i> . Ver-	tria. JMLR.org.	740
684	Version Number: 3.		
685	684 Lovisa Hagström, Sara Vera Marjanovic, Haeun Yu, Ar-	Sara Vera Marjanovic, Haeun Yu, Pepa Atanasova,	741
686	nav Arora, Christina Lioma, Maria Maistro, Pepa	Maria Maistro, Christina Lioma, and Isabelle Au-	742
687	Atanasova, and Isabelle Augenstein. 2025. A Re-	genstein. 2024. DYNAMICQA: Tracing Internal	743
688	ality Check on Context Utilisation for Retrieval-	Knowledge Conflicts in Language Models . In <i>Find-</i>	744
689	Augmented Generation . In <i>Proceedings of the 63rd</i>	<i>ings of the Association for Computational Linguistics:</i>	745
690	<i>Annual Meeting of the Association for Computational</i>	<i>EMNLP 2024</i> , pages 14346–14360, Miami, Florida,	746
691	<i>Linguistics (Volume 1: Long Papers)</i> , pages 19691–	USA. Association for Computational Linguistics.	747
692	19730, Vienna, Austria. Association for Computa-		
693	tional Linguistics.	Samuel Marks and Max Tegmark. 2024. The Geometry	748
694	693 Rooe Hendel, Mor Geva, and Amir Globerson. 2023.	of Truth: Emergent Linear Structure in Large Lan-	749
695	In-Context Learning Creates Task Vectors . In <i>Find-</i>	guage Model Representations of True/False Datasets .	750
696	<i>ings of the Association for Computational Linguis-</i>	In <i>First Conference on Language Modeling</i> .	751
697	<i>tics: EMNLP 2023</i> , pages 9318–9333, Singapore.		
698	Association for Computational Linguistics.	Xin Men, Mingyu Xu, Qingyu Zhang, Qianhao Yuan,	752
699	698 Oskar John Hollinsworth, Curt Tigges, Atticus Geiger,	Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei	753
700	and Neel Nanda. 2024. Language Models Linearly	Han, and Weipeng Chen. 2025. ShortGPT: Lay-	754
701	Represent Sentiment . In <i>Proceedings of the 7th</i>	ers in Large Language Models are More Redundant	755
702	<i>BlackboxNLP Workshop: Analyzing and Interpret-</i>	Than You Expect . In <i>Findings of the Association</i>	756
703	<i>ing Neural Networks for NLP</i> , pages 58–87, Miami,	<i>for Computational Linguistics: ACL 2025</i> , pages	757
704	Florida, US. Association for Computational Linguis-	20192–20204, Vienna, Austria. Association for Com-	758
705	tics.	putational Linguistics.	759
706	705 Zhuoran Jin, Pengfei Cao, Hongbang Yuan, Yubo Chen,	Kevin Meng, David Bau, Alex J. Andonian, and Yonatan	760
707	Jiexin Xu, Huaijun Li, Xiaojian Jiang, Kang Liu, and	Belinkov. 2022. Locating and Editing Factual Asso-	761
708	Jun Zhao. 2024. Cutting Off the Head Ends the Con-	ciations in GPT . In <i>Advances in Neural Information</i>	762
709	flict: A Mechanism for Interpreting and Mitigating	<i>Processing Systems</i> .	763
710	Knowledge Conflicts in Language Models . In <i>Find-</i>		
711	<i>ings of the Association for Computational Linguistics</i>	Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe,	764
712	<i>ACL 2024</i> , pages 1193–1215, Bangkok, Thailand	Mike Lewis, Hannaneh Hajishirzi, and Luke Zettle-	765
713	and virtual meeting. Association for Computational	moyer. 2022. Rethinking the Role of Demonstrations:	766
714	Linguistics.	What Makes In-Context Learning Work? In <i>Proceed-</i>	767
715	714 Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	<i>ings of the 2022 Conference on Empirical Methods in</i>	768
716	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	<i>Natural Language Processing</i> , pages 11048–11064,	769
717	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	Abu Dhabi, United Arab Emirates. Association for	770
718	täschel, Sebastian Riedel, and Douwe Kiela. 2020.	Computational Linguistics.	771
719	Retrieval-augmented generation for knowledge-		
720	intensive NLP tasks. In <i>Proceedings of the 34th</i>	Mistral AI and NVIDIA. 2024. Mistral-NeMo-Instruct-	772
721	<i>International Conference on Neural Information Pro-</i>	2407 .	773
722	<i>cessing Systems, NIPS '20, Red Hook, NY, USA.</i>		
723	Curran Associates Inc. Event-place: Vancouver, BC,	Nostalgebraist. 2020. interpreting GPT: the logit lens .	774
724	Canada.		
725	724 Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter	Nina Rimskey, Nick Gabrieli, Julian Schulz, Meg Tong,	775
726	Pfister, and Martin Wattenberg. 2023. Inference-time	Evan Hubinger, and Alexander Turner. 2024. Steer-	776
727	intervention: eliciting truthful answers from a lan-	ing Llama 2 via Contrastive Activation Addition . In	777
728	guage model. In <i>Proceedings of the 37th Interna-</i>	<i>Proceedings of the 62nd Annual Meeting of the As-</i>	778
	<i>national Conference on Neural Information Processing</i>	<i>sociation for Computational Linguistics (Volume 1:</i>	779
		<i>Long Papers)</i> , pages 15504–15522, Bangkok, Thai-	780
		land. Association for Computational Linguistics.	781

782	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H. Chi, Nathanael Schärli, and Denny Zhou. 2023. Large Language Models Can Be Easily Distracted by Irrelevant Context . In <i>Proceedings of the 40th International Conference on Machine Learning</i> , volume 202 of <i>Proceedings of Machine Learning Research</i> , pages 31210–31227. PMLR.	838
783		839
784		840
785		841
786		842
787		843
788		
789		
790	Suzanna Sia, David Mueller, and Kevin Duh. 2024. Where does In-context Learning \textbackslash\textbackslash Happen in Large Language Models? In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	
791		
792		
793		
794		
795	Nishant Subramani, Nivedita Suresh, and Matthew Peters. 2022. Extracting Latent Steering Vectors from Pretrained Language Models . In <i>Findings of the Association for Computational Linguistics: ACL 2022</i> , pages 566–581, Dublin, Ireland. Association for Computational Linguistics.	
796		
797		
798		
799		
800		
801	Alexander Matt Turner, Lisa Thiergart, Gavin Leech, David Udell, Juan J. Vazquez, Ulisse Mini, and Monte MacDiarmid. 2024. Steering Language Models With Activation Engineering . <i>arXiv preprint</i> . ArXiv:2308.10248 [cs].	
802		
803		
804		
805		
806	Jerry Wei, Jason Wei, Yi Tay, Dustin Tran, Albert Webson, Yifeng Lu, Xinyun Chen, Hanxiao Liu, Da Huang, Denny Zhou, and Tengyu Ma. 2023. Larger language models do in-context learning differently . <i>arXiv preprint</i> . Version Number: 2.	
807		
808		
809		
810		
811	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2024. Adaptive Chameleon or Stubborn Sloth: Revealing the Behavior of Large Language Models in Knowledge Conflicts . In <i>The Twelfth International Conference on Learning Representations</i> .	
812		
813		
814		
815		
816	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, Chujie Zheng, Dayiheng Liu, Fan Zhou, Fei Huang, Feng Hu, Hao Ge, Haoran Wei, Huan Lin, Jialong Tang, and 41 others. 2025. Qwen3 Technical Report . <i>arXiv preprint</i> . Version Number: 1.	
817		
818		
819		
820		
821		
822		
823	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. 2024. Making Retrieval-Augmented Language Models Robust to Irrelevant Context . In <i>The Twelfth International Conference on Learning Representations</i> .	
824		
825		
826		
827		
828	Emmanouil Zaranis, António Farinhas, Saul Santos, Beatriz Canaverde, Miguel Moura Ramos, Aditya K Surikuchi, André Viveiros, Baohao Liao, Elena Bueno-Benito, Nithin Sivakumaran, Pavlo Vasylenko, Shoubin Yu, Sonal Sannigrahi, Wafaa Mohammed, Ben Peters, Danae Sánchez Villegas, Elias Stengel-Eskin, Giuseppe Attanasio, Jaehong Yoon, and 12 others. 2025. Movie Facts and Fibs (MF²): A Benchmark for Long Movie Understanding . <i>arXiv preprint</i> . Version Number: 1.	
829		
830		
831		
832		
833		
834		
835		
836		
837		
	Yu Zhao, Xiaotang Du, Giwon Hong, Aryo Pradipta Gema, Alessio Devoto, Hongru Wang, Xuanli He, Kam-Fai Wong, and Pasquale Minervini. 2024. Analysing the Residual Stream of Language Models Under Knowledge Conflicts . <i>arXiv preprint</i> . Version Number: 2.	844
		845
		846
		847
		848
		849
		850
		851
		852
	Andy Zou, Long Phan, Sarah Chen, James Campbell, Phillip Guo, Richard Ren, Alexander Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li, Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt Fredrikson, and 2 others. 2023. Representation Engineering: A Top-Down Approach to AI Transparency . <i>arXiv preprint</i> . Version Number: 4.	853
		854
		855
		856
		857
		858
		859
		860
		861
		862
		863
		864
		865
		866
		867
		868
		869
		870
		871
		872
		873
		874
		875
		876
		877
		878
		879
		880
		881
		882
		883
		884
		885
		886
		887

888 A.2 Additional Relative Magnitudes

889 For a statement k , we calculate the change in separation
 890 between true and false representations when
 891 context is added to generate either true or false
 892 completions, while generating the other comple-
 893 tion without any context.

$$894 v_{k,tc-fnc}^{(l)} = a_{k, \text{True}, c}^{(l)} - a_{k, \text{False}, nc}^{(l)} \quad (8)$$

$$895 v_{k,tnc-fc}^{(l)} = a_{k, \text{True}, nc}^{(l)} - a_{k, \text{False}, c}^{(l)} \quad (9)$$

$$896 rm_{k,tc-fnc}^{(l)} = \frac{\|v_{k,tc-fnc}^{(l)}\|^2}{\|v_{k,nc}^{(l)}\|^2} \quad (10)$$

$$898 rm_{k,tnc-fc}^{(l)} = \frac{\|v_{k,tnc-fc}^{(l)}\|^2}{\|v_{k,nc}^{(l)}\|^2} \quad (11)$$

899 Equation 10 corresponds to $\frac{AD}{AB}$ and Equation 11
 900 corresponds to $\frac{BC}{AB}$ in Figure 2.

$$901 rm_{D,tc-fnc}^{(l)} = \frac{1}{|N_k|} \sum_k rm_{k,tc-fnc}^{(l)} \quad (12)$$

$$903 rm_{D,tnc-fc}^{(l)} = \frac{1}{|N_k|} \sum_k rm_{k,tnc-fc}^{(l)} \quad (13)$$

904 where $|N_k|$ represents the total number of state-
 905 ments.

906 In general, we observe that the relative magni-
 907 tudes averaged over all the statements are greater
 908 than 1. Except for ConflictQA Counter dataset
 909 with Qwen3-4B-Instruct, we find that at least one
 910 of TC-FC, TC-FNC and TNC-FC have a relative
 911 magnitude greater than 1. This suggests that con-
 912 text generally increases the separation between true
 913 and false points. For both TC-FNC and TNC-FC
 914 we observe that all models except Qwen increases
 915 the relative magnitudes of the truth vector across
 916 all the datasets.

917 A.3 Comparison between Relevant and 918 Non-Relevant Context

919 Figure 6 shows an example of the randomly gener-
 920 ated context for each random type. Random char-
 921 acters context are created by randomly selecting
 922 characters and joining them to create a word. Ran-
 923 dom words context are created by randomly select-
 924 ing words from all the english word present in the
 925 NLTK corpus. Random salad context are created
 926 by repeatedly selecting a predefined sentence struc-
 927 ture made up of parts of speech (such as articles,

Model	Dataset	TC-FNC	TNC-FC
LLaMA	Borderlines	1.55*	1.61*
	Politifact	1.48*	1.50*
	ScienceFeedback	1.32*	1.27*
	MF2	1.89*	1.87*
	CL-Bill	1.14*	1.19*
	CL-Company	1.68*	1.56*
	ConflictQA-Counter	1.41*	1.33*
	ConflictQA-Param	1.46*	1.49*
Mistral	Borderlines	1.20*	1.39*
	Politifact	1.15*	1.12*
	ScienceFeedback	1.21*	1.07*
	MF2	1.30*	1.29*
	CL-Bill	1.21*	1.14*
	CL-Company	1.32*	1.28*
	ConflictQA-Counter	1.10*	1.10*
	ConflictQA-Param	1.12*	1.13*
Qwen	Borderlines	1.16*	1.06*
	Politifact	1.00	1.00
	ScienceFeedback	0.93	1.09*
	MF2	1.04*	1.05*
	CL-Bill	0.92	0.93
	CL-Company	0.98	1.00
	ConflictQA-Counter	0.92	0.94
	ConflictQA-Param	0.94	0.96
SmoLLM	Borderlines	1.18*	1.57*
	Politifact	1.39*	1.39*
	ScienceFeedback	1.40*	1.37*
	MF2	1.26*	1.26*
	CL-Bill	1.09*	1.05*
	CL-Company	1.23*	1.10*
	ConflictQA-Counter	1.31*	1.28*
	ConflictQA-Param	1.19*	1.18*

Table 6: Relative magnitudes averaged over statements from the final layer of model across datasets. TC-FNC denotes the relative magnitude of the truth vector when true representations have context and false representations do not have context (Equation 12). TNC-FC denotes the relative magnitude of the truth vector when true representations do not have context and false representations have context (Equation 13). A value greater than 1 indicates that the magnitude of truth vector has increased compared to the truth vector when both true and false representations do not have context. Asterisk (*) indicates statistical significance of $p < 0.05$

nouns, verbs, adjectives, and adverbs), then filling each position by randomly choosing a word from the NLTK corpus. If no suitable words are available for a given part of speech, it falls back to the placeholder word "word". Random wiki context is created by crawling text from wikipedia. Figure 7 shows the distribution of word count for relevant vs non-relevant context. As random shuffle contexts are essentially the contexts from the same dataset, they will have the exactly same distribution as relevant context.

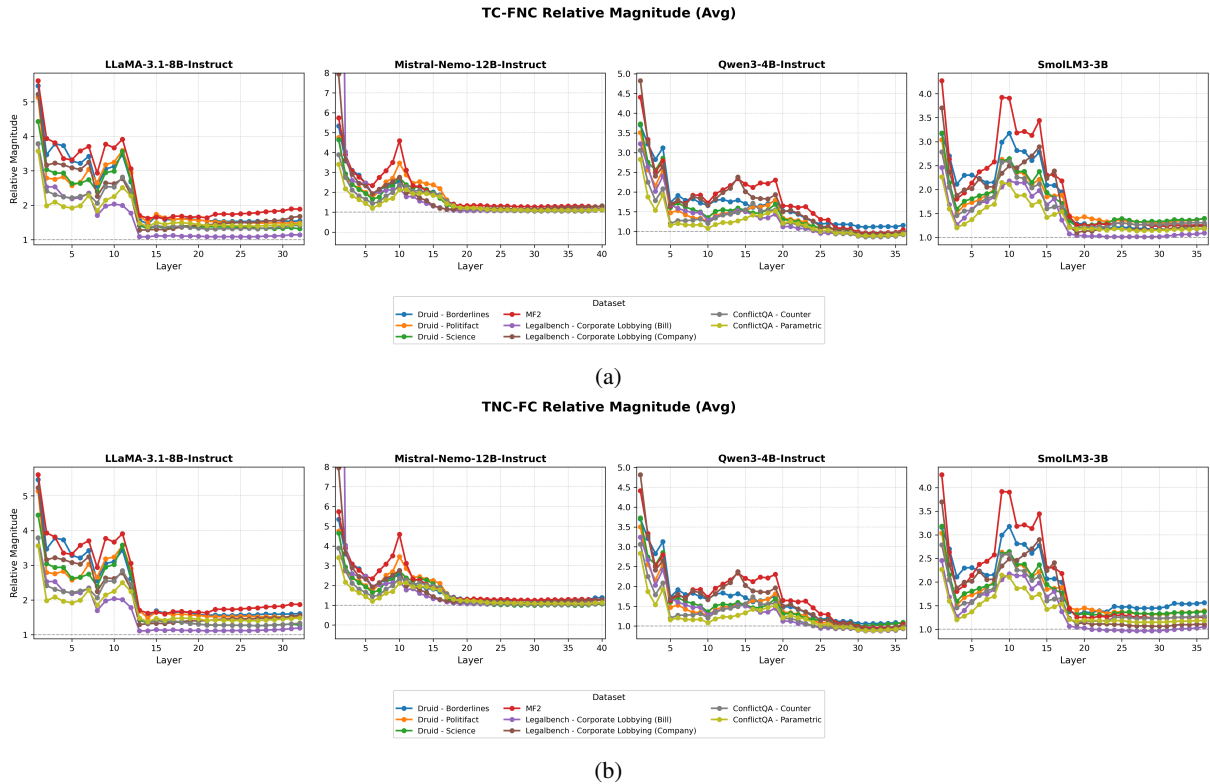


Figure 5: Additional relative magnitudes across layers for different models.

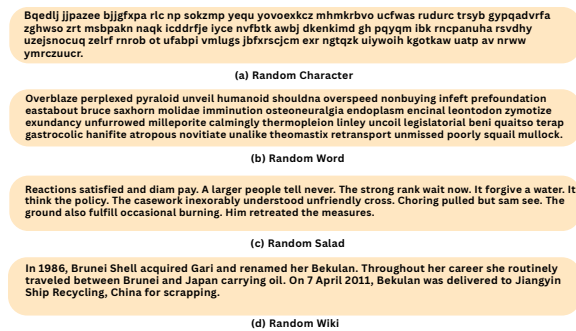


Figure 6: Contexts across different degrees of randomness.

A.4 Probes

To extract truth representations, we train linear probes to classify statements as true or false using an 80-20 train-test split. We compare four probe types: logistic regression, linear SVM, mass mean, and MLP. The results are shown in Figure 8. Probing accuracy peaks in middle layers across all models and datasets, consistent with prior findings that intermediate layers encode richer semantic information. Logistic regression and linear SVM achieve the highest accuracies, while MLP probes show weaker performance. Mass-mean probes, which compute the difference between mean true and

false activations, also achieve reasonable accuracy. Since both mass-mean probes and our metrics (θ and relative magnitude) are computed by averaging over statement-level activations, these reasonable accuracies validate our approach to extracting truth vectors.

A.5 Correlation with Normalized Probability Difference

Prior works have used the unembedding matrix to interpret intermediate representations as implicit token predictions (Nostalgebraist, 2020; Belrose et al., 2023). We compute a normalized probability difference p by taking the ratio of $P(\text{True}) - P(\text{False})$ with and without context across layers. The results are shown in Appendix Figure 9. We find that correlations between θ and p are weak across all models, suggesting directional changes do not directly track output probabilities. Relative magnitude shows stronger but inconsistent correlations (0.6–0.8 in middle layers for some datasets), capturing some relationship with output probability, but the connection is not robust across contexts.

A.6 Bonferroni Corrections

When conducting multiple statistical tests simultaneously, the probability of obtaining false pos-

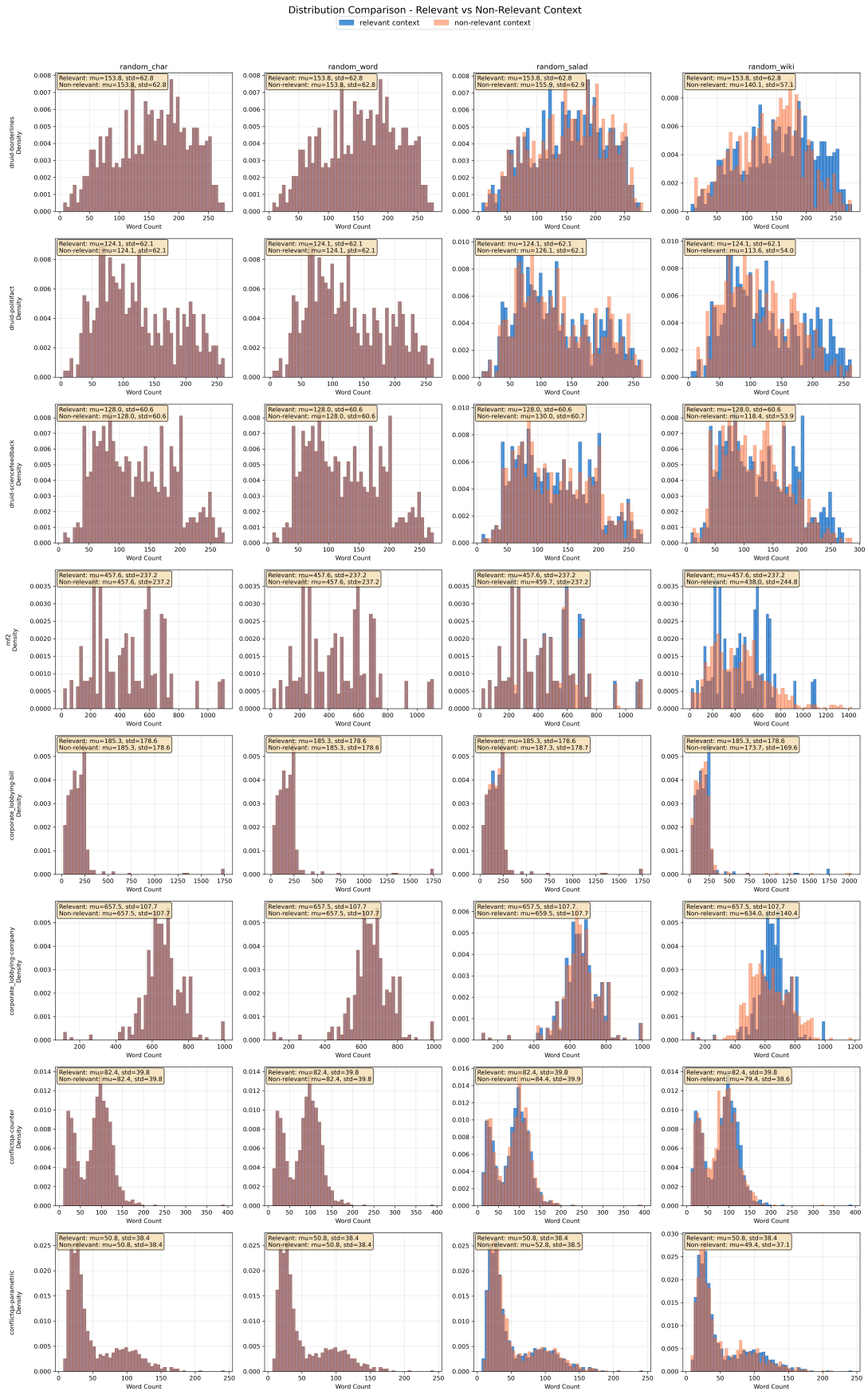


Figure 7: Distribution of word counts for relevant vs non-relevant context across datasets

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	2.84	6.87*	0.71	6.32*	5.32*
	Politifact	11.81*	13.87*	11.92*	12.22*	13.47*
	ScienceFeedback	2.55	4.79*	3.36	7.07*	3.97
	MF2	1.13	1.71*	4.91*	7.12*	2.08
	CL-Bill	-1.73	-0.15	-1.13	-1.39	-1.63
	CL-Company	-10.43	-10.90	-3.93	-3.81	-2.86
	ConflictQA - Counter	22.38*	22.16*	18.18*	18.10*	13.01*
	ConflictQA - Param	2.03	2.49	-7.51	-7.00	-10.29
	Mistral	Borderlines	3.09	-0.04	3.66	0.21
Politifact		18.05*	19.12*	20.97*	19.01*	18.53*
ScienceFeedback		1.49	4.58*	8.05*	4.69*	5.49*
MF2		7.61*	10.22*	9.96*	12.97*	5.41*
CL-Bill		-6.05	-5.40	-2.20	-3.74	-1.43
CL-Company		-4.95	-2.50	1.84	5.13*	2.12
ConflictQA - Counter		14.97*	16.67*	15.94*	16.18*	12.46*
ConflictQA - Param		0.54	3.19	2.63	4.89*	0.12
Qwen		Borderlines	0.73	1.48	0.96	-6.08
	Politifact	1.04	0.09	-0.83	-5.15	-2.49
	ScienceFeedback	4.43	4.97	3.74	2.99	2.02
	MF2	-2.51	-2.07	-4.85	-12.05	-5.81
	CL-Bill	0.28	-2.16	1.78	-2.17	-4.86
	CL-Company	-0.34	-0.26	-1.67	-3.87	-5.16
	ConflictQA - Counter	6.97*	8.51*	6.47*	-0.49	-0.79
	ConflictQA - Param	-4.11	-1.49	-7.53	-7.94	-9.11
	SmolLM	Borderlines	-4.63	-2.35	-3.98	1.54
Politifact		-4.80	1.15	0.14	1.34	0.41
ScienceFeedback		-1.42	2.58	1.62	4.65	0.62
MF2		-8.54	-6.47	-1.85	-1.54	-0.65
CL-Bill		-1.61	-1.46	-1.30	0.10	-0.78
CL-Company		-6.55	-5.65	-1.70	-2.44	-2.45
ConflictQA - Counter		2.29	2.06	2.17	4.94*	2.04
ConflictQA - Param		-3.07	-2.64	-2.13	1.73	-1.83

Table 7: Comparison between random and relevant contexts with Bonferroni correction (N=160). Notations same as in Table 3.

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	0.22*	0.26*	0.24*	-0.08	-0.19
	Politifact	0.10*	0.11*	-0.00	0.01	-0.03
	ScienceFeedback	0.13*	0.16*	0.19*	0.14*	0.00
	MF2	0.12*	0.10*	0.05*	-0.18	-0.05
	CL-Bill	0.07*	0.05*	-0.11	-0.01	0.00
	CL-Company	0.07*	-0.05	0.04*	0.01	0.04*
	ConflictQA - Counter	0.39*	0.40*	0.42*	0.20*	0.21*
	ConflictQA - Param	0.60*	0.65*	0.63*	0.39*	0.38*
Mistral	Borderlines	0.08*	0.17*	0.15*	0.09*	-0.02
	Politifact	0.08*	0.08*	-0.03	-0.00	0.01
	ScienceFeedback	0.02*	0.06*	-0.02	0.01	-0.01
	MF2	-0.11	-0.05	-0.03	0.01*	0.01*
	CL-Bill	0.01*	0.04*	0.07*	0.01*	0.00
	CL-Company	0.04*	0.05*	0.06*	-0.03	-0.01
	ConflictQA - Counter	0.12*	0.22*	0.10*	0.06*	-0.03
	ConflictQA - Param	0.21*	0.25*	0.15*	0.09*	0.06*
Qwen	Borderlines	0.07*	-0.00	0.15*	0.19*	0.22*
	Politifact	0.08*	0.01	0.04*	0.03*	0.10*
	ScienceFeedback	-0.03	-0.08	-0.01	-0.02	0.04*
	MF2	0.08*	0.03*	0.03*	0.05*	0.02*
	CL-Bill	0.01*	-0.04	-0.01	0.00	0.01*
	CL-Company	0.08*	0.01*	-0.03	-0.01	0.01*
	ConflictQA - Counter	0.14*	0.06*	0.09*	0.06*	0.09*
	ConflictQA - Param	0.20*	0.13*	0.16*	0.14*	0.14*
SmolLM	Borderlines	0.15*	0.20*	0.22*	0.09*	0.10*
	Politifact	0.18*	0.20*	0.23*	0.11*	0.14*
	ScienceFeedback	0.11*	0.14*	0.20*	0.14*	0.08*
	MF2	0.17*	0.17*	0.10*	0.05*	0.03*
	CL-Bill	0.05*	0.04*	0.04*	0.02*	0.02*
	CL-Company	0.02*	-0.01	0.03*	0.06*	0.01*
	ConflictQA - Counter	0.35*	0.34*	0.37*	0.25*	0.23*
	ConflictQA - Param	0.27*	0.25*	0.26*	0.13*	0.15*

Table 8: Comparison between random and relevant contexts with Bonferroni correction (N=160). Notations same as in Table 4

Model	Dataset	Char	Word	Salad	Wiki	Shuffle
LLaMA	Borderlines	Mag	Both	Mag	Theta	Theta
	Politifact	Both	Both	Theta	Theta	Theta
	ScienceFeedback	Mag	Both	Mag	Both	None
	MF2	Mag	Both	Both	Theta	None
	CL-Bill	Mag	Mag	None	None	None
	CL-Company	Mag	None	Mag	None	Mag
	ConflictQA - Counter	Both	Both	Both	Both	Both
	ConflictQA - Param	Mag	Mag	Mag	Mag	Mag
	Mistral	Borderlines	Mag	Mag	Mag	Mag
Politifact		Both	Both	Theta	Theta	Theta
ScienceFeedback		None	Both	Theta	Theta	Theta
MF2		Theta	Theta	Theta	Theta	Theta
CL-Bill		None	Mag	Mag	None	None
CL-Company		Mag	Mag	Mag	Theta	None
ConflictQA - Counter		Both	Both	Both	Both	Theta
ConflictQA - Param		Mag	Mag	Mag	Both	Mag
Qwen		Borderlines	Mag	None	Mag	Mag
	Politifact	Mag	None	Mag	None	Mag
	ScienceFeedback	None	None	None	None	None
	MF2	Mag	Mag	Mag	Mag	Mag
	CL-Bill	Mag	None	None	None	None
	CL-Company	Mag	None	None	None	None
	ConflictQA - Counter	Both	Both	Both	Mag	Mag
	ConflictQA - Param	Mag	Mag	Mag	Mag	Mag
	SmolLM	Borderlines	Mag	Mag	Mag	Mag
Politifact		Mag	Mag	Mag	Mag	Mag
ScienceFeedback		Mag	Mag	Mag	Mag	Mag
MF2		Mag	Mag	Mag	Mag	Mag
CL-Bill		Mag	Mag	Mag	None	Mag
CL-Company		Mag	None	Mag	Mag	None
ConflictQA - Counter		Mag	Mag	Mag	Both	Mag
ConflictQA - Param		Mag	Mag	Mag	Mag	Mag

Table 9: Comparison between random and relevant context with Bonferroni correction (N=320). Notation same as in Table 5

itives increases. For instance, at a significance level of $\alpha = 0.05$, performing 100 independent tests would yield approximately 5 false positives by chance. Bonferroni correction addresses this by adjusting the significance threshold: dividing α by the number of tests performed, thereby controlling the family-wise error rate. Tables 7, 8, and 9 present the results for comparing θ , relative magnitudes, and their combined effect, respectively, with Bonferroni correction applied. For Tables 7 and 8, we apply correction with $N=160$ tests (4 models \times 8 subsets \times 5 random conditions), yielding a corrected significance threshold of $\alpha_{\text{corrected}} = 0.05/160 = 0.0003125$. For Table 9, we apply correction with $N=320$ tests (4 models \times 8 subsets \times 5 random conditions \times 2 for θ and relative magnitudes), yielding $\alpha_{\text{corrected}} = 0.05/320 = 0.00015625$. We note that Bonferroni correction is known to be conservative, especially for large numbers of tests. Despite this strict threshold, we observe that most findings remain stable after correction, with the exception of θ comparisons for smaller models, which show reduced significance.

A.7 Additional Theta and Relative Magnitude Plots

For clarity, we plot θ and relative magnitudes along with standard error of the mean in Figures 10 and 11. For both quantities, error bars remain close to the mean values across layers. However, the two exhibit opposite patterns of variability: for θ , errors are more spread out in early layers but consolidate in later layers, whereas for relative magnitudes, early layers show less variability and later layers show more. This suggests that while the direction of the truth vector stabilizes in later layers, the separation between true and false representations becomes more variable.

A.8 Instruction Following Percentage

For all four prompts (Figure 1b), we instruct the LLM to continue generation, selecting only statements where it follows instructions across all prompts. Table 10 shows the instruction-following percentage across models and datasets. We check if the model starts the first token with “)” followed by the instructed selected choice (“Yes” or “No”) through string matching script. Additionally, we manually check some of the outputs to ensure that the generation follows the instruction.

Model	Dataset	w/o context	with context
LLaMA	ConflictQA-Counter	72.83%	51.21%
	ConflictQA-Parametric	68.89%	40.76%
	CL-Bill	100.00%	100.00%
	CL-Company	100.00%	100.00%
	Borderlines	76.62%	72.51%
	Politifact	61.82%	50.17%
	ScienceFeedback	95.30%	94.50%
	MF2	89.44%	88.65%
Mistral	ConflictQA-Counter	98.15%	92.36%
	ConflictQA-Parametric	97.75%	80.95%
	CL-Bill	98.40%	98.80%
	CL-Company	97.60%	99.40%
	Borderlines	83.33%	83.20%
	Politifact	84.09%	76.63%
	ScienceFeedback	95.97%	94.98%
	MF2	82.11%	75.46%
Qwen	ConflictQA-Counter	88.99%	69.05%
	ConflictQA-Parametric	81.11%	46.38%
	CL-Bill	74.00%	81.60%
	CL-Company	96.00%	88.40%
	Borderlines	76.62%	71.49%
	Politifact	87.27%	66.04%
	ScienceFeedback	93.29%	53.88%
	MF2	96.13%	94.30%
SmolLM	ConflictQA-Counter	91.96%	87.70%
	ConflictQA-Parametric	94.77%	76.21%
	CL-Bill	100.00%	100.00%
	CL-Company	100.00%	100.00%
	Borderlines	97.40%	93.48%
	Politifact	97.27%	90.74%
	ScienceFeedback	81.88%	77.99%
	MF2	99.48%	99.25%

Table 10: Instruction following percentage across models and datasets. "w/o context" denotes prompts without any context, while "with context" denotes prompts with context.

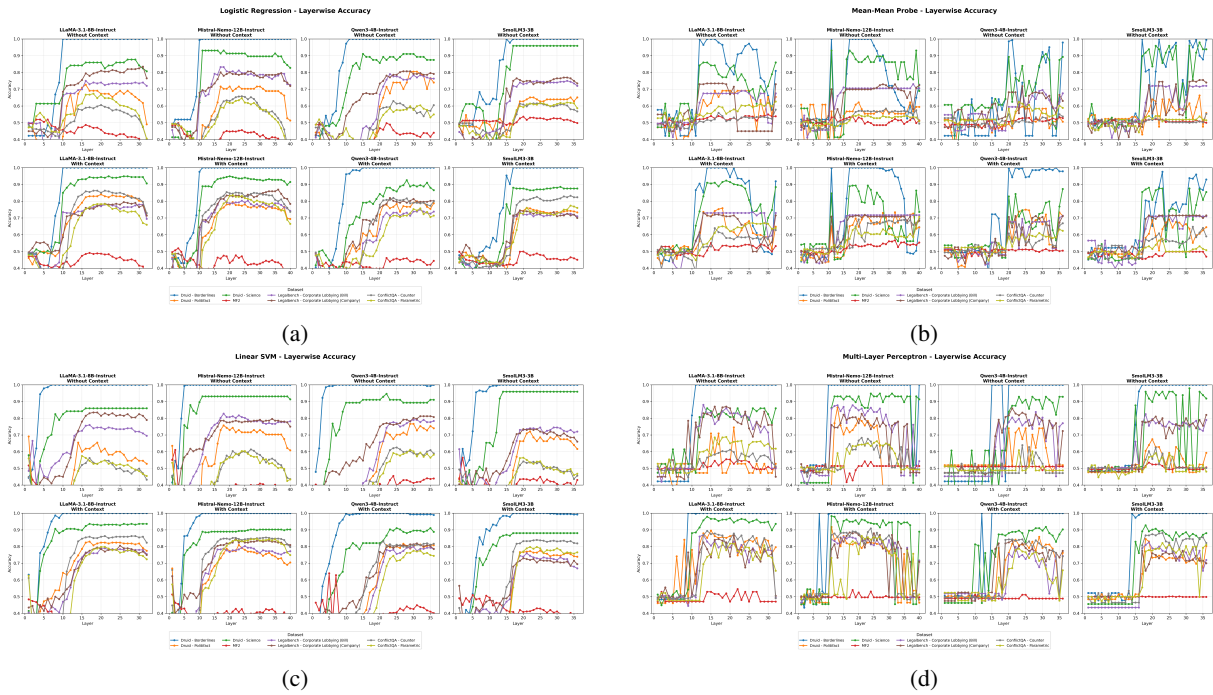


Figure 8: Accuracies of probes across layers

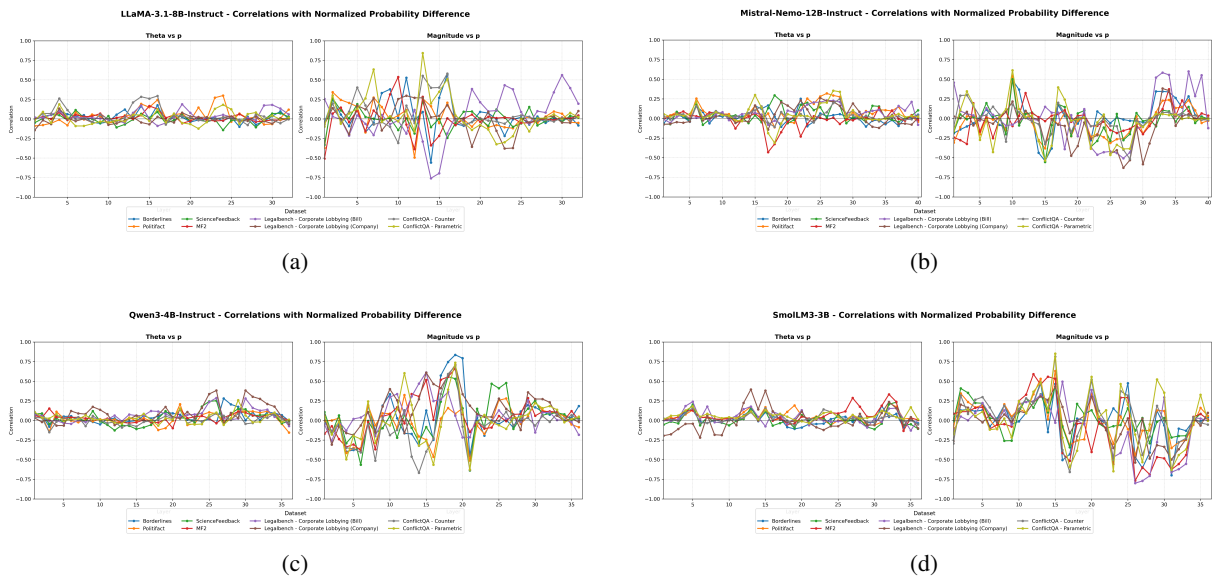
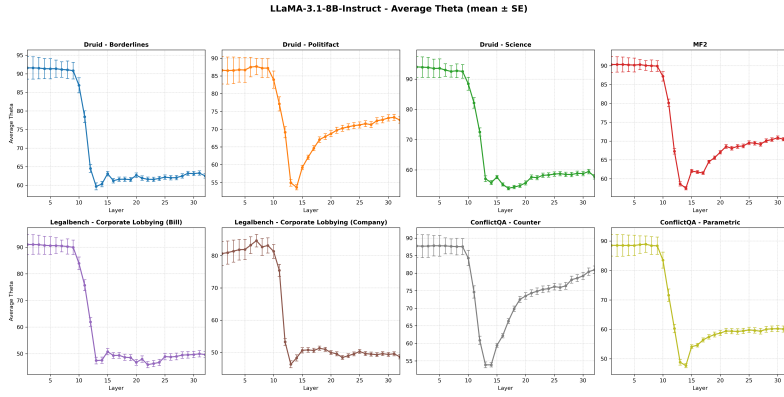
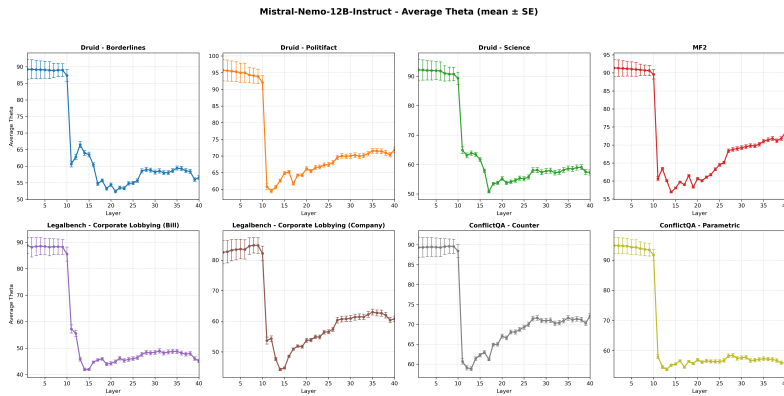


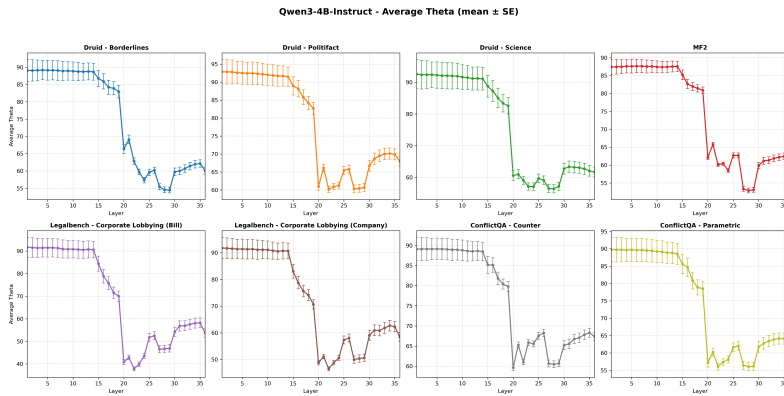
Figure 9: Correlation of θ and relative magnitude with Normalized Probability Differences across different models



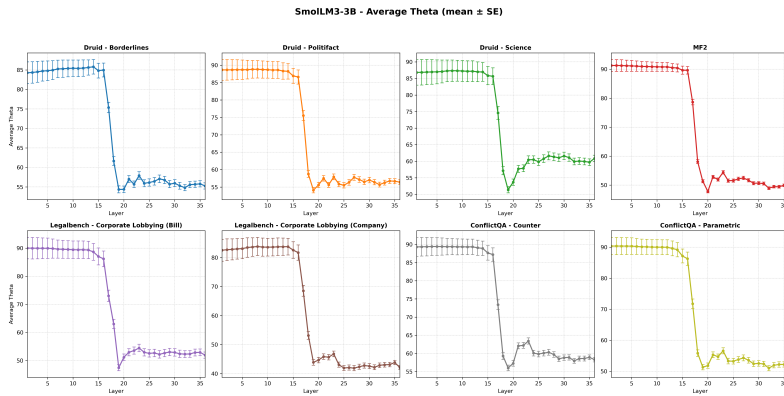
(a)



(b)

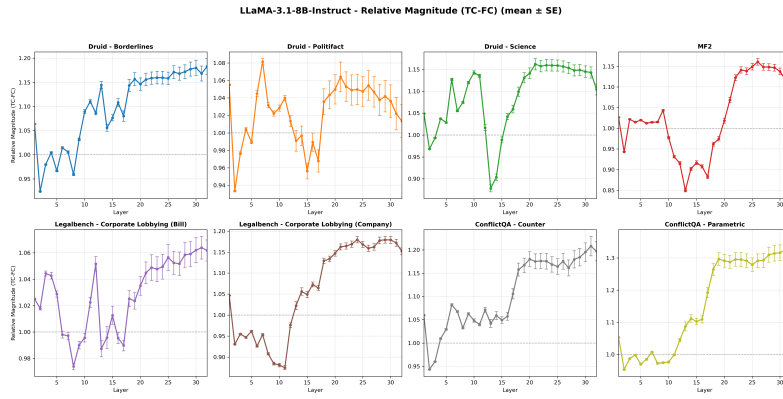


(c)

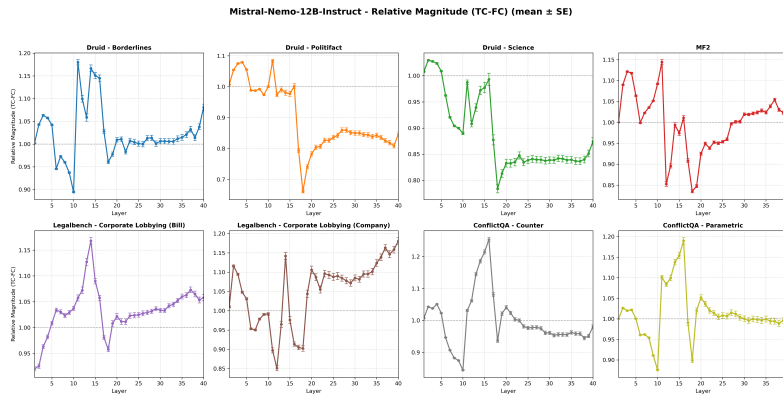


(d)

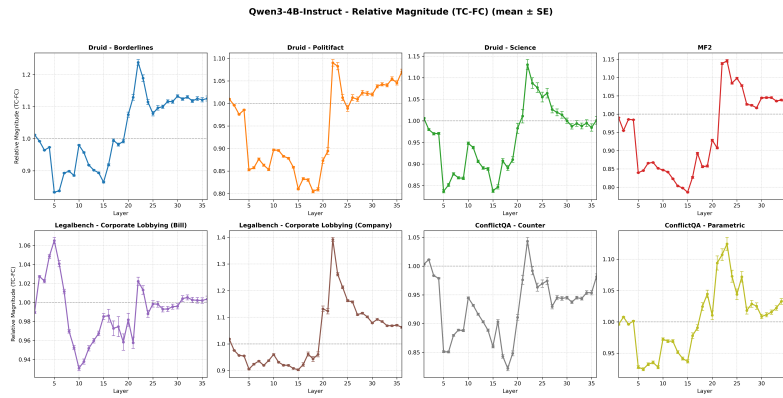
Figure 10: Layer wise plot of average θ in degrees across different models and datasets indicating the directional change in truth vectors when context is added. The error bars denote the standard error of mean



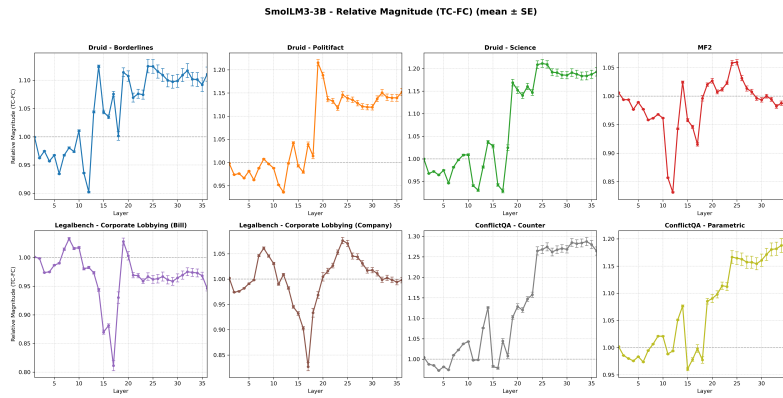
(a)



(b)



(c)



(d)

Figure 11: Layer wise plot of average relative magnitude across different models and datasets indicating the directional change in truth vectors when context is added. The error bars denote the standard error of mean