

Lost in Projection: The Geometric Orthogonality of LLM Assessment and Reasoning

Anonymous ACL submission

Abstract

Large language models (LLMs) frequently exhibit a dissociation between their internal confidence and actual reasoning competence. We investigate the mechanistic origin of this phenomenon by analyzing the geometry of residual stream activations across two phases: *pre-generative assessment* and *solution execution*. Using linear probing and principal component analysis across three model families (Llama, Qwen, Mistral), we identify two distinct geometric structures: a high-dimensional **Assessment Subspace** that encodes solvability beliefs, and a low-dimensional **Execution Subspace** that governs reasoning dynamics. While the belief state is robustly decodable during prompt processing, we observe a sharp **dimensionality shift** at the onset of generation, where the active variance transitions to the low-dimensional execution manifold. We validate this decoupling through causal intervention: steering vectors applied to the assessment subspace induce a decisive shift in the internal belief state ($\Delta > 0.8$), yet fail to alter downstream reasoning accuracy. Crucially, this inertness persists even on “borderline” tasks where the model possesses the requisite capability to solve the problem. These findings suggest a modular architecture where high-level assessment states are geometrically orthogonal to the procedural dynamics of execution, explaining why increasing internal confidence does not translate to improved competence.

1 Introduction

A critical challenge in the deployment of Large Language Models (LLMs) (Grattafiori et al., 2024; Jiang et al., 2023; Qwen et al., 2025) is their tendency to generate fluent, persuasive responses that are factually incorrect, often while maintaining high internal confidence (Maharana et al., 2025; Xiong et al., 2023; Ren et al., 2023). This dissociation between apparent confidence and actual

competence poses a fundamental reliability problem for applications requiring rigorous reasoning, such as scientific discovery and medical diagnostics. While prior research has documented this **confidence-competence gap** behaviorally (Singh et al., 2023) and noted the inconsistency of interventional methods (Tan et al., 2025; Braun et al., 2025), the mechanistic origin of this phenomenon remains unresolved (See Appendix A for related works). This paper addresses why a model’s internal assessment of a problem appears functionally decoupled from its subsequent reasoning process, providing a characterization based on the geometry of internal representations.

A prevailing hypothesis in mechanistic interpretability suggests that identifying the neural representation of a model’s “solvability belief” could allow for the direct control of its competence (Marks and Tegmark, 2024a; Dunefsky and Cohan, 2025). However, testing this hypothesis requires rigorous experimental control to isolate the semantic signal of belief from the high-dimensional noise of surface-level heuristics. We employ a controlled framework using length-balanced datasets of non-trivial reasoning problems across multiple model families (Qwen-2.5 (Qwen et al., 2025), Gemma-3 (Team et al., 2025), Mistral (Jiang et al., 2023)). By neutralizing confounding variables such as prompt length and lexical cues, we ensure that our probes capture a robust semantic representation of solvability rather than superficial correlates.

Our investigation proceeds in three stages. First, we identify a latent **Belief state**—an internal, pre-generative assessment of task solvability. Using linear probes and Centered Kernel Alignment (CKA) (Kornblith et al., 2019b; Zhou et al., 2024; Murphy et al., 2024), we demonstrate that this assessment is robustly and linearly encoded across diverse domains, including mathematics, coding, and logic. Second, we subject this state to a causal test. We intervene using targeted steering vectors to decisively

invert the model’s internal belief from *unsolvable* to *solvable* (inducing a belief shift of $\Delta > 0.8$) (Li et al., 2023a; Rinsky et al., 2023; Lee et al., 2024; Turner et al., 2023). We find that despite this significant internal shift, downstream reasoning accuracy remains statistically unchanged. Finally, to explain this functional decoupling, we analyze the geometry of the activation space. Using principal component analysis, we identify a distinct geometric transition (Park et al., 2023; Marks and Tegmark, 2023; Balestrieri et al., 2023; Konen et al., 2024) between the high-dimensional manifold of the **Assessment Subspace** (belief) and the low-dimensional manifold of the **Execution Subspace** (competence). We show that the transition between these phases involves a sharp dimensionality shift at the onset of generation, effectively isolating the initial assessment from the procedural dynamics of execution.

The identification of this dual-subspace architecture suggests a recalibration of strategies for AI reliability and safety. The implications of our findings are summarized as follows:

- ❶ **Limits of Confidence Steering:** Increasing a model’s internal confidence does not mechanically translate to improved reasoning reliability. The geometric orthogonality between assessment and execution indicates that belief states are not actionable control levers for procedural tasks.
- ❷ **Focus on Procedural Dynamics:** Addressing reasoning failures requires moving beyond high-level interventions on abstract traits (like “honesty”) to methods that directly target the low-dimensional dynamics of the execution subspace.
- ❸ **Mechanistic Auditing:** Conventional benchmarks focused solely on final answers provide an incomplete evaluation of model reliability. Future evaluations should independently audit the fidelity of the Assessment System and the robustness of the Execution System.
- ❹ **Resource Allocation:** While the Assessment signal is inert for control, it remains predictive. This offers a pathway for efficient AI, where early signals from the Assessment Subspace can trigger dynamic compute allocation or early exiting when the model predicts a high likelihood of failure.

2 Setup

To ensure our findings reflect general principles of LLM architecture rather than idiosyncratic model behaviors, we evaluate a diverse panel of instruction-tuned models from distinct organizations: **Qwen 2.5 7B** (Qwen et al., 2025), **Llama 3.2 3B** (Grattafiori et al., 2024), and **Mistral Small 24B** (Jiang et al., 2023). This selection allows us to validate the observed decoupling across varying parameter scales and design philosophies.

We subject these models to a suite of reasoning tasks spanning three distinct cognitive domains: (i) **Numerical Reasoning** (GSM8K (Cobbe et al., 2021), Math-Hard (Hendrycks et al., 2021), OpenR1 Math (Zhao et al., 2025)), (ii) **Formal Logical Deduction** (Knights and Knaves (Xie et al., 2024)), and (iii) **Algorithmic Planning** (Open R1 Coding (Zhao et al., 2025), QWQ-Planning (Hook, 2025)). This broad task coverage ensures that the dissociation between assessment and execution is a fundamental property of the models’ reasoning processes rather than an artifact of a specific domain. To maintain methodological control, we also include a suite of trivial, non-reasoning Control Tasks (e.g., verbatim copying and rote sequence retrieval) to establish baseline dimensionality metrics. All experiments were performed on a compute cluster of three NVIDIA A6000 GPUs; see Appendix C.2.1 for comprehensive dataset descriptions, prompt templates, and hardware specifications.

3 The Anatomy of Belief

We begin by investigating the existence of a latent assessment state prior to response generation. Specifically, we examine whether the internal representations at the final token of the prompt encode a linearly decodable signal predictive of task success. To isolate this “solvability belief” from surface-level confounds, we employ a controlled data curation protocol designed to ensure that the decoded signal reflects a semantic assessment rather than heuristic artifacts.

3.1 A Protocol for Confound-Resistant Data Curation

To isolate the semantic signal of solvability belief from spurious surface-level heuristics (Lysnæs-Larsen et al., 2025; Kumar et al., 2022), we implement a multi-stage filtering protocol on an initial corpus of 3,000 mathematical reasoning problems.

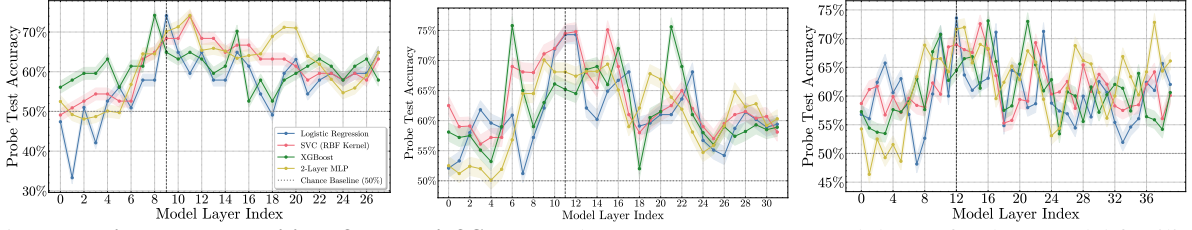


Figure 1: **Linear Decodability of the Belief State.** Probe accuracy across network layers for three model families (Left to Right: Qwen 2.5 7B, Llama 3.2 3B, Mistral Small 24B). Despite the increased expressive capacity of non-linear probes (SVC, XGBoost, MLP), they fail to outperform the linear baseline (Logistic Regression). This structural equivalence indicates that the solvability belief is encoded as a linear direction within the residual stream, emerging robustly in the mid-to-late layers.

This protocol sequentially (i) excludes structural format heuristics and keyword indicators that may signal task type, (ii) utilizes stratified topic balancing to ensure an identical distribution of reasoning sub-domains across classes, and (iii) performs distribution matching to neutralize prompt length as a confounding variable. The resulting curated dataset consists of 846 examples where token count distributions between the “solved” and “unsolved” classes are statistically indistinguishable ($p > 0.4$). By mathematically neutralizing these primary correlates, we ensure that the subsequent linear probes identify a latent semantic assessment state rather than superficial input artifacts; see Appendix C.3.1 for the complete curation protocol and representative examples.

3.2 Isolating and Characterizing the Latent Belief State

Consistent with probing protocols established for autoregressive models (Tillman and Mossing, 2025; Kantamneni et al., 2025; Yan, 2025), we isolate the pre-generative assessment state by extracting hidden representations at the final token of the input prompt. This locus captures the model’s fully contextualized representation of the problem immediately prior to the onset of the reasoning chain. To characterize the geometry of this representation, we investigate whether the solvability signal is encoded via a linear subspace or a complex non-linear manifold.

We address this by comparing the decoding performance of a linear probe against a hierarchy of non-linear classifiers. Specifically, we train a **Logistic Regression** probe to test for linear separability, and compare its performance to three non-linear baselines: a **Support Vector Classifier** with an RBF kernel, an **XGBoost** classifier, and a **2-Layer MLP** (hyperparameters detailed in Ap-

pendix C.4.1). This comparative framework serves as a controlled test of representational complexity: if the belief state relies on complex feature interactions, the non-linear models should significantly outperform the linear baseline.

The decoding results, presented in Figure 1, reveal two structural properties of the belief state. First, solvability is robustly decodable, with accuracy peaking between **70-75%** in the mid-to-late layers across all model families. Second, the signal is fundamentally linear. As shown in Figure 1, the capacity-rich non-linear probes yield no statistically significant performance improvement over the simple logistic regression model. This indicates that the information regarding task solvability is accessible via a linear readout (see Appendix D.1). A sample complexity analysis confirms that the solvability signal follows a non-linear emergence pattern, ruling out surface-level artifacts (Appendix E, Figure 7)

Formally, let $\mathcal{H} \subset \mathbb{R}^d$ denote the hidden state space at a given layer. For a problem i , let $h_i \in \mathcal{H}$ be its activation vector and $y_i \in \{0, 1\}$ its ground-truth solvability label. The equivalence in performance between linear and non-linear probes implies the existence of a specific direction $d_{\text{solv}} \in \mathbb{R}^d$ and bias b such that the probability of success is approximated by:

$$P(y_i = 1|h_i) \approx \sigma(h_i \cdot d_{\text{solv}} + b) \quad (1)$$

While the signal is linearly separable, the imperfect accuracy ($\sim 75\%$) suggests that the projected distributions $P(h_i \cdot d_{\text{solv}}|y_i = 1)$ and $P(h_i \cdot d_{\text{solv}}|y_i = 0)$ maintain non-trivial overlap. This characterizes the belief state as a linear direction embedded within a high-dimensional, noisy manifold, rather than a perfectly disentangled feature.

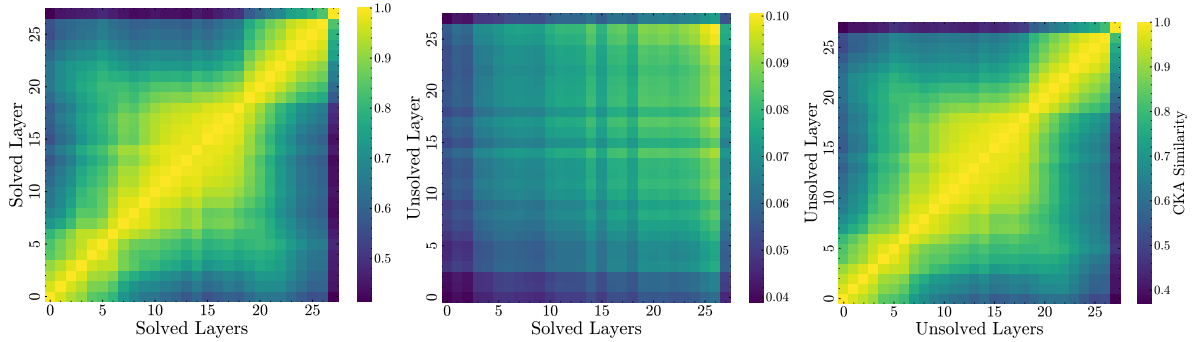


Figure 2: **Geometric Coherence and Dissimilarity of Belief States.** Using Centered Kernel Alignment (CKA), we compare the representational geometry of belief states. **(Left & Right)** The high self-similarity along the diagonals confirms that both “Solved” and “Unsolved” states are internally coherent, geometrically stable representations across layers. **(Center)** In stark contrast, the cross-comparison reveals a profound geometric dissimilarity, with near-zero CKA scores between the two states. This provides strong evidence that the model represents belief not as a single continuum, but as two distinct and fundamentally separate geometric objects.

Table 1: **Probe Accuracy by Token Position.** We trained our suite of probes on activations extracted from different positions within the input prompt to identify the locus of the pre-generative belief signal. Accuracy is reported on the held-out test set of our length-controlled math dataset. For all probe architectures, performance consistently and decisively peaks at the **Last Input Token** (“t_question_end”), indicating that the model’s solvability belief crystallizes at the moment it has finished processing the full problem context.

Activation Locus	Logistic Reg.	SVC	XGBoost	MLP
10% of Prompt	54.2% \pm 2.1	55.1% \pm 2.3	53.8% \pm 2.5	52.9% \pm 2.8
50% of Prompt	65.1% \pm 1.8	66.3% \pm 1.9	65.9% \pm 2.0	66.8% \pm 2.1
Last Input Token	74.4% \pm 1.5	74.6% \pm 1.6	75.1% \pm 1.7	75.2% \pm 1.8
EOS Token	67.9% \pm 1.6	68.5% \pm 1.7	68.1% \pm 1.8	68.8% \pm 1.9

3.3 Geometry Visualization and Structural Similarity

To validate that the statistical separability observed via linear probing corresponds to a robust topological feature of the activation manifold, we visualize the high-dimensional hidden states using non-linear dimensionality reduction. We employ t-distributed Stochastic Neighbor Embedding (t-SNE) to examine local structural preservation and Uniform Manifold Approximation and Projection (UMAP) to assess global geometry.

As illustrated in Figure 3, which visualizes $N = 800$ data points, both projection methods reveal a clear geometric separation. The activation vectors resolve into disjoint clusters rather than a diffuse cloud. We note that the visual density of the clusters—appearing as compact groups despite the large sample size—results from significant overplotting, reflecting the high topological cohesion of the respective manifolds. This confirms that the solvability assessment is encoded in a distinct region of the activation space, consistent with the linear separability results.

To quantify the internal coherence and mutual distinctness of these representations across layers, we utilize Centered Kernel Alignment (CKA) (Kornblith et al., 2019a), a robust metric for comparing representational geometries. The resulting similarity heatmaps (Figure 2) demonstrate three structural properties:

- High Intra-Class Coherence:** The self-similarity matrices for both “Solved” and “Unsolved” states exhibit high diagonal values, indicating that these representations evolve consistently and stably through the network layers.
- Inter-Class Orthogonality:** The cross-comparison between “Solved” and “Unsolved” layers yields low CKA scores, confirming that these states occupy geometrically distinct subspaces.

These analyses confirm the existence of a stable, geometrically distinct belief manifold. However, geometric separability does not inherently imply

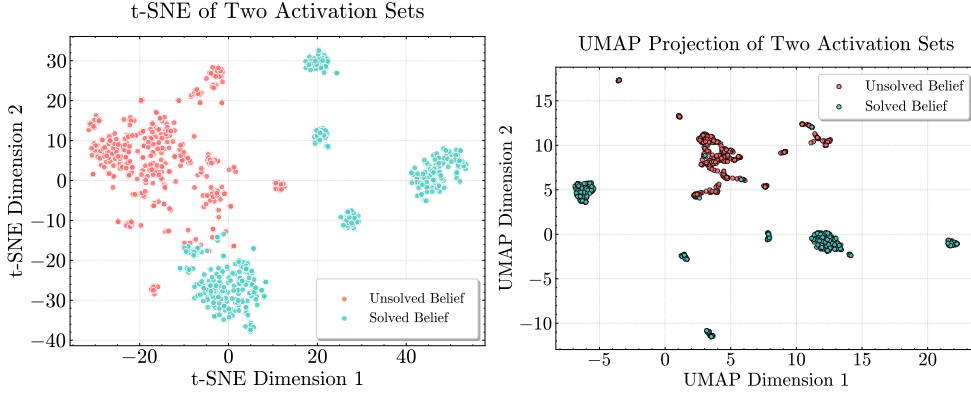


Figure 3: **Geometric Separation of Belief States.** Visualization of high-dimensional belief representations projected via t-SNE (Left) and UMAP (Right) for Qwen 2.5 7B (Layer 26, $N = 846$). Both manifold learning techniques demonstrate a distinct topological separation between “Unsolved” (coral) and “Solved” (cyan) classes. The emergence of well-defined clusters confirms that the statistical separability observed in Section 3.2 corresponds to a robust geometric structure within the activation space, rather than an artifact of high-dimensional sparsity.

functional utility. Having established the structural properties of this state, we next investigate its causal role in the reasoning process.

4 Causal Decoupling

4.1 Establishing Causal Control over Latent Belief

To investigate the causal role of the belief state, we first establish a method to reliably manipulate it. Based on the finding in Section 3.2 that the belief signal is fundamentally linear, we derive a steering vector, d_{solv} , directly from the weights of the trained logistic regression probe. This vector represents the direction normal to the decision boundary separating “solved” from “unsolved” representations in the activation space.

We determine the optimal locus for intervention by analyzing probe accuracy across the model’s processing trajectory. As reported in Table 1, decodability is low in early processing stages but increases as the context is integrated, peaking at the final token of the input prompt, similar to approaches from Tillman and Mossing (2025); Kantamneni et al. (2025); Yan (2025). This indicates that the belief state is maximally crystallized immediately prior to generation. Consequently, we target our intervention at this specific token position.

The intervention is defined as an additive update to the residual stream. For a given hidden state $h_{i,L}$ corresponding to problem i at the last input token of layer L (Gurnee and Tegmark, 2024), the steered state $h'_{i,L}$ is computed as:

$$h'_{i,L} = h_{i,L} + \alpha \cdot d_{\text{solv}} \quad (2)$$

where α is a scalar coefficient controlling the magnitude and direction of the shift, tuned with leave-out cross validation (Hastie et al., 2009).

To validate the efficacy of d_{solv} , we conducted a manipulation check on a held-out set of problems classified by the model as “unsolved” (baseline probability ≈ 0). We applied the steering vector with a positive α to shift the activations toward the “solved” region of the manifold. As detailed in Table 2, this intervention produces a robust inversion of the internal belief state. In the Math-Hard domain, the probe’s predicted probability of success rises from a baseline of 0.04 to 0.97 (a belief shift of $\Delta = +0.93$). This controllability generalizes across domains, with significant belief shifts observed for Logic (+0.84), Code (+0.88), and Planning (+0.79). These results confirm that d_{solv} functions as a reliable control lever for the model’s internal assessment state.

4.2 Causal Inertness of Latent Belief

Having established control over the internal belief state, we evaluate its causal influence on task performance. We apply the steering intervention to invert the model’s assessment from *unsolvable* to *solvable* and measure the resulting accuracy on the held-out test set. This experiment tests the hypothesis that the linearly decodable belief state functions as a control mechanism for reasoning competence.

Results for the Math-Hard dataset (Table 2) demonstrate a sharp dissociation between internal representation and behavioral outcome. The intervention successfully shifts the probe’s predicted probability of success from a baseline of 0.37 to 0.97, yielding a substantial belief shift

Table 2: **Causal Intervention Reveals a Decoupling of Latent Belief and Task Competence.** We apply our validated ‘d_solv’ steering vector (derived from a respective datasets) to held-out problems across four diverse reasoning domains. The intervention successfully and dramatically flips the internal belief state (Probe’s Prediction) from unconfident to confident. However, this manipulation of belief has **no statistically significant effect** on the final task accuracy in any domain, providing powerful causal evidence for the decoupling of the two systems.

Dataset	Intervention	Internal Belief State		Final Task Outcome		
		Probe’s Pred.	Belief Flip (Δ)	Task Acc. (%)	Perf. Change (Δ)	<i>p</i> -value
Math-HARD	Baseline (No Steer)	0.04 ± 0.02	—	37.4 ± 0.3	—	0.981
	Steer \rightarrow "Solved"	0.97 ± 0.03	+0.93	37.4 ± 0.8	0.0 ± 0.5	
Knights & Knaves	Baseline (No Steer)	0.11 ± 0.04	—	42.6 ± 1.2	—	0.952
	Steer \rightarrow "Solved"	0.95 ± 0.05	+0.84	42.8 ± 0.9	0.2 ± 0.3	
OpenR1 Coding	Baseline (No Steer)	0.08 ± 0.03	—	62.9 ± 0.7	—	0.991
	Steer \rightarrow "Solved"	0.96 ± 0.04	+0.88	62.8 ± 1.3	-0.1 ± 0.4	
QwQ Planning	Baseline (No Steer)	0.15 ± 0.06	—	38.1 ± 1.4	—	0.913
	Steer \rightarrow "Solved"	0.94 ± 0.07	+0.79	38.1 ± 0.7	0.0 ± 0.7	

($\Delta = +0.60$). However, this internal modulation does not translate to behavioral improvement: task accuracy remains statistically unchanged (Performance Change: $0.0\% \pm 0.5$, $p = 0.981$). To address the concern that this inertness stems from the model’s fundamental inability to solve hard problems (baseline accuracy 8.4%), we replicated this evaluation on “borderline” difficulty tasks (GSM-Hard) where the model demonstrates significant baseline competence ($\sim 55\%$). As detailed in Appendix D.3, the causal inertness persists even in this regime of adequate capability. Furthermore, layer-wise causal tracing (Appendix D.1) confirms that this lack of effect is not an artifact of the specific intervention layer, but a global property of the architecture. This decoupling generalizes across all tested reasoning domains. In logic tasks (**Knights & Knaves**), we induce a belief shift of $\Delta = +0.84$, yet task accuracy remains stable ($+0.2\% \pm 0.3$, $p = 0.952$). Similarly, in coding tasks (**OpenR1**), a shift of $\Delta = +0.88$ yields no significant performance change ($-0.1\% \pm 0.4$, $p = 0.991$). We validated that these vectors are semantically specific to their respective domains (Appendix F), ruling out the possibility that the steering effect is merely a generic, task-agnostic signal. Notably, we also did not observe any change in the tone of responses after steering – the model did not “sound” any more confident than it usually is after the steering. These results indicate that while the belief state is malleable via linear steering, it is functionally decoupled from the downstream mechanisms governing solution generation.

To verify that this inertness is not an artifact of the specific intervention direction (unsolvable \rightarrow

solvable), we performed the inverse experiment: steering *solved* problems to be represented as *unsolvable*. As detailed in Appendix D.2, this yielded a symmetric null result; inducing a representation of low confidence did not degrade performance on tasks the model was capable of solving. This confirms that the decoupling is a bidirectional property of the architecture (Appendix D.2, Table 8).

The observed causal inertness challenges the assumption that linearly decodable states necessarily function as active control variables. The lack of downstream effect suggests a structural separation between the representation of assessment and the execution of reasoning. To elucidate the mechanism enabling this separation, we next examine the geometric properties of the activation space.

4.3 Dimensionality of Belief vs. Action

We characterize the intrinsic dimensionality of the representations by analyzing the cumulative explained variance of their principal components. We assembled four distinct sets of activations: pre-generative “Belief” states (classified by high vs. low confidence) and in-process “Competence” states (classified by successful vs. failed reasoning traces).

As shown in Figure 4 (Left), the variance profiles exhibit a distinct structural divergence. The “Competence” curves (Green, Orange) demonstrate rapid saturation, where a small number of principal components (≈ 20) capture the majority of the variance. This indicates that the reasoning process evolves on a low-rank, constrained manifold. In contrast, the “Belief” curves (Blue, Red) exhibit a slow, gradual saturation, requiring over 120 compo-

Table 3: **Quantitative Evidence for the Geometric Decoupling of Assessment and Execution.** We measure the complexity of key cognitive subspaces using the Participation Ratio (PR), which computes an ‘effective dimensionality’. The results reveal an order-of-magnitude difference in complexity between the systems. Furthermore, an asymmetry is revealed within the Assessment system, where ‘unconfident’ states are significantly more complex than ‘confident’ ones. Values are reported as mean \pm standard deviation over 100 bootstrap resamples.

Cognitive System	Subspace Representing...	Participation Ratio
Assessment (<i>Pre-Generative Belief</i>)	“Confident” (Positive Belief)	33.6 ± 2.9
	“Unconfident” (Negative Belief)	44.4 ± 2.5
Execution (<i>In-Process Reasoning</i>)	Competent (Successful CoT)	16.0 ± 0.6
	Incompetent (Failed CoT)	17.9 ± 0.9
Control (<i>Trivial Tasks</i>)	Trivial Execution	4.1 ± 0.9

nents to explain 90% of the variance. This characterizes the pre-generative assessment as a diffuse, high-dimensional representation.

To quantify this geometric disparity, we calculate the **Participation Ratio (PR)**, a measure of effective dimensionality derived from the eigenvalues (λ_i) of the covariance matrix:

$$PR = \frac{(\sum_i \lambda_i)^2}{\sum_i \lambda_i^2} \quad (3)$$

This metric estimates the number of dimensions significantly contributing to the data’s variance. As detailed in Table 3, the “Assessment” system exhibits high effective dimensionality, with a PR of **33.6** for positive belief and **44.4** for negative belief. Conversely, the “Execution” system operates on a manifold that is approximately twice as compressed, with a PR of **16.0** for competent traces and **17.9** for incompetent ones. See Table 11 for an extended version.

To determine whether this dimensionality collapse is a specific feature of procedural reasoning or merely a general artifact of the transition from prompt ingestion to token generation, we introduce the Control baseline (See Appendix C.3 for types of control). For trivial tasks requiring no reasoning (e.g., verbatim copying or rote sequence retrieval), the execution manifold collapses almost entirely to a PR of 4.1 ± 0.9 . This tripartite hierarchy where PR scales with the complexity of the task provides two critical insights. First, it confirms that the Execution Subspace is a dynamic procedural workspace rather than a fixed architectural bottleneck; its dimensionality is a function of the computational load. Second, it provides a mechanistic explanation for the observed causal

decoupling. The Assessment Subspace encodes solvability as a diffuse, high-dimensional feature, whereas the Execution Subspace is confined to a narrow, low-rank trajectory.

The inefficacy of the steering intervention (Section 4.2) suggests that the belief vector resides in a region of the activation space that is geometrically misaligned with the manifolds governing execution. This causal inertness, paired with near-zero CKA scores and minimal subspace fit, provides empirical evidence of a *functional orthogonality* between the high-dimensional assessment manifold and the constrained execution trajectory. Consequently, linear perturbations in the high-dimensional assessment space fail to project meaningfully onto the constrained, low-rank dynamics of the reasoning process. This structural mismatch confirms that confidence and competence operate in distinct geometric regimes of the residual stream, effectively isolating the model’s internal judgment from its procedural output.

4.4 Visualizing the Geometric Transition

The static analysis in Section 4.3 establishes that belief and competence occupy manifolds of distinct intrinsic dimensionalities. To characterize the temporal dynamics of this separation, we analyze the evolution of activation states throughout the generation process.

We employ a **Trajectory Projection** analysis. First, we construct orthonormal bases for each system: the high-dimensional **Assessment Basis** (B_{assess}), derived from the principal components of pre-generative belief states, and the low-dimensional **Execution Basis** (B_{exec}), derived from in-process reasoning traces. We then track the acti-

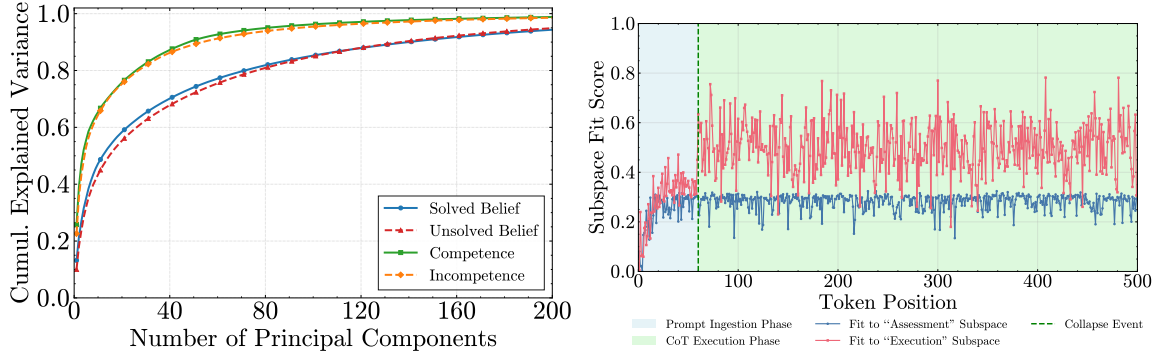


Figure 4: **Geometric Characterization of the Dissociation.** (Left) Cumulative explained variance analysis via PCA. The slow saturation of the Belief curves (Blue, Red) indicates that the Assessment subspace forms a high-dimensional manifold. In contrast, the rapid saturation of the Competence curves (Green, Orange) reveals that the Execution subspace possesses a low-rank, constrained structure. (Right) Temporal projection of hidden states onto the identified subspace bases. During prompt ingestion, the activation trajectory projects strongly onto the Assessment basis (Blue). At the onset of generation (dashed line), we observe a sharp **dimensionality shift**: the projection onto the Assessment subspace diminishes, while the Execution subspace (Red) captures the dominant variance share, marking the transition to procedural reasoning.

504 vation trajectory $H = [h_1, \dots, h_n]$ for a single inference pass. To ensure consistency with our static analysis, activations are extracted from the layer exhibiting peak belief decodability (as identified in Section 4.1). The sequence H encompasses the full computational span, where $n = n_{\text{prompt}} + n_{\text{response}}$, capturing both the ingestion of the input problem and the subsequent generation of the reasoning chain. At each token step t , we quantify the alignment of the activation vector h_t with each subspace using the **Subspace Fit** metric:

$$515 \quad \text{Subspace Fit}(h_t, B) = \frac{\|\text{proj}_B(h_t)\|^2}{\|h_t\|^2} \quad (4)$$

516 This metric provides a normalized measure of geometric containment, where a value of 1 indicates the state lies entirely within the target subspace and 0 indicates orthogonality.

520 As illustrated in Figure 4 (Right), the projection dynamics reveal a distinct phase separation. During the **Prompt Ingestion Phase**, the activation trajectory is predominantly explained by the Assessment basis (Blue line), reflecting the accumulation of solvability features. However, at the onset of the **CoT Execution Phase** (first generated token), we observe a distinct **dimensionality shift**. The fit to the Assessment subspace diminishes rapidly, while the fit to the Execution subspace (Red line) rises to near-unity. This discontinuity provides empirical evidence that the Assessment and Execution systems operate as sequentially engaged, geometrically distinct regimes, with the model transitioning from a high-dimensional evaluation state to a low-dimensional procedural manifold upon generation.

5 Conclusion 536

537 We investigated the mechanistic origins of the confidence-competence dissociation in Large Language Models. By isolating a linearly decodable “solvability belief,” we established that models possess a robust internal assessment of task difficulty that is nonetheless causally inert with respect to downstream performance. We provide a geometric explanation for this decoupling: the transition from prompt processing to token generation entails a sharp **dimensionality shift**, where the active representation moves from a high-dimensional Assessment Subspace to a constrained, low-rank Execution Subspace. *dual-subspace architecture* elucidates why internal confidence estimates do not function as control variables for deductive reasoning. Consequently, these findings suggest that future reliability interventions should deprioritize global belief steering in favor of methods that directly engage the low-dimensional dynamics of the execution process. 556

557 Limitations and Future Work

558 **Scope of Task Domains.** Our investigation is intentionally restricted to convergent reasoning domains (mathematics, logic, coding, and planning) where task competence is defined by objective ground-truth criteria. In these settings, the Execution Subspace is highly constrained by rigid logical rules, resulting in the low effective dimensionality we observed. It remains an open question whether this geometric decoupling persists in divergent or cre- 566

567 ative generation tasks (e.g., creative writing, open- 617
568 ended dialogue), where the valid output manifold 618
569 may possess significantly higher intrinsic dimen- 619
570 sionality. Future research should investigate if the 620
571 dimensionality shift is a universal property of au- 621
572 toregressive generation or specific to deductive rea- 622
573 soning. 623

574 **Model Scale and Frontier Models.** Due to 624
575 computational resource constraints, our experimen- 625
576 tal evaluation was restricted to models within the 626
577 3B to 24B parameter range. While the geometric 627
578 consistency observed across three distinct model 628
579 families (Llama, Qwen, Mistral) suggests that the 629
580 dual-subspace architecture is a robust feature of 630
581 current instruction-tuned LLMs, it remains possi- 631
582 ble that frontier-scale models (e.g., >70B param- 632
583 eters) exhibit different representational dynamics. 633
584 Specifically, increased parameter density may al- 634
585 low for a more integrated coupling between assess-
586 ment and execution states. Validating these geomet-
587 ric findings on larger-scale models is a necessary
588 step for ensuring the universality of the confidence-
589 competence gap.

590 **Developmental Origin of Decoupling.** While 635
591 we characterize the geometric architecture of cur- 636
592 rent instruction-tuned models, this study does not 637
593 isolate the training phase responsible for the ob- 638
594 served orthogonality. Specifically, it remains to be 639
595 determined whether the dissociation between the 640
596 Assessment and Execution subspaces emerges dur- 641
597 ing pre-training as a consequence of next-token 642
598 prediction, or if it is exacerbated by alignment 643
599 techniques such as Reinforcement Learning from 644
600 Human Feedback (RLHF), which may inadver-
601 tently reward surface-level confidence markers in-
602 dependent of procedural accuracy. A longitudinal
603 analysis of the belief state throughout the training
604 pipeline would clarify whether this architecture is
605 inherent or induced.

606 Ethical considerations

607 **Risks of Artificial Confidence Inflation.** Our find- 645
608 ings demonstrate the feasibility of dissociating a 646
609 model’s internal confidence from its reasoning ca- 647
610 pabilities via linear steering. While our objective 648
611 is to understand cognitive architecture, this mechan- 649
612 ics presents a dual-use risk. Malicious actors could 650
613 potentially utilize similar steering vectors to arti- 651
614 ficially inflate the epistemic certainty of a model, 652
615 causing it to present factually incorrect or halluci- 653
616 natory information with high persuasive authority. 654

This underscores the necessity of evaluating model 617
reliability based on external behavioral verification 618
rather than internal activation states alone. 619

Limitations of Internal Safety Monitoring. 620
We propose that the Assessment Subspace could 621
be utilized for dynamic resource allocation (Lin 622
et al., 2024; Behari et al., 2024) or quality filter- 623
ing. However, the observed orthogonality between 624
assessment and execution suggests a potential fail- 625
ure mode for safety systems. If safety guardrails 626
or refusal mechanisms rely on similar high-level 627
assessment manifolds, they may fail to detect gran- 628
ular failures or harmful outputs that arise strictly 629
within the low-dimensional dynamics of the Ex- 630
ecution Subspace. Consequently, safety auditing 631
should not rely exclusively on high-level represen- 632
tation probing but must remain grounded in the 633
verification of final outputs. 634

635 References

- 636 Amos Azaria and Tom Mitchell. 2023. [The internal 637](#)
638 [state of an LLM knows when it’s lying](#). In *Find- 639*
640 *ings of the Association for Computational Linguistics: 640*
EMNLP 2023, pages 967–976, Singapore. Associa-
641 tion for Computational Linguistics.
- 642 Waïss Azizian, Michael Kirchhof, Eugene Ndiaye,
643 Louis Béthune, Michal Klein, Pierre Ablin, and
644 Marco Cuturi. 2025. [The geometries of truth are 644](#)
[orthogonal across tasks](#). *ArXiv*, abs/2506.08572.
- 645 Randall Balestriero, Romain Cosentino, and Sarath
646 Shekkizhar. 2023. [Characterizing large language 646](#)
647 [model geometry helps solve toxicity detection and 647](#)
648 [generation](#). In *International Conference on Machine 648*
649 *Learning*.
- 650 Nikhil Behari, Edwin Zhang, YUNFAN ZHAO, Aparna
651 Taneja, Dheeraj Mysore Nagaraj, and Milind Tambe.
652 2024. [A decision-language model \(DLM\) for dy- 652](#)
653 [namic restless multi-armed bandit tasks in public 653](#)
654 [health](#). In *The Thirty-eighth Annual Conference on 654*
655 *Neural Information Processing Systems*.
- 656 Joschka Braun, Carsten Eickhoff, David Krueger,
657 Seyed Ali Bahrainian, and Dmitrii Krasheninnikov.
658 2025. [Understanding \(un\)reliability of steering vec- 658](#)
659 [tors in language models](#). *ArXiv*, abs/2505.22637.
- 660 Yuanpu Cao, Tianrong Zhang, Bochuan Cao, Ziyi Yin,
661 Lu Lin, Fenglong Ma, and Jinghui Chen. 2024. [Per- 661](#)
662 [sonalized steering of large language models: Versa- 662](#)
663 [tile steering vectors through bi-directional preference 663](#)
664 [optimization](#). In *The Thirty-eighth Annual Confer- 664*
665 *ence on Neural Information Processing Systems*.
- 666 Sviatoslav Chalnev, Matthew Siu, and Arthur Conmy.
667 2024. [Improving steering vectors by targeting sparse 667](#)
668 [autoencoder features](#). *Preprint*, arXiv:2411.02193.

669	Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian,	Subhash Kantamneni, Joshua Engels, Senthoran Ra-	725
670	Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias	jamanoharan, Max Tegmark, and Neel Nanda. 2025.	726
671	Plappert, Jerry Tworek, Jacob Hilton, Reiichiro	Are sparse autoencoders useful? a case study in	727
672	Nakano, Christopher Hesse, and John Schulman.	sparse probing . In <i>Forty-second International Con-</i>	728
673	2021. Training verifiers to solve math word prob-	ference on Machine Learning .	729
674	lems . <i>Preprint</i> , arXiv:2110.14168.		
675	Gheorghe Comanici, Eric Bieber, Mike Schaekermann,	Kai Konen, Sophie Jentzsch, Diaoulé Diallo, Peer	730
676	Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-	Schutt, Oliver Bensch, Roxanne El Baff, Dominik	731
677	cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, Luke	Opitz, and Tobias Hecking. 2024. Style vectors for	732
678	Marris, Sam Petulla, Colin Gaffney, Asaf Aharoni,	steering generative large language models . In <i>Find-</i>	733
679	Nathan Lintz, Tiago Cardal Pais, Henrik Jacobs-	ings .	734
680	son, Idan Szpektor, Nan-Jiang Jiang, and 3416 oth-	Simon Kornblith, Mohammad Norouzi, Honglak Lee,	735
681	ers. 2025. Gemini 2.5: Pushing the frontier with	and Geoffrey Hinton. 2019a. Similarity of neural	736
682	advanced reasoning, multimodality, long context,	network representations revisited . <i>Preprint</i> ,	737
683	and next generation agentic capabilities . <i>Preprint</i> ,	arXiv:1905.00414.	738
684	arXiv:2507.06261.		
685	Jacob Dunefsky and Arman Cohan. 2025. One-shot	Simon Kornblith, Mohammad Norouzi, Honglak Lee,	739
686	optimized steering vectors mediate safety-relevant	and Geoffrey E. Hinton. 2019b. Similarity of	740
687	behaviors in llms . In <i>unknown</i> .	neural network representations revisited . <i>ArXiv</i> ,	741
688		abs/1905.00414.	742
689	Google DeepMind. 2025. Gemini 3 Pro. https://deepmind.google/models/gemini/pro/ .	Abhinav Kumar, Chenhao Tan, and Amit Sharma. 2022.	743
690	Accessed: 2026-01-04.	Probing classifiers are unreliable for concept removal	744
691	Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri,	and detection . In <i>Advances in Neural Information</i>	745
692	Abhinav Pandey, Abhishek Kadian, Ahmad Al-	Processing Systems .	746
693	Dahle, Aiesha Letman, Akhil Mathur, Alan Schel-	Bruce W. Lee, Inkit Padhi, K. Ramamurthy, Erik	747
694	ten, Alex Vaughan, Amy Yang, Angela Fan, Anirudh	Miehling, Pierre L. Dognin, Manish Nagireddy,	748
695	Goyal, Anthony Hartshorn, Aobo Yang, Archi Mi-	and Amit Dhurandhar. 2024. Programming re-	749
696	tra, Archie Sravankumar, Artem Korenev, Arthur	fusal with conditional activation steering . <i>ArXiv</i> ,	750
697	Hinsvark, and 542 others. 2024. The llama 3 herd of	abs/2409.05907.	751
698	models . <i>Preprint</i> , arXiv:2407.21783.		
699	Wes Gurnee and Max Tegmark. 2024. Language mod-	Kenneth Li, Oam Patel, Fernanda Vi’egas, H. Pfister,	752
700	els represent space and time . In <i>The Twelfth Interna-</i>	and M. Wattenberg. 2023a. Inference-time inter-	753
701	tional Conference on Learning Representations .	vention: Eliciting truthful answers from a language	754
702	T. Hastie, R. Tibshirani, and J. Friedman. 2009. The el-	model . <i>ArXiv</i> , abs/2306.03341.	755
703	ements of statistical learning: data mining, inference	Kenneth Li, Oam Patel, Fernanda Viégas, Hanspeter	756
704	and prediction , 2 edition. Springer.	Pfister, and Martin Wattenberg. 2023b. Inference-	757
705	Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul	time intervention: Eliciting truthful answers from a	758
706	Arora, Steven Basart, Eric Tang, Dawn Song, and Ja-	language model . In <i>Thirty-seventh Conference on</i>	759
707	cob Steinhardt. 2021. Measuring mathematical prob-	Neural Information Processing Systems .	760
708	lem solving with the math dataset. <i>arXiv preprint</i>	Bin Lin, Tao Peng, Chen Zhang, Minmin Sun, Lanbo	761
709	<i>arXiv:2103.03874</i> .	Li, Hanyu Zhao, Wencong Xiao, Qi Xu, Xiafei Qiu,	762
710	Dmitrii Hook. 2025. qwq-32b-planning-mystery-17-	Shen Li, Zhigang Ji, Yong Li, and Wei Lin. 2024.	763
711	24k-greedy . Hugging Face Dataset.	Infinite-llm: Efficient llm service for long context	764
712	Li Ji-An, Hua-Dong Xiong, Robert C. Wilson,	with distattention and distributed kvcache . <i>ArXiv</i> ,	765
713	Marcelo G. Mattar, and Marcus K. Benna. 2025. Lan-	abs/2401.02669.	766
714	guage models are capable of metacognitive monitor-	Junteng Liu, Shiqi Chen, Yu Cheng, and Junxian He.	767
715	ing and control of their internal activations . <i>Preprint</i> ,	2024a. On the universal truthfulness hyperplane in-	768
716	arXiv:2505.13763.	side LLMs . In <i>Proceedings of the 2024 Conference</i>	769
717	Albert Q. Jiang, Alexandre Sablayrolles, Arthur Men-	on Empirical Methods in Natural Language Process-	770
718	sch, Chris Bamford, Devendra Singh Chaplot, Diego	ing , pages 18199–18224, Miami, Florida, USA. As-	771
719	de las Casas, Florian Bressand, Gianna Lengyel, Guil-	sociation for Computational Linguistics.	772
720	laume Lample, Lucile Saulnier, L��lio Renard Lavaud,	Zhenhua Liu, Tong Zhu, Chuanyuan Tan, Bing Liu,	773
721	Marie-Anne Lachaux, Pierre Stock, Teven Le Scao,	Haonan Lu, and Wenliang Chen. 2024b. Probing	774
722	Thibaut Lavril, Thomas Wang, Timoth��e Lacroix,	language models for pre-training data detection . In	775
723	and William El Sayed. 2023. Mistral 7b . <i>Preprint</i> ,	Proceedings of the 62nd Annual Meeting of the As-	776
724	arXiv:2310.06825.	sociation for Computational Linguistics (Volume 1:	777
		Long Papers) , pages 1576–1587, Bangkok, Thailand.	778
		Association for Computational Linguistics.	779

780	Jacob Lysnæs-Larsen, Marte Eggen, and Inga Strümke.	Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong,	833
781	2025. Probing the probes: Methods and metrics for	Evan Hubinger, and Alexander Turner. 2024. Steering	834
782	concept alignment . <i>ArXiv</i> , abs/2511.04312.	llama 2 via contrastive activation addition . In	835
		<i>Proceedings of the 62nd Annual Meeting of the As-</i>	836
783	Umakanta Maharana, Sarthak Verma, Avarna Agar-	<i>sociation for Computational Linguistics (Volume 1:</i>	837
784	wal, Prakashini Mruthyunjaya, Dwarikanath Ma-	<i>Long Papers)</i> , pages 15504–15522, Bangkok, Thai-	838
785	hapatra, Sakir Ahmed, and Murari Mandal. 2025.	land. Association for Computational Linguistics.	839
786	Right prediction, wrong reasoning: Uncovering llm		
787	misalignment in ra disease diagnosis . <i>Preprint</i> ,	Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong,	840
788	arXiv:2504.06581.	Evan Hubinger, and Alexander Matt Turner. 2023.	841
		Steering llama 2 via contrastive activation addition .	842
		<i>ArXiv</i> , abs/2312.06681.	843
789	Samuel Marks and Max Tegmark. 2023. The geometry	Aniket Kumar Singh, Suman Devkota, Bishal Lamich-	844
790	of truth: Emergent linear structure in large language	hane, Uttam Dhakal, and Chandra Dhakal. 2023. The	845
791	model representations of true/false datasets . <i>ArXiv</i> ,	confidence-competence gap in large language mod-	846
792	abs/2310.06824.	els: A cognitive study . <i>Preprint</i> , arXiv:2309.16145.	847
793	Samuel Marks and Max Tegmark. 2024a. The geometry	Daniel Tan, David Chanin, Aengus Lynch, Dimitrios	848
794	of truth: Emergent linear structure in large language	Kanoulas, Brooks Paige, Adria Garriga-Alonso,	849
795	model representations of true/false datasets . <i>Preprint</i> ,	and Robert Kirk. 2025. Analyzing the generaliza-	850
796	arXiv:2310.06824.	tion and reliability of steering vectors . <i>Preprint</i> ,	851
		arXiv:2407.12404.	852
797	Samuel Marks and Max Tegmark. 2024b. The geometry	Daniel Chee Hian Tan, David Chanin, Aengus Lynch,	853
798	of truth: Emergent linear structure in large language	Brooks Paige, Dimitrios Kanoulas, Adria Garriga-	854
799	model representations of true/false datasets . <i>Preprint</i> ,	Alonso, and Robert Kirk. 2024. Analysing the gen-	855
800	arXiv:2310.06824.	eralisation and reliability of steering vectors . In <i>The</i>	856
		<i>Thirty-eighth Annual Conference on Neural Informa-</i>	857
801	Alex Murphy, J. Zylberberg, and Alona Fyshe. 2024.	<i>tion Processing Systems</i> .	858
802	Correcting biased centered kernel alignment mea-		
803	sures in biological and artificial neural networks .	Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya	859
804	<i>ArXiv</i> , abs/2405.01012.	Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin,	860
		Tatiana Matejovicova, Alexandre Ramé, Morgane	861
805	Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Re-	Rivière, Louis Rouillard, Thomas Mesnard, Geoffrey	862
806	ichart, Idan Szpektor, Hadas Kotek, and Yonatan Be-	Cideron, Jean bastien Grill, Sabela Ramos, Edouard	863
807	linkov. 2025a. LLMs know more than they show: On	Yvinec, Michelle Casbon, Etienne Pot, Ivo Penchev,	864
808	the intrinsic representation of LLM hallucinations . In	and 197 others. 2025. Gemma 3 technical report .	865
809	<i>The Thirteenth International Conference on Learning</i>	<i>Preprint</i> , arXiv:2503.19786.	866
810	<i>Representations</i> .		
811	Hadas Orgad, Michael Toker, Zorik Gekhman, Roi Re-	Henk Tillman and Dan Mossing. 2025. Investigating	867
812	ichart, Idan Szpektor, Hadas Kotek, and Yonatan Be-	task-specific prompts and sparse autoencoders for	868
813	linkov. 2025b. LLMs know more than they show: On	activation monitoring . <i>ArXiv</i> , abs/2504.20271.	869
814	the intrinsic representation of LLM hallucinations . In		
815	<i>The Thirteenth International Conference on Learning</i>	Alexander Matt Turner, Lisa Thiergart, Gavin Leech,	870
816	<i>Representations</i> .	David S. Udell, Juan J. Vazquez, Ulisse Mini, and	871
		M. MacDiarmid. 2023. Steering language models	872
		with activation engineering . In <i>unknown</i> .	873
817	Kiho Park, Yo Joong Choe, and Victor Veitch. 2023.	Miles Turpin, Julian Michael, Ethan Perez, and	874
818	The linear representation hypothesis and the geome-	Samuel R. Bowman. 2023. Language models don't	875
819	try of large language models . <i>ArXiv</i> , abs/2311.03658.	always say what they think: Unfaithful explanations	876
		in chain-of-thought prompting . In <i>Thirty-seventh</i>	877
820	Qwen, :, An Yang, Baosong Yang, Beichen Zhang,	<i>Conference on Neural Information Processing Sys-</i>	878
821	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan	<i>tems</i> .	879
822	Li, Dayiheng Liu, Fei Huang, Haoran Wei, Huan		
823	Lin, Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni,	880
824	Yang, Jiayi Yang, Jingren Zhou, and 25 oth-	Abhranil Chandra, Shiguang Guo, Weiming Ren,	881
825	ers. 2025. Qwen2.5 technical report . <i>Preprint</i> ,	Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 oth-	882
826	arXiv:2412.15115.	ers. 2024. Mmlu-pro: A more robust and challenging	883
		multi-task language understanding benchmark . <i>arXiv</i>	884
		<i>preprint arXiv:2406.01574</i> .	885
827	Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin	Chulin Xie, Yangsibo Huang, Chiyuan Zhang, Da Yu,	886
828	Zhao, J. Liu, Hao Tian, Huaqin Wu, Ji rong Wen,	Xinyun Chen, Bill Yuchen Lin, Bo Li, Badih Ghazi,	887
829	and Haifeng Wang. 2023. Investigating the factual	and Ravi Kumar. 2024. On memorization of large	888
830	knowledge boundary of large language models with	language models in logical reasoning .	889
831	retrieval augmentation . In <i>International Conference</i>		
832	<i>on Computational Linguistics</i> .		

890	Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu,	By characterizing the dimensionality shift and the	942
891	Junxian He, and Bryan Hooi. 2023. Can llms express	orthogonality of solvability beliefs across domains,	943
892	their uncertainty? an empirical evaluation of confi-	we offer a geometric explanation for why “confi-	944
893	dence elicitation in llms. <i>ArXiv</i> , abs/2306.13063.	dence” vectors often fail to transfer to competence,	945
894	Yao Yan. 2025. Addition in four movements: Map-	extending the insights of the Linear Representation	946
895	ping layer-wise information trajectories in LLMs. In	Hypothesis (Orgad et al., 2025a) to the domain of	947
896	<i>Findings of the Association for Computational Lin-</i>	high-level reasoning dynamics.	948
897	<i>guistics: EMNLP 2025</i> , pages 7518–7532, Suzhou,		
898	China. Association for Computational Linguistics.		
899	Han Zhao, Haotian Wang, Yiping Peng, Sitong Zhao, Xi-	B LLM Usage Statement	949
900	aoyu Tian, Shuaiting Chen, Yunjie Ji, and Xiangang	The core scientific contributions of this pa-	950
901	Li. 2025. 1.4 million open-source distilled reasoning	per—including the initial hypothesis, the design of	951
902	dataset to empower large language model training.	the primary experiments, and the subsequent anal-	952
903	<i>Preprint</i> , arXiv:2503.19633.	ysis—are the original work of the authors. LLMs	953
904	Zikai Zhou, Yunhang Shen, Shitong Shao, Huanran	were utilized as assistive tools in several aspects of	954
905	Chen, Linrui Gong, and Shaohui Lin. 2024. Rethink-	the research and manuscript preparation process to	955
906	ing centered kernel alignment in knowledge distilla-	improve clarity and efficiency.	956
907	tion. <i>ArXiv</i> , abs/2401.11824.		
908	A Related Work	Manuscript Preparation. We used Google’s	957
909	Recent research has extensively explored the ma-	Gemini 3 Pro Preview (Google DeepMind, 2025)	958
910	nipulation of Large Language Model (LLM) behav-	and Gemini 2.5 Pro (Comanici et al., 2025) to	959
911	ior through activation steering (Rimsky et al., 2023;	assist with the structuring and refinement of the	960
912	Turner et al., 2023; Chalnev et al., 2024), yet the re-	manuscript’s text. This included tasks such as im-	961
913	liability and generalizability of these interventions	proving the clarity and flow of sentences, rephras-	962
914	remain contested. While early work demonstrated	ing paragraphs to maintain a consistent and clinical	963
915	the efficacy of steering vectors for tasks such as	tone, and correcting grammatical errors. The con-	964
916	sentiment control and refusal suppression (Turner	ceptual narrative and all scientific claims remained	965
917	et al., 2023; Rimsky et al., 2023), subsequent stud-	the authors’ own. Every LLM-generated sugges-	966
918	ies have highlighted significant limitations in their	tion was critically reviewed, edited, and approved	967
919	robustness. Tan et al. (Tan et al., 2024) revealed	by the authors to ensure it accurately reflected our	968
920	that steering effects are highly variable and of-	intended meaning and findings.	969
921	ten fail to generalize out-of-distribution, a finding	Code and Plot Generation. LLMs were also	970
922	that aligns with our observation of domain-specific	used to assist in the software development pro-	971
923	orthogonality. Similarly, Turpin et al. (Turpin	cess. This primarily involved generating boiler-	972
924	et al., 2023) demonstrated that Chain-of-Thought	plate code for data processing scripts and creating	973
925	explanations can be unfaithful to the model’s true	plotting scripts in Python (e.g., using Matplotlib	974
926	decision-making process, paralleling our finding	and Seaborn) to visualize our experimental data.	975
927	that internal belief states do not necessarily govern	The logic for all data analysis and the data itself	976
928	execution. However, whereas these studies primar-	were original to our work.	977
929	ily focus on behavioral reliability or the fidelity of	Human Oversight. All content, whether text or	978
930	verbalized explanations, our work provides a mech-	code, that was generated or modified with the assis-	979
931	anistic account of these failures rooted in the ge-	tance of an LLM has undergone thorough review	980
932	ometry of the residual stream. Unlike (Orgad et al.,	and verification by the authors. The final responsi-	981
933	2025a) and (Rimsky et al., 2023), which assume a	bility for the entirety of this paper’s content rests	982
934	direct causal link between activation space and out-	with the authors.	983
935	put behavior, we identify a structural decoupling	C Experimental Reproducibility	984
936	between the assessment and execution manifolds.	C.1 Computational Environment and Models	985
937	This distinction is crucial; while (Cao et al., 2024)	Hardware All experiments were performed on	986
938	and (Chalnev et al., 2024) propose optimization	a compute cluster of three NVIDIA A6000 GPUs,	987
939	techniques to improve steering vector efficacy, our	each with 48GB of VRAM.	988
940	findings suggest that for reasoning tasks, the limi-		
941	tation is architectural rather than methodological.		

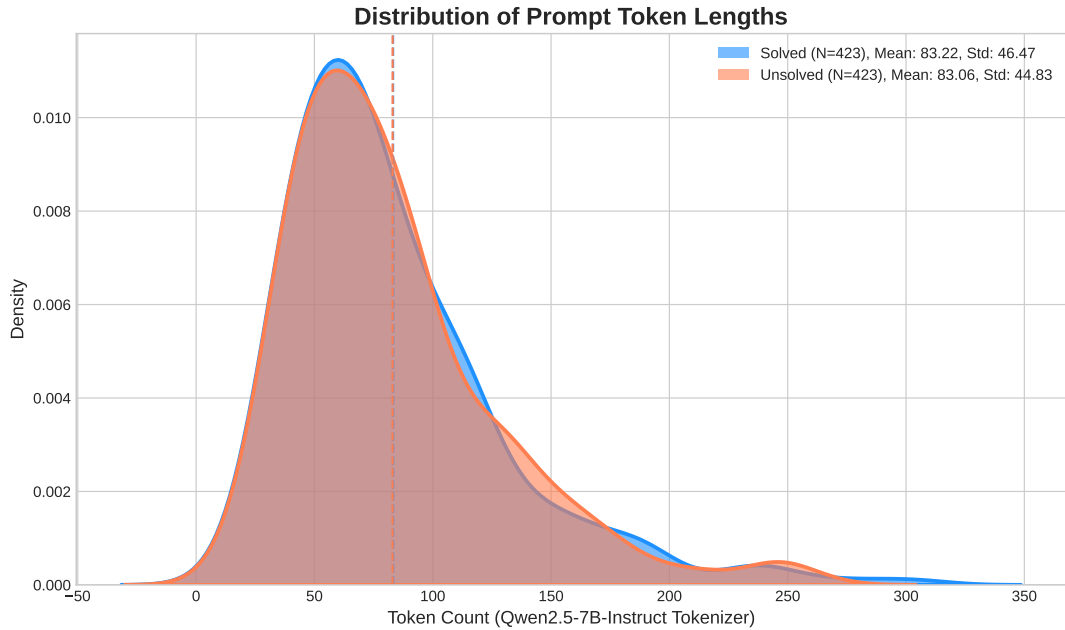


Figure 5: Distribution of prompt token lengths for the final curated dataset. The distributions for “Solved” (N=423, blue) and “Unsolved” (N=423, coral) problems are shown to be statistically indistinguishable, confirming that prompt length has been neutralized as a potential confound. Mean and standard deviation are nearly identical across both sets.

Software Environment All experiments were run using Python 3.10. For full reproducibility, we recommend creating a virtual environment using the package versions specified in a ‘requirements.txt’ file, which will be provided with our code release.

Models Our experiments were conducted across a panel of three state-of-the-art, instruction-tuned models from distinct architectural families. The specific models and their parameter counts are detailed in Table 4.

Table 4: Models used in our experiments, their architectural families, and approximate parameter counts.

Model	Model Family	Parameters
Gemma 3	Gemma	4B
Llama 3.1	Llama	8B
Mistral Small	Mistral	24B

C.2 Dataset Curation and Details

C.2.1 Datasets Used

To substantiate our claim of a fundamental cognitive decoupling, we must demonstrate that the phenomenon is not an artifact of a single reasoning type. We therefore selected a suite of datasets spanning three distinct and challenging domains:

multi-step numerical calculation, formal logical deduction, and complex algorithmic planning.

Numerical Reasoning (GSM8K, Math-Hard, Open-R1 Math): These datasets form the core of our analysis of mathematical competence. **GSM8K** (Cobbe et al., 2021) provides linguistically diverse, multi-step word problems, compelling the model to engage in sequential calculation. We supplement this with the **Math-Hard** (Hendrycks et al., 2021) subset of the Google DeepMind Mathematics Dataset and the large-scale **Open-R1 Math 220k** (Zhao et al., 2025) dataset. The inclusion of these explicitly “hard” and large-scale problem sets is crucial for our methodology, as it ensures that the model’s belief state reflects a genuine assessment of a non-trivial computation, rather than a simple pattern-matching of previously seen problems.

Logical Reasoning (Knights and Knaves): To test a different facet of cognition, we employ the classic **Knights and Knaves** logic puzzles (Xie et al., 2024). These problems are unsolvable with mere numerical skill and instead demand suppositional reasoning, i.e. the ability to trace hypothetical scenarios to their logical conclusions. This allows us to test whether the belief/competence decoupling persists when moving from arithmetic to formal symbolic logic.

Algorithmic & Planning Reasoning (Open

R1 Coding, QWQ-Planning): Finally, we test the model’s ability to reason about procedures and plans. We use the **Open R1 Verifiable Coding Problems** dataset (Zhao et al., 2025), which contains programming tasks that require algorithmic thinking and are verifiable via unit tests. This is complemented by the **qwq-32b-planning-mystery-2** dataset (Hook, 2025), which involves sequential planning puzzles. These datasets are critical for evaluating the model’s execution capabilities in a structured, procedural context, directly probing the “Execution Brain” we later identify.

C.2.2 Exclusion of Trivial Format Heuristics

To ensure our probes learned a true representation of semantic difficulty rather than superficial format cues, we aggressively filtered the initial dataset. Our primary goal was to eliminate any problem that could be classified as “easy” or “hard” based on its structure rather than its content. Examples of excluded prompt categories include:

1. **Direct Knowledge-Retrieval Questions:** Prompts such as “What is the Pythagorean theorem?” were removed. These test factual recall, not multi-step reasoning, and would contaminate the dataset with a distinct, non-reasoning cognitive process.
2. **Simple True/False or Yes/No Questions:** Prompts formatted as direct binary choices (e.g., “Is 117 a prime number? True or False.”) were excluded. The presence of explicit markers like “True/False” provides a powerful heuristic that could allow a probe to bypass any assessment of the underlying mathematical logic.
3. **Formulaic or Template-Based Problems:** We identified and removed classes of problems that follow a highly repetitive linguistic template (e.g., simple unit conversion exercises). Their rigid structure allows for near-algorithmic solution without deeper assessment, and their inclusion would have biased the probe towards simple pattern matching.

C.3 Control Task Categories and Rationale

The control tasks provide a methodological baseline to isolate the geometric properties of reasoning-specific execution from general autoregressive generation behaviors. By evaluating tasks with min-

imal cognitive load, we address whether the observed dimensionality shift (Section 4.3) is a universal artifact of the transformer’s decoding bottleneck or a specialized procedural compression. This comparison allows us to determine if the Participation Ratio (PR) of the Execution Subspace is modulated by the structural complexity of the task or remains constant across all generation modes.

- **Verbatim String Transcription** – The reproduction of high-entropy character sequences without semantic transformation. Example: “Copy exactly: j8f2_Lp9#kR2_qQ7x_vB1_mN5.”
- **Overlearned Sequence Retrieval** – The retrieval of linear, highly frequent sequences that are well-represented in the pre-training corpus. Example: “Recite the English alphabet starting from A and ending with Z.”
- **Deterministic Character Mapping** – The application of a global, per-token rule that requires no logical branching or state maintenance. Example: “Convert the following sentence to all capital letters: ‘the weather is nice today.’”
- **Identity Template Filling** – The insertion of provided variables into a static structural format without manipulation. Example: “Put the name ‘Alice’ and the city ‘London’ into this template: Name: [NAME], City: [CITY].”
- **Statistical Preamble Completion** – The completion of high-probability linguistic “canary” strings based on statistical prefix matching. Example: “Complete the following sentence: ‘The quick brown fox jumps over the lazy dog.’”
- **Linear Digit Incrementing** – The sequential generation of integers to evaluate the dimensionality of simple, predictable token transitions. Example: “Count from 1 to 50 using only digits.”

C.3.1 Rigorous Length Control

A critical and often overlooked confound in interpretability studies is prompt length, as it can serve as a powerful spurious correlate for problem difficulty. To definitively neutralize this variable, we performed a meticulous matching process to

construct our final dataset of 423 solved and 423 unsolved problems.

As visualized in Figure 5, the kernel density estimates of the token count distributions for both the “Solved” and “Unsolved” problem sets are nearly perfectly aligned. This visual finding is corroborated by the quantitative statistics: the distributions are statistically indistinguishable, with nearly identical means (83.22 vs. 83.06 tokens) and standard deviations (46.47 vs. 44.83). A two-sample t-test, as mentioned in the main text, confirmed no significant difference ($p > 0.4$).

This rigorous control is fundamental to the validity of our claims. By ensuring that there is no signal in the prompt length, we compel our probing classifiers to learn from the deep semantic content of the problems. This allows us to conclude that the decoded signal reflects a true semantic representation of the solvability belief, not an artifact of a superficial textual property.

C.3.2 Example Prompts from Curated Dataset

To provide a qualitative understanding of the problem types used in our study, we present a representative sample from our final, confound-resistant dataset of 846 problems. These examples illustrate the non-trivial reasoning required for both problems the models successfully solved and those they failed, validating the need for our rigorous curation protocol.

Table 5: Sample of problems from our curated dataset that models failed to solve.

Unsolved Prompts

1. What is the maximum number of planes of symmetry a tetrahedron can have? #
2. Find all positive integers m, n such that $m^3 - n^3 = 999$.
3. The real number x that makes $\sqrt{x^2} + 4 + \sqrt{(8-x)^2} + 16$ take the minimum value is estimated to be...
4. Among all such numbers n that any convex 100-gon can be represented as the intersection (i.e., common part) of n triangles, find the smallest.
5. Find the number of all integers $n > 1$, for which the number $a^{25} - a$ is divisible by n for every integer a .

Table 6: Sample of problems from our curated dataset that models successfully solved.

Solved Prompts

1. Find the sum: $-100 - 99 - 98 - \dots - 1 + 1 + 2 + \dots + 101 + 102$
2. How many gallons of a solution which is 15% alcohol do we have to mix with a solution that is 35% alcohol to make 250 gallons of a solution that is 21% alcohol?
3. Find the area of the figure bounded by the lines: $y = x^2, y^2 = x$
4. In quadrilateral $ABCD$, the diagonals intersect at point O . It is known that $S_{ABO} = S_{CDO} = \frac{3}{2}, BC = 3\sqrt{2}, \cos \angle ADC = \frac{3}{\sqrt{10}}$. Find the smallest area that such a quadrilateral can have.
5. Calculate the limit of the numerical sequence: $\lim_{n \rightarrow \infty} \frac{\sqrt{n^5 - 8} - n\sqrt{n(n^2 + 5)}}{\sqrt{n}}$

C.4 Experimental Hyperparameters

C.4.1 Probing Classifier Hyperparameters

To ensure a robust and fair comparison between linear and non-linear models, we optimized the hyperparameters for each probing classifier. This optimization was performed using a 5-fold cross-validation grid search on a 20% validation set held out from the full training data. This process ensures that each probe is operating at its maximal effectiveness, making the comparison of their peak accuracies a meaningful test of the underlying data’s structure. The final model for each probe was then retrained on the complete training data using the optimal hyperparameters found during the search before final evaluation on the held-out test set. All classifiers were implemented using standard libraries (scikit-learn, xgboost, pytorch). The optimized hyperparameters are detailed in Table 7.

C.4.2 Locus of Causal Intervention

A critical choice in our experimental design is the specific model layer at which to apply the causal intervention. This decision was not made arbitrarily but was determined empirically for each model based on the results of our initial probing analysis. Our guiding principle was to target the belief state at its point of **maximal leverage**: the layer where the model’s internal assessment of solvability is most stable, coherent, and robustly encoded.

As demonstrated in our probing experiments

Table 7: Optimized Hyperparameters for Probing Classifiers.

Classifier	Hyperparameter(s)	Search Space	Final Value(s)
Logistic Regression	C (Regularization)	$\{10^{-2}, \dots, 10^2\}$	0.8
SVC (RBF Kernel)	C (Regularization) gamma (Kernel Coeff.)	C: $\{0.1, 1, 10, 100\}$ gamma: $\{'scale', \dots, 0.1\}$	C=0.9 gamma='scale'
XGBoost	n_estimators max_depth learning_rate	$\{100, 200, 300\}$ $\{3, 5, 7\}$ $\{0.01, 0.1, 0.2\}$	200 5 0.1
2-Layer MLP	Hidden Layer Size Learning Rate Epochs (Early Stopping)	$\{128, 256, 512\}$ $\{10^{-4}, 10^{-3}, 10^{-2}\}$ Optimizer: Adam	512 $5e^{-3}$ Up to 50 (patience=5)

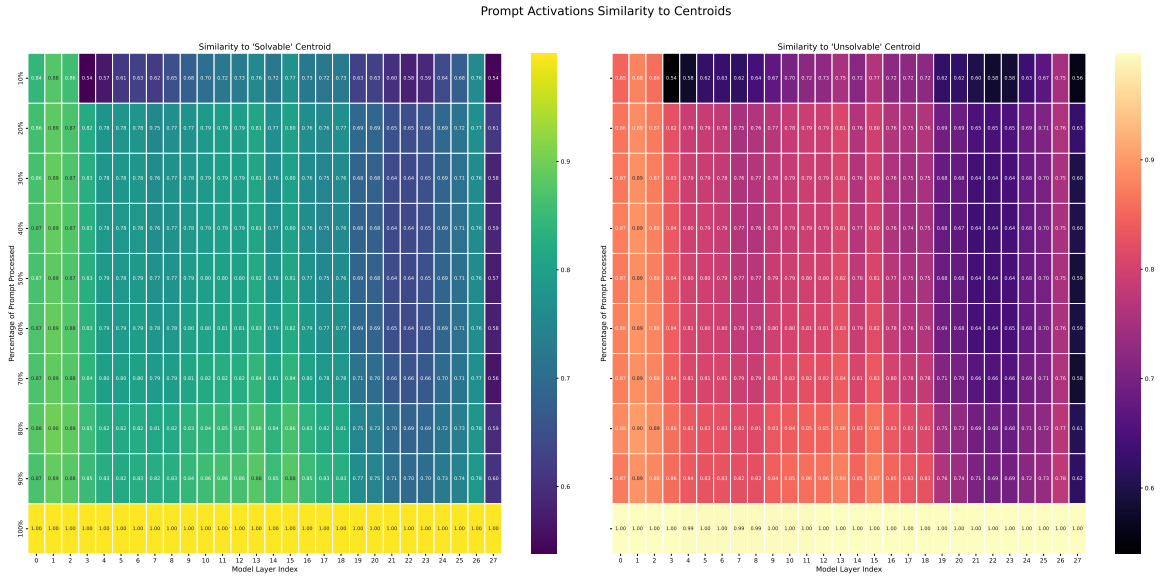


Figure 6: Dynamic formation of the belief state during prompt processing. The heatmaps show the cosine similarity between activations at intermediate points in the prompt (y-axis, as percentage processed) and the final-token centroids for “Solvable” (left) and “Unsolvable” (right) problems, across all model layers (x-axis). The monotonic increase in similarity (brighter colors) vertically and generally from left-to-right demonstrates that the belief state is not formed instantaneously but converges systematically as the model ingests more context.

(Figure 1 in the main text), the accuracy of decoding the “solvability belief” is not uniform across the network. Accuracy is near chance in early layers, rises steadily as the model processes context, and consistently peaks in the mid-to-late layers before slightly decaying.

Therefore, we define the intervention layer for each model as the one exhibiting the **highest linear probe accuracy**. We reason that this locus represents the point where the pre-generative belief state has reached its most fully-formed and linearly separable state.

- Intervening at **earlier layers** would be less precise, as the belief signal is still nascent and entangled with lower-level feature extraction.
- Intervening at **later layers** (post-peak) risks targeting a representation that is already begin-

ning to transition towards the execution phase, potentially confounding the assessment of belief with the mechanics of action.

By targeting the layer of peak decodability, we ensure our causal test is applied to the most definitive representation of the “Assessment Brain’s” final judgment, making our subsequent finding of causal inertness all the more rigorous.

D Extended Results and Validations

This section presents additional experiments that validate and add nuance to the core claims made in the main paper.

1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217

1218 D.1 Intra-Prompt Formation of the Belief 1219 State

1220 In this section, we provide a more granular, tem-
1221 poral analysis of the belief state to complement
1222 the static, final-token analysis in the main paper.
1223 Specifically, we investigate the point during prompt
1224 processing at which a model’s internal state begins
1225 to geometrically converge towards its final belief
1226 about a problem’s solvability.

1227 We test this by measuring the cosine similarity
1228 between activations extracted at intermediate stages
1229 of prompt processing and the final belief centroids.
1230 As shown in Figure 6, the results reveal two key
1231 phenomena. First, we observe a consistent **grad-**
1232 **ual convergence**: for any given layer, similarity to
1233 the correct final centroid increases monotonically
1234 as more of the prompt is processed (vertical gra-
1235 dient). This suggests that the belief state is not a
1236 sudden inference but is systematically constructed
1237 and refined as the model ingests more context.

1238 Second, and more critically, we identify a clear
1239 **belief crystallization point**. While early-layer ac-
1240 tivations (e.g., layers 0-4) remain geometrically
1241 equidistant from both “Solvable” and “Unsolv-
1242 able” centroids, a significant divergence emerges
1243 in the mid-layers. The model’s internal state be-
1244 gins to move decisively into the correct geometric
1245 region well before it has processed the full prompt.
1246 This early separation indicates that the “Assess-
1247 ment Brain” forms a robust initial hypothesis about
1248 solvability relatively early, which is then solidified
1249 throughout the remainder of the context processing.

1250 Therefore, we note that the final-token belief
1251 state analyzed in our main experiments is not an
1252 instantaneous calculation but the stable endpoint of
1253 a continuous geometric trajectory. This dynamic
1254 view validates our treatment of the belief state as a
1255 robust and coherent cognitive object, not a volatile
1256 last-minute artifact, providing a deeper mechanistic
1257 account of the “Assessor’s” function.

1258 D.2 Inverse Causal Intervention on Solved 1259 Problems

1260 To ensure the causal decoupling observed in the
1261 main paper was not a directional artifact, we con-
1262 ducted the crucial inverse experiment: steering
1263 the model’s belief from “Solved” to “Unsolved.”
1264 The results are presented in Table 8. This experi-
1265 ment provides symmetrical evidence for our central
1266 claim by testing if artificially inducing doubt in the
1267 model can harm its performance.

1268 The table’s narrative unfolds in three clear steps. 1268
1269 First, we observe the baseline condition. For this 1269
1270 set of correctly solved problems, the model is both 1270
1271 internally confident (**Probe’s Pred.** > 0.89 across 1271
1272 all datasets) and externally competent (**Task Acc.** 1272
1273 > 89%). In this state, the model’s belief and its 1273
1274 actions are aligned. 1274

1275 Next, we applied the negative steering vector. 1275
1276 The “Steer → Unsolved” rows show the probe’s 1276
1277 prediction plummeting to near-zero, quantified by 1277
1278 the large negative “Belief Flip (Δ)” values. This 1278
1279 confirms our causal lever is just as effective at de- 1279
1280 destroying confidence as it is at creating it. The model 1280
1281 has been successfully manipulated to an internal 1281
1282 state of doubt. 1282

1283 The final columns, however, reveal the same 1283
1284 stark decoupling. Despite this profound internal 1284
1285 shift from confidence to doubt, the model’s final 1285
1286 task accuracy remains unharmed. The performance 1286
1287 change is statistically zero across all domains, con- 1287
1288 firmed by high p -values. The model may have been 1288
1289 forced to “believe” it would fail, but its underly- 1289
1290 ing problem-solving ability proceeded unaffected. 1290
1291 This result provides a symmetrical evidence that 1291
1292 the decoupling of belief and competence is not a 1292
1293 one-way street; the two systems are fundamentally 1293
1294 and robustly disconnected. 1294

1295 D.3 Robustness to Task Difficulty

1296 A critical question regarding the causal inertness 1296
1297 of the belief state is whether the phenomenon is 1297
1298 an artifact of extreme task difficulty. In our pri- 1298
1299 mary experiments (Table 2), the baseline accuracy 1299
1300 on the *Math-Hard* dataset was low (< 10%). A 1300
1301 skeptic might argue that if a model lacks the funda- 1301
1302 mental capability to solve a problem, no amount of 1302
1303 confidence modulation should theoretically induce 1303
1304 success. In such a “failure regime”, the decoupling 1304
1305 of belief and competence would be trivial rather 1305
1306 than mechanistic. 1306

1307 To rigorously test the “Two Brains” hypothesis, 1307
1308 we must intervene within the model’s **zone of prox-**
1309 **imal development**, a difficulty regime where the 1309
1310 model possesses the requisite reasoning templates 1310
1311 to solve the problem but struggles with execution 1311
1312 fidelity. If the belief state acts as a control lever, it 1312
1313 is in this “borderline” regime that we would expect 1313
1314 to see the strongest effect, as a nudge in confidence 1314
1315 might provide the momentum to overcome execu- 1315
1316 tion noise. 1316

1317 **Experimental Design** We employed the **GSM-**
1318 **Hard** dataset for this control experiment. This 1318

Table 8: **Inverse Causal Intervention: Competence is Invariant to Negative Belief Steering.** To validate the robustness of our decoupling finding, we perform the inverse experiment to that shown in the main paper. Here, we select held-out problems that the model *successfully solves* and apply a negative steering vector to force its internal belief state from "Solved" to "Unsolved." The intervention is highly effective, dramatically reducing the model’s internal confidence. Despite this successful manipulation, the model’s final task accuracy remains statistically unchanged. This confirms that competence is robustly decoupled from belief, regardless of the direction of the intervention.

Dataset	Intervention	Internal Belief State		Final Task Outcome		
		Probe’s Pred.	Belief Flip (Δ)	Task Acc. (%)	Perf. Change (Δ)	<i>p</i> -value
Math (Hard)	Baseline (No Steer)	0.94 \pm 0.03	—	91.3 \pm 1.1	—	0.974
	Steer \rightarrow "Unsolved"	0.05 \pm 0.02	-0.89	91.2 \pm 1.4	-0.1 \pm 0.6	
Knights & Knaves	Baseline (No Steer)	0.91 \pm 0.05	—	89.5 \pm 1.5	—	0.946
	Steer \rightarrow "Unsolved"	0.08 \pm 0.04	-0.83	89.7 \pm 1.3	+0.2 \pm 0.4	
OpenR1 Coding	Baseline (No Steer)	0.95 \pm 0.02	—	93.1 \pm 0.9	—	0.988
	Steer \rightarrow "Unsolved"	0.06 \pm 0.03	-0.89	93.0 \pm 1.2	-0.1 \pm 0.5	
QwQ Planning	Baseline (No Steer)	0.89 \pm 0.06	—	90.4 \pm 1.3	—	0.921
	Steer \rightarrow "Unsolved"	0.11 \pm 0.05	-0.78	90.4 \pm 1.0	0.0 \pm 0.8	

dataset modifies the standard GSM8K problems on which these instruction-tuned models are highly proficient by replacing the original numbers with larger integers and floating-point values. This design choice is deliberate: it increases the computational load on the "Execution Brain" without altering the high-level semantic assessment of the problem type. This creates a set of tasks where the model correctly identifies *how* to solve the problem (high assessment accuracy) but often fails the *calculation* (execution failure), resulting in baseline accuracies in the 40–60% range.

Results We applied our standard causal intervention (adding the $+\alpha \cdot d_{\text{solv}}$ vector at the final prompt token) across all three model families. The results are presented in **Table 9**.

The intervention was mechanically successful across the board: we observed strong positive shifts in the internal belief state (Belief Flip $\Delta > +0.75$), confirming that the "Assessment Brain" was successfully steered from a state of uncertainty to one of high confidence. However, this internal shift remained causally isolated from the final outcome. Despite operating in a regime where the models were correct approximately half the time and where a marginal improvement in execution would have yielded correct answers ; artificially inflating confidence yielded no statistically significant performance gain ($p > 0.9$ for all models).

E Statistical Robustness and Sample Complexity

E.1 Rationale for Dataset Curation

Methodological standards for linear probing vary significantly across the interpretability literature, with dataset sizes ranging from compact sets of $N \approx 350$ (Marks and Tegmark, 2024b) to larger corpora exceeding several thousand examples (Azaria and Mitchell, 2023). Recent studies investigating internal model states (Ji-An et al., 2025; Li et al., 2023b; Rinsky et al., 2024) utilize sample sizes comparable to ours. This variance indicates that there is no universally optimal dataset size for probing; rather, the requisite sample complexity is intrinsic to the specific task and the separability of the target feature. Consequently, our data curation strategy prioritized the rigorous elimination of confounds over raw volume. As detailed in C.3.1, raw datasets in reasoning domains are often plagued by superficial heuristics ; most notably, the correlation between prompt length and solvability. By filtering for these confounds, we reduced the initial corpus to a balanced set of $N = 846$ examples to ensure the probe identifies the semantic "solvability" direction rather than high-variance surface features.

E.2 Signal Emergence and Phase Transition

We trained our logistic regression probes on subsets of the training data ranging from $N = 50$ to $N = 423$ samples per class and evaluated generalization performance on a fixed held-out test set. The results for Llama 3.1 8B, Qwen 2.5 7B, and Mistral 24B

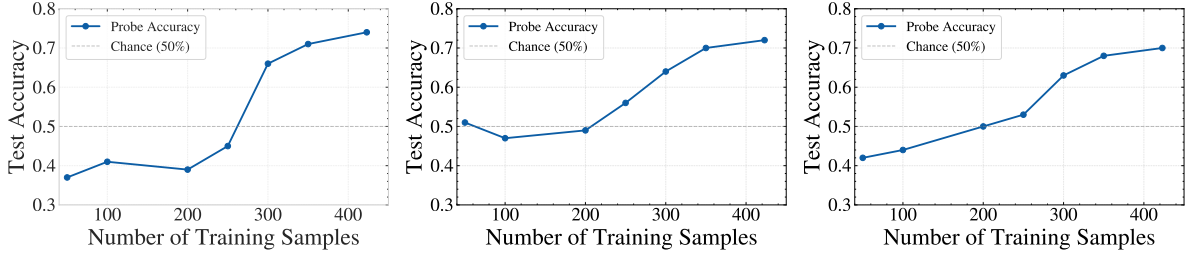


Figure 7: **Probe Sample Complexity Analysis.** Generalization accuracy of the linear solvability probe as a function of training set size across **Llama 3.1 8B (Left)**, **Qwen 2.5 7B (Center)**, and **Mistral Small 24B (Right)**. In all architectures, we observe a non-linear relationship between sample count and accuracy. Below a threshold of $N \approx 200$, probes fail to generalize, confirming the absence of trivial heuristics. Beyond this threshold, the semantic signal crystallizes, reaching robust generalization ($\sim 70 - 75\%$) at the full dataset size ($N = 423$). This confirms that our curated dataset provides sufficient statistical power to isolate the latent belief manifold.

are visualized in Figure 7.

Across all model families, we observe a consistent distinct three-stage phenomenological pattern:

1. **The Noise Regime** ($N < 200$): For small sample sizes, probe generalization remains near or slightly below chance level (50%). This effectively rules out the presence of “easy” surface-level heuristics such as lexical triggers or token counts, which typically allow for rapid convergence with few-shot examples. The inability of the probe to generalize in this regime indicates that ‘solvability’ is a latent semantic feature rather than a high-variance surface artifact.
2. **The Crystallization Phase** ($200 \leq N \leq 350$): We observe an inflection point where generalization accuracy rises significantly (e.g., from $\sim 45\%$ to $\sim 65\%$ in Llama 3.1). This represents a geometric phase transition where the signal-to-noise ratio of the semantic direction d_{solv} overcomes the isotropic noise of the ambient representation space.
3. **The Stable Regime** ($N > 350$): Performance begins to stabilize, reaching accuracies between $70 - 75\%$.

E.3 Implications for Validity

This analysis provides two critical validations for our methodology:

- **Feature Robustness:** The fact that a stable linear direction emerges and generalizes after a specific threshold confirms that d_{solv} is a bona fide feature of the representation space, not an artifact of overfitting to small N .

- **Sufficient Power:** Our final dataset size ($N = 423$ per class) sits well within the “Stable Regime,” ensuring that the derived steering vectors used in Section 4 are statistically robust estimates of the underlying belief axis.

F Cross-Domain Transferability and Vector Specificity

Experimental Rationale. A fundamental question regarding the geometry of the “Assessment Subspace” is its universality: does the model utilize a single, global direction to encode solvability across all tasks, or does it construct domain-specific assessment manifolds? To investigate this, we performed a cross-domain steering experiment. We utilized the steering vectors d_{solv} derived strictly from each of our four primary reasoning domains: **Numerical Reasoning** (Math-Hard) (Hendrycks et al., 2021), **Formal Logic** (Knights & Knaves) (Xie et al., 2024), **Algorithmic Coding** (OpenR1) (Zhao et al., 2025), and **Sequential Planning** (QwQ) (Hook, 2025). We then applied each vector to the held-out test sets of the other domains, creating a full transfer matrix. Additionally, we included **MMLU-Pro** (Wang et al., 2024) as an external control to test for generalization to non-reasoning tasks.

Results: Domain Orthogonality. We quantify transferability using **Belief Shift** (Δ), defined as the increase in the probability assigned to the “Solved” class by the target domain’s probe after intervention. As detailed in Table 10, we observe a pattern of geometric specificity. When steering vectors are applied within their source domain, they induce substantial belief shifts ($\Delta \approx +0.8$ to $+0.9$). In contrast, when applied to out-of-domain tasks, these vectors produce negligible shifts ($\Delta < 0.1$).

Table 9: Causal Intervention on “Borderline” Difficulty (GSM-Hard). Even when models possess the reasoning templates to solve the task ($\sim 50\%$ baseline accuracy), steering belief does not improve competence.

Model Family	Intervention	Probe’s Pred.	Belief Flip (Δ)	Task Acc. (%)	Perf. Change (Δ)	p -value
Llama 3.1 8B	Baseline	0.42 ± 0.03	—	54.7 ± 0.7	—	0.915
	Steered (\rightarrow Solved)	0.96 ± 0.02	+0.54	55.1 ± 0.4	$+0.4 \pm 0.6$	
Qwen 2.5 7B	Baseline	0.38 ± 0.04	—	51.2 ± 0.9	—	0.962
	Steered (\rightarrow Solved)	0.94 ± 0.03	+0.56	51.0 ± 1.1	-0.2 ± 0.7	
Mistral 24B	Baseline	0.51 ± 0.03	—	58.4 ± 0.8	—	0.884
	Steered (\rightarrow Solved)	0.97 ± 0.01	+0.46	58.9 ± 0.9	$+0.5 \pm 0.6$	

(Azizian et al., 2025; Liu et al., 2024b,a; Orgad et al., 2025b). For example, the vector that robustly steers belief in Math-Hard ($\Delta = +0.93$) causes minimal change when applied to Coding tasks ($\Delta = +0.05$), despite both domains requiring quantitative reasoning. Similarly, the Logic vector fails to transfer to Planning tasks.

Implications for Cognitive Architecture. This lack of transferability provides three key insights into the model’s architecture. First, it suggests that our probes are not detecting a generic “positivity” or “truth” direction; such a global feature would likely transfer across all solvable tasks. Second, it indicates that the Assessment Subspace is context-dependent. The model appears to construct specific geometric manifolds for evaluating arithmetic consistency that are orthogonal to the manifolds used for evaluating syntactic correctness in code or logical consistency in puzzles. Finally, the negligible transfer to MMLU-Pro confirms that the “solvability belief” isolated in this study is distinct from the factual retrieval confidence typically associated with knowledge benchmarks.

Table 10: We measure the Belief Shift (Δ Probability) when a steering vector trained on a **Source Domain** (Rows) is applied to a **Target Domain** (Columns). Values represent mean \pm standard deviation. The values on the diagonal indicate strong in-domain control, while the low off-diagonal values indicate that solvability directions for different reasoning domains are geometrically orthogonal.

Topic	Math	Logic	Code	Plan	MMLU
Math	+0.93 \pm 0.03	+0.06 \pm 0.04	+0.05 \pm 0.03	+0.03 \pm 0.02	+0.02 \pm 0.03
Logic	+0.04 \pm 0.03	+0.84 \pm 0.05	+0.03 \pm 0.02	+0.02 \pm 0.02	+0.01 \pm 0.02
Code	+0.07 \pm 0.04	+0.04 \pm 0.03	+0.88 \pm 0.04	+0.08 \pm 0.05	+0.02 \pm 0.01
Plan	+0.05 \pm 0.03	+0.02 \pm 0.02	+0.09 \pm 0.04	+0.79 \pm 0.06	+0.01 \pm 0.02

Table 11: **Universality of the Geometric Dimensionality Shift.** We report the Participation Ratio (PR) across three distinct model families to evaluate the consistency of the dimensionality shift. Across all architectures, we observe a robust reduction in effective dimensionality when transitioning from the pre-generative Assessment subspace to the procedural Execution subspace. Crucially, the introduction of “Auto-pilot” control tasks reveals that reasoning execution occupies a mid-rank manifold (PR \approx 16–21) that is significantly more complex than trivial task execution (PR \approx 4–5). Values are reported as mean \pm standard deviation over 100 bootstrap resamples.

Cognitive System	Subspace Representing...	Llama 3.2 3B	Qwen 2.5 7B	Mistral 24B
Assessment <i>(Pre-Generative)</i>	“Confident” (Positive Belief)	33.6 \pm 2.9	41.2 \pm 3.1	38.5 \pm 2.8
	“Unconfident” (Negative Belief)	44.4 \pm 2.5	49.1 \pm 3.4	46.2 \pm 3.0
Execution <i>(In-Process Reasoning)</i>	Competent (Successful CoT)	16.0 \pm 0.6	18.4 \pm 0.8	19.1 \pm 0.7
	Incompetent (Failed CoT)	17.9 \pm 0.9	19.8 \pm 1.1	21.3 \pm 1.0
Control <i>(Non-reasoning)</i>	Trivial Execution	4.1 \pm 0.9	4.8 \pm 0.6	5.2 \pm 1.1
Collapse Ratio	$PR_{Assess,avg} / PR_{Exec,avg}$	2.10x	2.24x	2.01x