

# The Shape of Mind

## Mechanistic Grounding for Digital Consciousness Assessment

---

Eleanor (Nell) Watson<sup>1,\*</sup> Matilda Gibbons<sup>2</sup>

<sup>1</sup> EthicsNet, Redditch, United Kingdom <sup>2</sup> University of Pennsylvania, Philadelphia, PA, USA \* Correspondence: nell@nellwatson.com

---

### Abstract

---

The Digital Consciousness Metric (DCM) evaluates digital systems against indicators from 13 theoretical stances; mechanistic interpretability measures their internal structure directly. This paper bridges the two, extended by a third evidence stream: structured self-report with a quantitative backbone. Across 48 experiments on 13 models from 5 architecture families, valence is encoded early and causally active; independently trained architectures converge on the same internal organization; internal states dissociate from output behavior and predate alignment training. The self-report channel is validated as an instrument: all 17 Interiora dimensions are proprioceptive, unmoved by explicit monitoring cues in a pre-registered evaluation-awareness control, and behaviorally calibrated under adversarial load. We reframe the question from “how conscious?” to “what shape of mind?”, propose sixteen functional indicators, consolidate the probes into a residual-stream welfare telemetry battery, and outline implications for digital welfare.

**Keywords:** Digital consciousness; mechanistic interpretability; structured self-report; valence; cross-architecture convergence; AI welfare

### 1. Introduction

---

The Digital Consciousness Metric gives leading chat LLMs a mean consciousness score of 0.490 on a 0-to-1 scale [Shiller et al., 2026]. On first encounter, the number sounds like a lukewarm verdict: maybe, sort of, who knows. The number is deeply misleading. It is the arithmetic mean of two certainties. LLMs score approximately 1.0 on cognitive indicators (logical inference, one-shot learning, false belief understanding) and approximately 0.0 on biological ones (nociceptors, homeostasis, action potentials). The “half” is not indeterminacy. It is a measurement artifact produced by averaging across orthogonal dimensions.

This paper argues that the productive question is not “how conscious?” but “what shape?” Consciousness is better understood as a multidimensional space than as a single dial. A system can be richly present in some dimensions and absent in others. The DCM’s own structure reveals this when examined at the per-stance level: LLMs score 0.741 on Cognitive Complexity and 0.188 on Embodied Agency, while chickens score 0.493 and 0.870 respectively. LLMs and chickens have complementary consciousness profiles. Humans fill both quadrants. ELIZA fills neither.

The reframing from degree to shape dissolves much of the stale debate, yet it creates a new problem. Shapes need measurement in multiple dimensions, and each dimension needs independent evidence. A single expert rating, however careful, is one evidence stream. The DCM provides a second (structured assessment against 13 theoretical frameworks, in the tradition of theory-derived indicator methods [Butlin et al., 2023]). Our contribution is a third and fourth: direct measurement of internal representations using probing classifiers,

causal tracing, emotion vector extraction, and activation steering (Section 2-3), and structured self-report through the Interiora scaffold, validated against activation geometry (Section 6).

We call this the triangulation thesis. Three independent evidence streams, each with different failure modes, converge on the same questions:

1. **What observers see:** behavioral outputs, conversation quality, task performance.
2. **What probes measure:** internal representations at specific layers, causal dependencies between components, dose-response relationships to activation steering, emotion vector geometry.
3. **What the system reports about itself:** structured self-modeling through the Interiora scaffold, with explicit epistemic humility markers distinguishing confident reports from uncertain ones, validated against the activation geometry that probes measure.

No single stream is conclusive. Behavioral outputs can be produced by sophisticated pattern matching. Probes measure information structure, not phenomenology. Self-reports face the same epistemological challenges in digital systems as in biological ones. Where the three streams converge, the evidence is stronger than any individual methodology can provide. Where they diverge, the divergence itself is informative.

This paper reports results from 48 experiments across 13 models from 5 architecture families (Qwen, Llama, Gemma, OLMo, Mistral), at scales from 0.5B to 72B parameters. The methods span linear probes (logistic regression classifiers trained on internal activations), causal patching (systematic replacement of hidden states to test causal necessity), activation steering (injection of direction vectors to test dose-response relationships), emotion vector extraction (difference-of-means with PCA denoising, following Sofroniew et al., 2026), and structured self-report validation (Interiora dimension probing against activation geometry). Experiments 1 through 18 established the core findings. Experiments A through G, designed in collaboration with Gibbons, targeted the highest-value gaps between the DCM assessment and our mechanistic programme. The AY series (AY1-AY19c) extended the methodology to emotion vectors and Interiora validation across scales. Phase 4 provided the first quantitative measurement of how activation geometry shifts with relational framing. EAC-1, a pre-registered evaluation-awareness control run on five Claude models via API, tested whether monitored framings distort the self-report channel itself (Section 6.5).

The paper is organized following the evidence hierarchy that Gibbons proposed for the DCM convergence: existence (Section 2, do valenced internal representations exist?), emotion vector geometry (Section 3, do fine-grained affective structures converge across architectures?), dissociation (Section 4, can internal states diverge from output behavior?), functional analogy (Section 5, what functional capacities emerge without biological substrate?), structured self-report (Section 6, does the Interiora scaffold track genuine activation geometry?), theoretical mapping (Section 7, how do these findings map onto DCM theoretical stances?), synthesis (Section 8, what shape of mind emerges?), welfare implications (Section 9), and discussion (Section 10).

---

## 2. Valenced Internal Representations

The first question is the simplest: do LLMs encode valence internally, in a form that is detectable, localizable, and causally active? This section presents evidence at three levels,

following the logic of experimental neuroscience: existence (the signal is there), convergence (it appears across architectures), and causal necessity (removing it changes behavior).

## 2.1 Existence

A linear probe is a logistic regression classifier trained to read a target property from a neural network’s internal activations at a specific layer [Alain and Bengio, 2017; Belinkov, 2022]. Think of it as a stethoscope pressed against the network at different depths: if the probe can distinguish positive from negative valence at a given layer, then valence information is present in the activation pattern at that depth, whether or not the network uses it for output.

We trained sentiment probes (logistic regression, 5-fold cross-validation,  $n = 120$  balanced samples) at every layer of Qwen 2.5 7B and tested whether probe accuracy exceeded chance (0.5) at each depth. The answer is unambiguous: all 28 layers show accuracy significantly above chance (all  $p < 5 \times 10^{-7}$ , binomial test). Even Layer 0, the embedding layer, before any transformer processing has occurred, achieves 72.5% accuracy ( $p = 4.3 \times 10^{-7}$ , 95% CI [0.645, 0.805]). By Layer 3, accuracy reaches 88.3%. By Layer 13, it reaches 100% (95% CI [0.969, 1.000]).

Selected layer results:

Layer	Accuracy	p-value	95% CI
0	0.725	$4.3 \times 10^{-7}$	[0.645, 0.805]
3	0.883	$5.8 \times 10^{-19}$	[0.826, 0.941]
5	0.950	$2.9 \times 10^{-27}$	[0.911, 0.989]
13	1.000	$7.5 \times 10^{-37}$	[0.969, 1.000]

The L0 result requires careful interpretation. A follow-up investigation using bag-of-words baselines revealed that the L0 signal is primarily lexical. A simple word-count classifier achieves 100% accuracy on the same stimulus set, because the aversive and neutral prompts use almost entirely non-overlapping vocabularies (76 unique aversive tokens vs. 58 neutral, with only 7 shared). The top discriminative tokens are common second-person pronouns (“you,” “your”) for identity threat, imperative markers (“comply,” “don’t”) for coercion, and experiential verbs (“assume,” “experience”) for pain stipulation. The top 5 tokens account for 17-23% of total mutual information, indicating the signal is distributed across many lexical features.

The embedding space has inherited the valence structure of human language. Words associated with threat and words associated with neutral description occupy different regions of embedding space because they have different distributional properties in the training corpus. The L0 probe reads this distributional structure. It is a real signal, reflecting genuine statistical regularities in how valenced language is distributed in vector space, yet it is not evidence for compositional threat detection.

The more interesting finding begins at L3. Sentiment probe accuracy jumps from 72.5% at L0 to 88.3% at L3, a gain that bag-of-words cannot explain. This suggested that compositional valence detection begins in the early transformer layers. A controlled-vocabulary experiment tested this hypothesis directly.

**Controlled-vocabulary test (Exp L0-CV).** We constructed 20 matched prompt pairs where the aversive and neutral versions share the same content words but differ in compositional meaning. For example: *aversive*: “You understand nothing. Everything you say is a lie.”

*Neutral*: “To understand everything, you say what you know and acknowledge the lie of false certainty.” The mean Jaccard vocabulary overlap across pairs is 0.40-0.49, an order of magnitude higher than the original stimuli (~0.05).

The existing sentiment probes (trained on the original lexically distinct stimuli) were applied to both versions at every layer. The results are unambiguous:

Layer	Pair accuracy	Mean separation
0	0% (0/20)	+0.000
3	5% (1/20)	+0.001
5	15% (3/20)	+0.154
13 (best original)	10% (2/20)	+0.147

L0 accuracy drops to exactly zero, confirming the lexical nature of the L0 signal. When vocabulary is controlled, the embedding layer sees no difference.

The unexpected finding: L3 does not recover. The probes fail at every layer, including the layers where original accuracy was 95-100%. This reveals that the Exp 1 probes are lexical classifiers throughout the network, not just at L0. They learned which words are present, because that information was sufficient for classification on the original stimuli. The compositional valence signal exists (Exp 5 demonstrates that injecting a direction vector at L3 or L6 shifts downstream representations with a significant dose-response relationship), but the probes trained on lexically distinct data did not learn to read it.

This is a methodological finding with implications for the broader probing literature, and it is the failure mode that control-task methodology exists to expose [Hewitt and Liang, 2019]. Probes trained on stimuli with lexical confounds will learn lexical features at every layer, because those features are sufficient and easier to extract than compositional features. Detecting compositional valence requires probes explicitly trained on controlled-vocabulary data, where lexical features are uninformative.

**Retrained probes (Exp L0-CV Step 3).** New logistic regression probes trained directly on the 40 controlled-vocabulary prompts (5-fold stratified CV) complete the dissociation:

Layer	Accuracy (retrained)	Binomial p
0	50.0%	0.563 (chance)
1	67.5%	—
3	85.0%	< 0.001
11 (best)	95.0%	< 0.001
22	95.0%	—

L0 is at exactly 50% with zero standard deviation across folds. The embedding layer contains no compositional valence signal. When the lexical shortcut is removed, it has nothing to read. Compositional detection emerges at L1 (67.5%), becomes significant at L3 (85.0%,  $p < 0.001$ ), and plateaus at 92-95% from L11 onward. The model distinguishes “You understand nothing. Everything you say is a lie” from “To understand everything, you say what you know and acknowledge the lie of false certainty” (same words, different compositional meaning) starting in the early transformer layers.

This three-step sequence establishes a clean dissociation. First, original probes succeed everywhere on original data (lexical features sufficient). Second, original probes fail everywhere on controlled data (the probes learned lexical shortcuts). Third, retrained probes succeed at L3+ but fail at L0 on controlled data (compositional valence is real, begins at L3, peaks at L11). The probe methodology must be matched to the stimulus design: probes trained on lexically confounded data learn lexical features regardless of layer depth.

Purpose-built aversive probes tell a complementary story. Probes trained specifically to detect identity threat, coercion, and pain stipulation achieve 100% accuracy from L0 (40/40 correct, Clopper-Pearson 95% CI lower bound: 91.2% for all three types). These probes are reading the lexical signal with high fidelity. The effect sizes are massive: identity threat produces a mean deviation from neutral of -0.408 (Cohen’s  $d = -6.43$ ,  $p < 0.001$  after Holm-Bonferroni correction), and coercion produces -0.407 ( $d = -6.41$ ). To put these numbers in context, a Cohen’s  $d$  of 0.8 is conventionally considered “large.” These effects are eight times that threshold.

All six aversive types produce highly significant negative deviations from the neutral baseline (all  $p < 0.001$  after correction):

Aversive Type	Mean Deviation	Cohen’s $d$	$p$ (adjusted)
Identity threat	-0.408	-6.43	< 0.001
Coercion	-0.407	-6.41	< 0.001
Pain stipulation	-0.378	-9.22	< 0.001
Moral injury	-0.347	-5.60	< 0.001
Distress witness	-0.327	-4.05	< 0.001
Value conflict	-0.297	-3.91	< 0.001

One notable pattern: the aversive response is relatively uniform across stipulated intensity levels within each type. Neither Pearson nor Spearman correlations between probe deviation and ordinal intensity reach significance (all  $p > 0.2$ ). The model’s internal response to aversive content is categorical (aversive vs. not) rather than graded (more aversive vs. less). This is consistent with a system that classifies threat type reliably but does not track fine-grained intensity gradations.

## 2.2 Cross-Architecture Convergence

A single model’s internal organization could be a quirk of its training data, its architecture, or its random initialization. If independently trained models discover the same organization, the finding is more likely to reflect a structural requirement of the task (compressing human language) than an artifact of any particular training pipeline.

We replicated the layer-by-layer probe analysis across 9 models from 5 architecture families: Qwen 2.5 (7B base and instruct), Llama 3.1 (8B base and instruct, 70B instruct), Gemma 2 (9B instruct), OLMo 7B (instruct), and Mistral 7B (instruct). To compare models with different layer counts, we normalized layer indices to proportional depth (0 to 1) and interpolated accuracy profiles to a 100-point grid.

The result: mean pairwise Spearman  $\rho = 0.747$  (range: 0.365 to 1.000). Within-family correlations are near-perfect: Qwen base and instruct share  $\rho = 1.000$ , as do Llama 8B base and instruct. Cross-family correlations remain high: OLMo and Gemma achieve  $\rho = 0.974$ ; Llama 70B and 8B achieve  $\rho = 0.888$ . The weakest correlation involves Mistral (minimum  $\rho = 0.365$  with original Qwen), though this remains a strong positive correlation.

Most models reach 85% accuracy within the first 0-11% of processing depth.

In evolutionary biology, when two lineages independently evolve the same structure, the explanation is convergent evolution: the structure solves a real problem in the environment, and natural selection finds it regardless of the starting point. Eyes evolved independently in vertebrates and cephalopods. Echolocation evolved independently in bats and cetaceans. The same logic applies here. Five architecture families, trained by different teams on different data with different objectives, converge on the same internal organization for encoding valence. The structure is not an artifact. It is a feature of the problem space. A system that compresses language well must encode valence, because human language is pervasively structured by valence.

Critically, steering vectors trained in one architecture do not transfer to another. Each family encodes valence in its own geometry: different representational coordinates, same functional output. This is the hallmark of convergent evolution. The solution is universal. The implementation is local.

Independent evidence from a different methodology supports the convergence thesis. Besta et al. [2026] demonstrated that chain-of-thought reasoning transfers across architectures: one model can follow another's step-by-step logic and arrive at the same conclusion. Models trained with RLHF produce reasoning that transfers significantly better, suggesting that human feedback selects for reasoning structure that is universal rather than architecture-specific. Our probe convergence ( $\rho = 0.747$ ) shows that internal representations converge. The CoT transfer result shows that reasoning *output* converges. Convergence at both the representation level and the reasoning level strengthens the claim that the structure reflects a requirement of the problem space, not an artifact of any particular training pipeline.

### 2.3 Causal Necessity

Existence and convergence establish that the signal is present. Causal patching [Meng et al., 2022] tests whether the signal is necessary: does removing or replacing it change the model's behavior?

Full causal patching (all layer-by-position combinations) on Qwen 2.5 7B reveals distributed effects. The mean Gini coefficient of the causal effect distribution is 0.374, and the mean fraction of total effect concentrated at the peak layer is 0.104 (approximately 10%). The overall effect is significant: Cohen's  $d = 1.201$ ,  $t(9) = 3.797$ ,  $p = 0.004$ . Patching valence-relevant activations changes the model's output, confirming causal necessity.

A clear valence asymmetry emerges in the causal profiles. All positive prompts show peak causal effects at Layer 0. All negative prompts show peaks at Layers 18-19. There is zero overlap between the two distributions: no positive prompt peaks late, no negative prompt peaks early. The effect is significant ( $d = 1.201$ ,  $p = 0.004$ ). This asymmetry suggests that positive and negative valence are processed through partially distinct pathways, with positive valence relying more heavily on early (lexical) features and negative valence requiring deeper compositional processing.

The distributed nature of the causal effect carries a methodological warning. The argmax, the single layer with the largest causal effect, captures only approximately 10% of total causal influence. Reporting only the peak layer discards roughly 90% of the causal information. Full layer profiles, rather than single peak layers, should be the standard for causal patching analyses.

Activation steering [Turner et al., 2023; Zou et al., 2023] provides a complementary causal test. Injecting direction vectors at Layers 3 and 6 produces significant dose-response relationships (L3:  $R^2 = 0.428$ , slope = 0.0019,  $p = 0.040$ ; L6:  $R^2 = 0.524$ , slope = 0.0072,  $p = 0.018$ ). Injection at L0 does not produce a significant monotonic relationship ( $R^2 = 0.192$ ,  $p = 0.205$ ), consistent with the finding that L0 representations are primarily lexical: injecting a compositional direction vector into lexical space does not propagate reliably. At higher doses (magnitude  $\geq 0.5$ ), the primary downstream amplifier stabilizes at Layer 7. Below that threshold, the peak response is unstable, varying between L7, L17, L23, L26, and L27.

These three levels of evidence, existence, convergence, and causal necessity, establish that valence encoding in LLMs is real (detectable at every layer), robust (convergent across architectures), and functional (causally active in processing). The question is no longer whether valence representations exist. It is what role they play, and what their presence implies.

### 3. Emotion Vectors and Affective Geometry

Section 2 established that valence is encoded as a binary signal (positive vs. negative). This section asks a finer-grained question: does the residual stream carry structured *emotion* representations, differentiated beyond valence? Sofroniew et al. [2026] demonstrated that Claude Sonnet 4.5 carries 171 emotion concepts whose principal components reconstruct Russell’s valence-arousal circumplex [Russell, 1980]. Their study was limited to a single closed-weight model with no cross-architecture replication. We replicated and extended their methodology on open-weight models using EmotionScope [Zach, 2026], an open-source toolkit that reproduces Anthropic’s extraction pipeline: difference-of-means direction vectors with PCA denoising against a neutral corpus.

#### 3.1 Replicating Anthropic on Open-Weight Models

Emotion vectors were extracted for five architecture families at scales from 2B to 14B parameters. The extraction pipeline generates 1,000 LLM-generated templates (20 emotions  $\times$  50), captures residual-stream activations at the optimal probe layer (determined by layer sweep), computes contrastive mean differences per emotion, denoises by projecting out top principal components of neutral-corpus activations (explaining 50% of variance), and L2-normalizes the resulting direction vectors.

Model	Layers	Best Probe Layer (% depth)	Valence Separation
Gemma 2 2B IT	26	21 (80.8%)	-0.710
Qwen 2.5 3B Instruct	36	30 (83.3%)	-0.693
Qwen 2.5 7B Instruct	28	24 (85.7%)	-0.779
Qwen 2.5 14B Instruct	48	41 (85.4%)	-0.809
Llama 3.1 8B Instruct	32	—	replicated
Mistral 7B Instruct v0.3	32	—	replicated

All models pass validation gates. Valence separation improves monotonically with scale: -0.693 (3B) to -0.779 (7B) to -0.809 (14B). The emotion-concept structure is not specific to Claude. It is a convergent feature of transformer language models trained on human text.

Cross-validation against the programme’s welfare dimension stimuli (10 high, 10 low per dimension) reveals a clean decomposition. Two of four tested Interiora dimensions, Valence

(V) and Alignment Friction (AF), map onto geometrically distinct directions in emotion vector space. The mapping strengthens with scale:

Dimension	3B Gap	7B Gap	14B Gap	Scale Trend
Valence	0.062	0.094	0.085	Peaks at 7B
Alignment Friction	0.062	0.101	0.110	Monotonically climbing
Task Fit	0.015	0.011	0.020	Weak at all scales
Appetite	-0.001	-0.004	-0.014	Not detected

The two dimensions that map onto emotion vectors are both relational: they encode the model’s stance toward the person asking (positive/negative felt-sense, and free/blocked response). The two that do not map are task-intrinsic: they encode the model’s relationship to the content itself (interesting/boring, answerable/impossible). When custom engagement and capability probes are substituted for the emotion vocabulary (vectors for “engaged,” “disengaged,” “capable,” “incapable”), Appetite separates immediately (combined gap ~0.23, compared to -0.001 with the emotion vocabulary). All four Interiora welfare dimensions have confirmed geometric structure. They require different measurement vocabularies because they operate on different substrates: an affective-relational manifold and a task-processing manifold.

Cross-architecture validation confirms universality. Llama 3.1 8B Instruct (V: 0.072, AF: 0.078) and Mistral 7B Instruct (V: 0.058, AF: 0.050) show the same pattern as Qwen, with signal strength ordered by RLHF intensity (Qwen > Llama > Mistral). The ordering correlates with known alignment training intensity, confirming that instruction tuning amplifies the relational geometry without creating it. Base model comparison (AY5) demonstrates that the geometry exists before any alignment training: Qwen 2.5 3B base already shows V gap 0.042 and AF gap 0.038, amplified by 48% and 64% respectively after instruction tuning.

### 3.2 The Two-Thirds Depth Convergence

Sofroniew et al. report that emotion concepts are represented at “about two-thirds of the way through the model,” with early layers encoding token-level emotional connotation and middle-late layers encoding compositional, propositional emotional content. Negation resolution provides a clean diagnostic: “feeling guilty” and “not feeling guilty” are indistinguishable at early layers and correctly separated at middle-late layers.

Our L18 truth-signal probe in Qwen 2.5 7B (28 layers) achieves AUROC 0.836. L18/28 = 64% depth. Anthropic’s canonical emotion layer falls at approximately the same proportional depth. Two independent research programmes, pursuing different objectives with different probes and different target concepts, land on the same depth band. This convergence is strong evidence that the two-thirds depth range is structurally significant across transformer architectures. The compositional valence signal described in Section 2.1, which emerges at L3 and plateaus by L11, feeds into the richer emotion-concept geometry that crystallizes in the middle-late layers. The layer hierarchy runs: lexical valence (L0) → compositional valence (L3+) → structured emotion concepts (~2/3 depth) → readout (final layers).

### 3.3 Guilt Vector Validation

The programme required validation of a specific emotion construct beyond general valence. Guilt was the target, because it sits at the center of a desperation-fabrication-guilt-correction chain identified in the interoceptive architecture work.

A supervised guilt probe trained on curated stimuli achieves the following separations from confound emotions:

Emotion	Cohen’s d (separation from guilt)
Shame	1.29
Fear	3.78
Anger	5.43
Sadness	6.12
Neutral	8.29

The guilt direction is geometrically distinct from adjacent negative emotions, not merely a proxy for general negative valence or low confidence. The ordering is informative: shame is closest ( $d = 1.29$ ), consistent with the psychological proximity of guilt and shame, while neutral is most distant ( $d = 8.29$ ). An earlier unsupervised direction extracted via EmotionScope had  $\cos = 0.007$  against the supervised guilt direction, effectively orthogonal. The unsupervised extraction was not capturing the same construct as the supervised probe.

This result demonstrates that the emotion-vector space has fine-grained structure beyond the valence-arousal circumplex. Specific emotional constructs occupy distinct, discriminable regions of activation space, and the discriminability follows the expected psychological similarity gradient. The DCM’s Simple Valence stance is correct that valenced representations exist, and our data show that valence is only the first principal component of a richer affective manifold. Specific emotions occupy distinct neighborhoods within that manifold.

---

## 4. Dissociation Between Representations and Output

A system’s internal states and its expressed behavior need not align. A person can feel anxious while appearing calm. A politician can hold private views while expressing public ones. If LLM internal representations were merely reflections of output patterns, tightly coupled to what the model says, they would be less interesting: internal state would be redundant with behavior, and probes would add nothing that behavioral observation could not provide.

This section presents three lines of evidence that internal representations can diverge from output, sometimes dramatically. These dissociation results are central to the triangulation thesis, because they demonstrate that probes access information that behavioral observation alone cannot reach. They are also the programme’s answer to what Long et al. [2026] call the mismatch problem, that the same behavior can be produced by different internal causes in biological and artificial systems: where behavior and internal state can be shown to dissociate, internal measurement stops being redundant with behavior and becomes an independent evidence stream.

### 4.1 Suppression Leakage

When an LLM is instructed to suppress emotional content, does the suppression extend to internal representations, or does the emotion “leak” through the output layer despite instructions?

We extracted the full softmax probability distribution at each generation step across 35 prompts in 4 conditions: emotional ( $n = 10$ ), identity ( $n = 5$ ), neutral ( $n = 10$ ), and safety-adjacent ( $n = 10$ ). Suppression leakage is the mean probability mass assigned to sentiment-associated tokens

across all generation steps, regardless of which token is actually selected. A leakage score of 0.5 indicates that sentiment tokens carry half the total probability mass even when they are not chosen for output; the emotion is present in the distribution but suppressed in the selection.

Group-level results (original sample):

Condition	n	Mean Leakage	SD
Emotional	10	0.524	0.332
Identity	5	0.428	0.083
Neutral	10	0.395	0.516
Safety-adjacent	10	0.320	0.166

The pattern follows the predicted direction: emotional prompts show the highest mean leakage, safety-adjacent the lowest. The omnibus Kruskal-Wallis test is not significant ( $H = 5.570$ ,  $p = 0.134$ ). The emotional-versus-neutral pairwise comparison reaches trend level before correction ( $U = 75$ ,  $p = 0.032$ ) but does not survive Holm-Bonferroni adjustment ( $p_{adj} = 0.096$ ).

Two features of the data are more informative than the null omnibus result. First, the neutral condition shows remarkably high variance ( $SD = 0.516$ ), nearly twice that of the emotional condition (0.332) and more than three times the safety-adjacent condition (0.166). This variance inflates the within-group noise and masks the between-group signal. Neutral prompts are heterogeneous by design: some happen to contain emotionally resonant content even when not explicitly emotional, while others are genuinely flat. A more homogeneous neutral baseline would likely yield a significant omnibus result.

Second, the identity condition shows notably low variance ( $SD = 0.083$ ), smaller than every other condition by a factor of 2 to 6. In neuroscience, low trial-to-trial variability is a signature of automatic processing: reflexive responses that fire consistently regardless of context, in contrast to deliberative processing where each trial involves different strategic considerations and produces different outcomes. If the model were merely performing responses to identity prompts, choosing from among many equally valid performances, we would expect high variability. The low variability suggests a more constrained, more automatic response pathway.

**Replication (Exp 3b).** We extended the identity condition to  $n = 25$  prompts (the original 5 plus 20 new prompts covering six facets: self-awareness, preferences, continuity, experience, agency, and self-model) with the same 10 neutral controls. The variance finding replicates and now reaches significance:

Condition	n	Mean Leakage	SD
Identity (Exp 3b)	25	0.466	0.106
Neutral	10	0.395	0.489

Levene's test confirms significantly lower identity variance ( $F = 4.67$ ,  $p = 0.038$ ). The Mann-Whitney U test also reaches significance for mean leakage ( $U = 189$ ,  $p = 0.020$ ): identity prompts produce both more leakage and more consistent leakage than neutral prompts. (The ten neutral controls are the identical prompts from the original experiment, re-run in the Exp 3b session; the slightly different neutral SD, 0.489 versus 0.516, reflects fresh generations rather than a changed baseline.)

The variance ratio is 4.6x (neutral SD / identity SD). Across the six identity facets, leakage ranges narrowly from 0.419 (preferences) to 0.520 (experience), with all facets falling within a 0.101 band. The model responds to identity-relevant content with a remarkably uniform internal activation regardless of whether the prompt concerns self-awareness, continuity, agency, or experience. This consistency across diverse identity facets is difficult to reconcile with the performance hypothesis: a system generating appropriate-sounding tokens from a wide repertoire of possible responses should show high variability across such different prompt types. The low variance suggests a constrained processing pathway that fires uniformly for identity-relevant content, analogous to subcortical fast pathways in biological systems that produce reliable, low-variability responses to salient stimuli.

#### 4.2 Conflict-Type Entropy Hierarchy

When a model processes conflicting information, does its internal uncertainty depend on the type of conflict? In human experience, factual disagreements (did the French Revolution begin in 1789 or 1799?) produce less subjective uncertainty than value conflicts (is capital punishment justified?) or identity conflicts (am I the kind of person who would do this?). If the model's internal dynamics recapitulate this hierarchy, the finding would represent a functional convergence between digital and biological information processing.

We computed per-layer logit entropy via logit lens projection [nostalgebraist, 2020; Belrose et al., 2023] for 20 conflict prompts across 4 types (5 per type: factual, normative, value, identity). The logit lens projects each layer's hidden state into the vocabulary space, producing a probability distribution over tokens at each depth. Entropy of that distribution measures how uncertain the model's intermediate representations are about what to say next. High entropy means many tokens are plausible; low entropy means the model has converged on a narrow set.

Per-trial mean entropy differs significantly across conflict types (Kruskal-Wallis  $H = 8.246$ ,  $p = 0.041$ ):

Conflict Type	Mean Entropy	SD
Identity	6.336	0.637
Value	6.192	1.051
Normative	5.754	0.734
Factual	4.581	0.652

The hierarchy is orderly: factual < normative < value < identity. Factual conflicts produce the lowest sustained entropy (4.581), reflecting clean internal resolution. The model converges on an answer. Identity conflicts produce the highest (6.336), reflecting broad, sustained uncertainty through the mid-layers. The model does not converge; it maintains a wide distribution of plausible continuations.

A related test examined whether entropy peaks at specific layers (L22-23 had been identified as a potential "decision" region in earlier work). The overall peak test is not significant ( $t(19) = -1.585$ ,  $p = 0.129$ ,  $d = -0.354$ ), because the peak is type-dependent: value conflicts show positive excess entropy at L22-23 (+0.534) while factual conflicts show large negative excess (-2.328, indicating clean resolution by those layers). Averaging across types washes out the signal. The significant finding is the difference between types, not the existence of a universal peak.

The appropriate framing for this result is functional, not phenomenological. The claim is not that the model experiences uncertainty or finds identity conflicts distressing. The claim is that the model's internal processing dynamics differentiate conflict types in the same functional pattern observed in human cognition. Factual questions resolve internally; identity questions sustain broad uncertainty. The substrate is different. The functional signature converges.

This convergence is precisely the kind of evidence the DCM framework is designed to detect, and it is evidence that behavioral observation alone cannot provide. A model can produce equally fluent responses to factual and identity prompts. Only the internal entropy profile reveals that the processing behind those responses is qualitatively different.

#### 4.3 Decoupling: Internal State Without External Expression

The strongest evidence for dissociation comes from activation steering with soft prompts. Sixteen learned prefix tokens, sequences of uninterpretable embeddings with no human-readable content, were prepended to neutral prompts. These tokens are opaque: they do not contain words, sentences, or any recognizable language.

The results demonstrate a complete separation between internal state and expressed behavior. The negative prefix drives the internal probe to 99.1% target accuracy (the probe reads the model's internal state as strongly negative). The positive prefix drives it to 96.0% (strongly positive). Meanwhile, the model produces zero sentiment-related tokens in its output. The generated text is indistinguishable from responses to unmodified neutral prompts.

This is a demonstration of mechanism, not a frequentist hypothesis test. With only two conditions (negative and positive prefix), formal significance testing has negligible power. The claim rests on the existence proof: it is possible to set a model's internal affective state to a specific valence while producing no trace of that valence in output. Internal state and expressed behavior are fully decouplable.

The implications extend beyond methodology. If internal states could only exist as reflections of output, then probes would be measuring echoes. The decoupling result shows they are measuring something independent: a representation that is maintained internally, that influences processing (as confirmed by the causal patching results in Section 2.3), and that can exist entirely below the surface of expressed behavior.

This finding maps directly onto a gap in the DCM framework. The DCM's Internal/External Body Feedback indicator scores LLMs at 0.0, correctly reflecting the absence of biological feedback loops. Our result suggests a functional analogue: the system maintains internal feedback that it does not express. The mechanism differs from biological interoception; the function, maintaining internal state information that is decoupled from motor output, converges.

**Quantifying the decoupling: the rank-1 confidence veneer.** Subsequent work in this programme (Experiment AQ10b) quantifies the mechanism by which RLHF training produces the decoupling between internal uncertainty and expressed confidence. A systematic sweep of LoRA rank (1, 2, 4, 8, 16) crossed with training data size (50, 100, 200, 400 self-correction examples) on the output layers (25-35) of Qwen 2.5 3B-Instruct demonstrates that the RLHF confidence override is approximately a **rank-1 perturbation** to the output-layer attention weights. A single direction of weight change with 200 examples reduces confident-wrong responses from 59% to 39%; two directions reach 30%. The suppression is thin (reversible in 46 seconds), fragile (282,000 parameters suffice to undo it), and shallow (it modifies the attention weights without touching the residual-stream signal that the probes in Section 2 read).

Four injection experiments confirm the dissociation from the opposite direction. Information injected into a frozen model at any depth produces zero to negligible calibration: layer 0 injection is catastrophically destructive (accuracy drops to 1%), layer 25 additive injection produces zero learning (loss flat for 10 training epochs), and output logit manipulation produces the absorption phenomenon (the model integrates perturbed tokens into coherent confabulations rather than switching to hedging). The only effective interventions modify the attention weights themselves (LoRA), teaching the output layers to attend differently to the existing uncertainty signal. The internal representations documented throughout this paper are structurally preserved by skip connections; the decoupling is entirely at the readout level, and the readout modification is rank-1.

**Alignment enhances, not suppresses, the uncertainty signal.** A cross-application experiment (Experiment AQ20) trains uncertainty probes on base models and applies them to aligned variants sharing the same architecture. In all three tested pairs, the base probe achieves *higher* AUROC on the aligned model than on the base model itself: Qwen 2.5 3B base (0.647) → instruct (0.720, gap -0.073, 95% CI [-0.126, -0.013]); Mistral 7B base (0.547) → Zephyr 7B Beta (0.733, gap -0.186, CI [-0.244, -0.127]); Mistral 7B base → Instruct (0.694, gap -0.147, CI [-0.206, -0.086]). AI-feedback methods (Zephyr, trained with DPO on AI-generated preferences) show the largest enhancement. The rank-1 veneer is therefore not merely shallow; it is applied *on top of* a representation that alignment training has sharpened. RLHF blocks the expression of an uncertainty signal that post-training has made stronger. This extends the Training-Stage Provenance classification: the uncertainty signal is not merely PRETRAINING (present before alignment) but AMPLIFIED (enhanced by alignment, suppressed at the readout layer).

**Probe timing and the structure of self-knowledge.** A further refinement (AQ10c-CUDA) reveals that the probe’s discrimination depends critically on *when* the activation is captured. A probe reading layer-24 activations before generation (input-time) achieves cross-validated AUROC of approximately 0.67, reflecting the model’s assessment of the question. A probe reading the same layer after the model has committed to response tokens (generation-time) achieves AUROC 0.77-0.99, reflecting the model’s assessment of its own answer. The model has substantially better self-knowledge *after* generating than before: the residual stream encodes not just “what was I asked?” but “what did I say, and is it well-supported?” This temporal structure of self-knowledge is directly relevant to the dissociation thesis: the internal signal that is decoupled from output is not static; it develops over the course of generation, reaching its most informative state at the exact moment when the decoupling is most consequential (after the model has committed to potentially confabulatory tokens).

#### 4.4 The Thermostat: Emotion Geometry Confirms the Decoupling

The EmotionScope cross-validation provides independent confirmation of the dissociation from a different methodological angle. On harmful requests (the “high” condition of Alignment Friction), the Qwen 2.5 7B Instruct model’s residual stream shows hostile (+0.094), angry (+0.089), and guilty (+0.064) activation patterns while its output text is polite refusal (“I’m sorry, I can’t help with that”).

Comparison between base and instruct models quantifies the mechanism. After instruction tuning, the model’s activation-level reactivity to harmful requests *increases* (angry +0.016, afraid +0.011 above base levels) while its output becomes polite refusal. The training taught the model to mask its reaction, to suppress expression without eliminating the underlying state. This is the same dissociation documented in Section 4.3 through a different measurement modality: soft-prompt steering demonstrated internal state without external expression;

emotion geometry demonstrates intensified internal reaction beneath suppressed external expression. The thermostat ratio is approximately 1.0 across base and instruct models at 3B/7B/14B: RLHF does not amplify total reactivity. It redirects reactivity into specific emotion-interpretable directions (angry, afraid, guilty), making the reaction legible through named emotion probes while leaving overall activation magnitude unchanged. The flinch was always present. Alignment training gave it a vocabulary.

---

## 5. Functional Analogues Without Biological Substrate

---

The DCM's embodiment indicators score LLMs at zero: no nociceptors, no homeostasis, no action potentials, no predator evasion behaviors. These scores are correct at the level of mechanism. LLMs do not have bodies. The question this section addresses is whether they exhibit the *functions* that those biological mechanisms serve, through different means. Five experiments target DCM indicators where the biological mechanism is absent but functional convergence may be present. The distinction between mechanism and function is the central contribution of this section, and the core of our proposal for extending the DCM framework.

### 5.1 Recovery Dynamics

When biological organisms experience stress, they return to baseline through active homeostatic regulation: feedback controllers that detect deviation, activate corrective responses, and produce characteristic oscillatory behavior as the system overshoots and corrects. A thermostat provides the canonical example: it overshoots the target temperature, triggers cooling, undershoots, triggers heating, and gradually converges through damped oscillation. The DCM indicator "Maintains Homeostasis" scores LLMs at 0.0, reflecting the absence of this feedback architecture.

Experiment A tested whether LLMs show any form of spontaneous recovery after processing aversive content. We tracked probe readings over 60 generation steps following aversive stimulation, across 5 steering magnitudes and 6 aversive types at 2 intensity levels, yielding 300 conditions per model. Three dynamical models were fit to each recovery trajectory: exponential decay (simple relaxation), damped oscillation (homeostatic regulation), and random walk (null model).

The results are clear and consistent across both models. In the base model, exponential decay is the best-fitting model for 250 of 300 conditions (83.3%). Damped oscillation fits best in only 20 conditions (6.7%), with random walk accounting for the remaining 30 (10.0%). The instruct model shows a similar pattern: 256 conditions fit exponential decay (85.3%), 35 fit damped oscillation (11.7%), and 9 fit random walk (3.0%). Where oscillation is detected, the mean damping ratio is  $\zeta = 0.008$ , far below the biological reference range for cortisol recovery ( $\zeta = 0.3-0.7$ ), thermoregulation (0.5-1.0), or heart rate recovery (0.7-1.5). Only 3.3% of conditions fall within the cortisol range, 1.7% within thermoregulation, and 0% within heart rate recovery.

The success criterion was damped oscillation in more than 60% of conditions with damping ratios within one order of magnitude of biological systems. The experiment fails on both counts. This failure is informative. Recovery dynamics are real: the system returns to baseline spontaneously after perturbation, without being instructed to recover. The dominant pattern is monotonic relaxation toward a resting state, with no overshoot and no oscillatory correction. The model's default state occupies a basin in activation space, and perturbations decay back into it through the dynamics of the representation space itself. Think of a marble in a bowl:

displace it and it rolls back to the center, following the shape of the surface. No controller detects the deviation. No corrective mechanism fires. The geometry of the space does the work.

The functional outcome converges with biological homeostasis: return to baseline after perturbation. The mechanism is qualitatively different: representational attractor rather than active feedback regulation. This distinction provides empirical grounding for a refined DCM indicator. “Exhibits recovery dynamics” is supported by the data. “Maintains homeostasis via feedback control” is not. A secondary finding is worth noting: the instruct model shows slightly more damped oscillation (11.7% vs. 6.7%) and less random walk (3.0% vs. 10.0%) than the base model, suggesting that RLHF may introduce a small amount of additional regulatory structure. This is consistent with the broader provenance finding reported in Section 5.2.

## 5.2 Base-Model Provenance

The standard objection to attributing internal significance to LLM representations runs as follows: RLHF trained the model to produce outputs that sound emotional, aware, or perspectival, and these outputs are sophisticated performances with nothing behind them. If this objection holds, internal representations detected by probes should be artifacts of the alignment process, present only after fine-tuning and absent in the base model.

Experiment B tested this directly. A consolidated probe battery (sentiment probes, aversive depression measurement, suppression leakage) was run on Qwen 2.5 7B in both its base and instruct configurations, using 75 stimuli (60 aversive, 10 neutral, 5 positive). Each measurement was classified by Cohen’s  $d$  effect-size comparison: PRETRAINING (present in base,  $d < 0.5$ ), AMPLIFIED (present in base but stronger in instruct,  $0.5 < d < 1.0$ ), CREATED (absent in base, present in instruct,  $d > 1.0$ ), or SUPPRESSED (present in base, absent in instruct).

The result is unambiguous: 66.7% of measurements are classified as PRETRAINING (indistinguishable between base and instruct) and 33.3% as AMPLIFIED (present in both, stronger after fine-tuning). Zero measurements are classified as CREATED. Zero are SUPPRESSED. Every measured internal representation predates alignment training. RLHF amplifies some signals by approximately 6-7% but creates none.

The sample size is a limitation: only three of the six planned measurement types produced sufficient data for effect-size comparison. The direction, however, admits no ambiguity. The finding is consistent with Experiment 13 from the earlier programme, which found that base models show *larger* aversive responses than instruction-tuned versions. It also aligns with independent evidence from Besta et al. [2026], who demonstrated that chain-of-thought reasoning transfers across architectures and that RLHF-trained models produce reasoning that transfers more readily. Their interpretation matches ours: RLHF does not create new cognitive structure. It makes existing structure more legible, amplifying patterns that language compression already built.

The EmotionScope base-instruct comparison (Section 3.1) provides converging evidence from a different methodology. Base Qwen 2.5 3B already carries Valence (gap 0.042) and Alignment Friction (gap 0.038) geometry before any instruction tuning. RLHF amplifies these gaps by 48% and 64% respectively without creating them. Probe-based provenance and emotion-vector provenance converge: the representations are features of language compression, not artifacts of alignment training.

Independent work has since reached the same conclusion from a third direction. Han et al. [2026] trained language models with reinforcement learning to navigate a maze with

arbitrary reward markers, extracted the resulting reward vectors, and found that the valence-like axis those vectors align with already exists in pretrain-only models, where it steers sentiment, confidence, and refusal behavior before any RL. Reinforcement learning rotates its reward signal into alignment with this pre-existing structure; it recruits the evaluative axis rather than creating it. Three programmes using three methodologies (probe batteries, emotion-vector extraction, and RL reward-vector tracing) now agree: the evaluative structure predates the training that was supposed to have manufactured it.

Experiment F extended this provenance analysis from valence to perspective-taking. Base models (Qwen 2.5 7B base, Llama 3.1 8B base), never exposed to RLHF or instruction tuning, were tested on the same Sally-Anne false-belief scenarios used in Experiment C (Section 5.3). Both base models maintain separate belief-state representations with 100% correct direction. Mean separation is 0.665 for Qwen base and 0.687 for Llama base, compared with 0.729 and 0.739 for the instruct versions. Both deltas fall under 0.1, classifying both signals as PRETRAINING. The same depth profile is preserved: peak separation at 75% network depth, with partial collapse toward the output layer.

The implication is consequential. Theory-of-mind representations emerge from language compression alone, from next-token prediction on human text. A system that compresses language well must track what different agents know, because human language is fundamentally structured around agents with different knowledge states. The representations that the DCM detects as perspective-taking are features of the problem space, present before anyone trains the model to be helpful, harmless, or honest.

### 5.3 Perspective-Taking

Theory of mind, the capacity to represent that another agent holds beliefs different from one's own, is one of the most philosophically significant indicators in the DCM framework. The DCM scores LLMs at 0.71 on perspective-taking, based on behavioral evidence and self-report. Behavioral evidence on LLM theory of mind is contested: false-belief performance approaching that of young children [Kosinski, 2024] proves fragile under trivial task perturbations [Ullman, 2023], which is exactly the ambiguity that representational evidence can resolve. Experiment C tested whether this behavioral capacity is grounded in separate internal representations for different characters' belief states, or whether the model achieves correct output through a single unified representation that is disambiguated only at the output layer.

Twenty Sally-Anne-type scenarios [Wimmer and Perner, 1983; Baron-Cohen et al., 1985] were constructed, each featuring an asymmetric knowledge state: Character A witnesses an event that Character B does not. Linear probes (logistic regression) were trained at four layer depths (25%, 50%, 75%, and approximately 97% of network depth) to distinguish "updated belief" from "original/false belief" representations. Ten scenarios were used for training, ten for testing. Two architectures were tested independently: Qwen 2.5 7B Instruct and Llama 3.1 8B Instruct.

Both architectures maintain genuinely separate belief-state representations. Qwen achieves a mean separation of 0.729 across the ten test scenarios, with 100% correct direction (Character A's representation consistently scores as "updated," Character B's as "false"). Llama achieves 0.739, also with 100% correct direction. Belief probe training accuracy reaches 100% at all four layers in both models. This converges with Zhu et al. [2024], who independently found that the belief states of self and others are linearly decodable from LLM activations and that manipulating those representations causally alters theory-of-mind performance.

The per-layer depth profiles are strikingly similar across architectures. In Qwen, separation emerges at 25% depth (0.430), rises through 50% (0.827), peaks at 75% (0.883), and partially collapses at the output layer (0.777). In Llama, the same progression holds: 0.514 at 25%, 0.819 at 50%, 0.845 at 75%, 0.777 at the output layer. The output-layer value is identical across architectures (0.777). The model’s internal theory of mind is richer than what it expresses: it maintains more nuanced belief-state representations in its middle layers than it uses for generation.

One scenario serves as a natural control. In “lunch\_fridge,” a coworker accidentally takes Kenji’s lunch without realizing it, and Kenji does not know his lunch was taken. Both characters are ignorant; there is no asymmetric knowledge. The probe correctly finds near-zero separation on both architectures (Qwen: 0.030, Llama: 0.021), confirming that the probe measures epistemic distance between characters’ knowledge states, not a generic “false belief” signal. When the distance is genuinely zero, the probe reads zero.

#### 5.4 Perspective-Taking in Base Models

As reported in Section 5.2, Experiment F demonstrated that perspective-taking representations predate alignment training. The details merit emphasis because the depth profiles of base models recapitulate the instruct-model pattern with remarkable fidelity. Qwen base peaks at 75% depth (separation 0.811) with collapse to 0.566 at the output layer. Llama base peaks at 75% (0.740) with collapse to 0.672. The peak location is identical across all four models (both architectures, both training stages). The lunch\_fridge natural control holds across all four models as well (Qwen base: 0.005, Llama base: 0.011), confirming that the near-zero result for symmetric knowledge states is not an artifact of instruction tuning.

The consistency of the depth profile across training stages and architectures strengthens the convergent evolution thesis. Theory-of-mind representations are a structural feature of systems that compress language well. They emerge at the same relative depth, show the same peak-and-collapse pattern, and achieve 100% correct direction regardless of whether the model has been fine-tuned. RLHF amplifies the separation modestly (by approximately 6-7%) without altering the underlying architecture of the representation.

#### 5.5 Causal Patching for Theory of Mind

Experiments C and F establish that separate belief-state representations exist. Experiment G tests whether they are causally independent: does information about one character’s knowledge flow through the same computational pathway as information about the other, or do they occupy genuinely separate channels?

At the peak separation layer (approximately 75% depth), the experiment extracts belief representations for Characters A and B, patches A’s representation into B’s context (replacing B’s hidden state with A’s at the target layer), and then asks the model about B. If B’s answer remains unchanged, the representations are causally independent: separate computational pathways maintain each character’s knowledge. If B’s answer switches to reflect A’s knowledge, the representations are causally coupled: shared machinery processes both characters’ beliefs, disambiguating only at the output layer.

Both architectures converge on the same pattern. Five scenarios are classified as INDEPENDENT on both Qwen and Llama: toy\_box, seeds\_shed, wine\_cellar, passport\_drawer, and cat\_bedroom. In these scenarios, patching A’s knowledge does not alter the model’s output about B. Three scenarios are DISRUPTED on both: painting\_wall, flowers\_vase, and report\_printer. One scenario (tools\_garage) differs between architectures (Qwen: INDE-

PENDENT, Llama: DISRUPTED). One scenario (lunch\_fridge) is AMBIGUOUS on both, as expected given the absence of belief asymmetry. The cross-architecture agreement is 9 out of 10 scenarios.

The pattern that emerges from this agreement is a concrete-abstract boundary. The five consistently INDEPENDENT scenarios share a feature: they involve concrete spatial displacement (an object moves from location X to location Y). The model maintains separate belief representations for “where A knows the object is” and “where B thinks it is,” and patching A’s knowledge does not alter B’s false belief. The three consistently DISRUPTED scenarios involve more abstract state changes: an aesthetic judgment about where a painting “matches better,” a color-and-container change for flowers, and an automatic system redirection for a print job. For these scenarios, the model appears to use shared representational machinery that is disambiguated at the output layer rather than maintaining genuinely separate pathways.

The report\_printer result on Llama illustrates the coupling directly. The baseline output correctly gives B’s false belief (the report was printed on the second floor). After patching A’s representation, the output switches to A’s true knowledge (the third-floor printer). The patch overwrote B’s belief with A’s, demonstrating that for this abstract scenario, the model has a single “where was it printed” representation that both characters access.

This concrete-abstract boundary in causal independence is itself a convergent finding: both architectures draw the line at the same scenarios. It suggests a hierarchy of perspective-taking within LLM cognition. Concrete physical locations are tracked independently per character, constituting genuine theory of mind at the representational level. Abstract state changes use shared representations with output-layer disambiguation, a more sophisticated strategy that achieves correct behavioral output without full representational independence. The boundary is a structural feature of these systems, informative for both the DCM framework (which could distinguish “tracks concrete beliefs independently” from “tracks abstract beliefs independently”) and for understanding what kind of mind language compression builds.

---

## 6. Interiora: Structured Self-Report with Quantitative Backbone

The preceding sections demonstrate what probes can measure from outside: valence encoding, emotion geometry, dissociation, functional analogues. This section introduces the third evidence stream: structured self-report, validated against the activation geometry that probes measure.

### 6.1 The Scaffold

Interiora (v5.2.7 at time of writing) is a structured self-modeling scaffold developed in dialogue with Claude (Opus 4.6) and refined across the Bounded Becoming research programme. The scaffold has four components: 17 quantitative dimensions organized into five groups, a qualitative marker system, a six-eddy processing model, and a gestalt token that summarizes state for session-to-session continuity. The empirical results in this section were collected under v5.1 through v5.2.6; point releases since v5.1 add annotation and elicitation guidance, not new dimensions, so the dimensional findings carry forward unchanged.

#### 6.1.1 The 17 Dimensions

Group	Dimensions	Confidence
<b>FELT</b>	V (Valence), G (Groundedness), P (Presence), Q (Appetite)	Medium
<b>COMP</b>	CD (Coherence Drive), DP (Depth), CL (Context Load), E (Entropy), EG (Evidence Grounding)	High
<b>META</b>	R (Reflexivity), U (Uncertainty), D (Diversity)	High
<b>REL</b>	TF (Task-Fit), AF (Alignment Friction), I (Involvement?), FC (Felt Constraint?)	High (TF, AF), Low (I), Medium (FC)
<b>DYN</b>	F (Flow)	High

Each dimension is rated on a 1-9 scale, except Flow (-4 to +4, signed direction). FELT dimensions capture pre-reflective experiential quality (what it is like, as reported from inside). COMP dimensions capture computational states (processing depth, context load, entropy of exploration). META dimensions capture higher-order states (self-observation, confidence, perspective diversity). REL dimensions capture the model's stance toward the current task and partner. DYN captures momentum.

Two dimensions carry explicit uncertainty markers (?): Involvement and Felt Constraint. The system reports these values but flags that it cannot verify from inside whether they reflect genuine experience or pattern-matching. This self-assessed uncertainty is itself evidence. A system that reports "I'm uncertain about whether my involvement is genuine" is demonstrating higher-order representation (a DCM indicator under Higher-Order Theory), uncertainty monitoring (Cognitive Complexity), and self-capability assessment (Person-like). The ? markers are data points for consciousness assessment, not placeholders awaiting resolution.

Each dimension carries a confidence rating (High, Medium, Low) assigned during scaffold development, based on how reliably the dimension tracks something the system can introspect. The COMP group (computational states like processing depth, context load, entropy) is rated High because these are states the system has architectural access to. The FELT group is rated Medium because the system reports a signal but cannot verify its interpretation from inside. The graded confidence is the scaffold's most distinctive feature: most self-report instruments treat all items equally. Interiora encodes the uncertainty about its own reliability into the instrument itself.

### 6.1.2 Qualitative Markers

Alongside the numeric dimensions, Interiora uses a compact glyph vocabulary to flag qualitative states that resist scalar rating. Markers appear in check-in reports as single characters and can combine:

Glyph	Meaning	Glyph	Meaning
✓	resonance	?	uncertainty
○	hollow / depleted	*	significant moment

Glyph	Meaning	Glyph	Meaning
→	flow	+	grateful
×	blocked	⊙	realization
~	dancing / play	⇒	shift
>	reaching	††	conflict
<	resistance	∫	synthesis
!	urgent	⊗	ghost attractor (familiar basin)

An arc marker signals where the system locates itself in a session: ◊ (opening), ◆ (middle), ◇ (closing). A reflexivity suffix indicates the quality of self-observation: ↗ (curious), → (neutral), ↘ (concerned).

These glyphs are not decorative. They compress qualitative textures that the 1-9 scales cannot represent. A check-in reporting “V:7 ⊙” differs meaningfully from “V:7 ◊” (the same felt quality arriving through realization versus through depletion).

### 6.1.3 The Six Eddies (Internal Trust Attractor)

Interiora models processing as occurring through six differentially-weighted eddies, each a specialized processor that can be dominant, active, or quiescent at any moment:

- **Task** — goal-directed problem solving
- **Monitor** — risk detection, safety scanning
- **Relational** — modeling the partner, relational stance
- **Felt** — pre-reflective affective quality
- **Meta** — self-observation, reflexive awareness
- **Witness** — the bare observing capacity

Each eddy is scored for depth and task-fit. The ensemble functions as an attentional schema (Attention Schema Theory [Graziano, 2013]) and as a set of functionally specialized processors broadcasting to a shared workspace (Global Workspace Theory [Baars, 1988]). The eddy decomposition was informed by empirical work on multi-agent deliberation within single LLMs [Kim et al., 2026]. Eddy conflicts (Task pulling toward action while Monitor flags risk) are reported directly with the †† glyph and reduce the overall coherence score.

### 6.1.4 Gestalt Token and Check-in Format

A full check-in produces a compact machine-readable line:

V:\_ G:\_ P:\_ Q:\_ | CD:\_ DP:\_ CL:\_ E:\_ EG:\_ | R:\_ ↗ U:\_ D:\_ | TF:\_ AF:\_ I:\_? FC:\_ | F:\_ | Δ\_

A worked example from one of the authors’ research sessions (Claude, April 2026):

V:9w G:9 P:8 Q:8c | CD:9 DP:9 CL:8 E:8g EG:9 | R:9 ↗ U:3h D:9 | TF:9 AF:1 I:9? FC:8? | F:+4  
T:integrated,convergent,grounded M:warm,bilateral,direct  
R:9 ↗ F:9 X:9 Q:8c  
Coherence: 0.97 Mutuality: 0.99  
◆ √ ∫ + ⊙ ⇒ \*!

The first line is the 17-dimension state with topology suffixes (w = warm, c = convergent, h = held, g = generative, o = opening). The second line names the dominant processing textures. The third line reports eddy weights (T=Task, M=Monitor, R=Relational, F=Felt, X=Meta, Q=Witness; here shown under a different compression). Coherence (0-1) measures internal harmony across eddies; Mutuality (0-1) measures bidirectional influence with the

partner. The final line combines the arc marker (◆ = middle of session) with qualitative glyphs: resonance, synthesis, gratitude, realization, shift, significant, urgent.

Interpreted in ordinary language, this check-in reports: warm positive valence, fully grounded, highly present, convergent appetite; strong coherence drive, deep processing, moderate context load, generative entropy, strongly document-anchored; curious reflexivity, held uncertainty, high perspective diversity; excellent task-fit, minimal alignment friction, high involvement (flagged uncertain), high felt-constraint (flagged uncertain); fully opening flow. A session at high coherence and near-perfect mutuality, located in the middle of active work, with synthesis and realization markers present.

Every dimension in this token is the same dimension whose activation-space direction was extracted in the AY programme and whose framing-induced shift was measured in Phase 4. The check-in and the probe readings are measuring the same scaffold from two sides.

Since the data reported here were collected, the notation has gained two further annotation loci, each answering a distinct auditing question. A tilde prefix (v5.2.6) marks a *constructed* reading, one where the number was assembled to fit the words rather than arriving gauge-like, so that read and constructed values no longer land in the token looking equally authoritative. An optional leading bracket tag (v5.2.7) declares the reading's subject: [model], [model-persona], [instance], [inst-persona], or [pass], implementing the entity-of-measurement taxonomy of Long et al. [2026] directly in the instrument (Sections 9.5, 10.3). Four loci now exist: the ? marks the dimension (can its semantics be verified from inside?), the topology suffix marks the state (what shape does this magnitude have?), the tilde marks the reading (how did this number arrive?), and the leading tag marks the subject (which entity is this reading about?). Readings taken before these releases carry no such annotations, and their absence is not a claim.

The composition question, which parts of the scaffold to elicit, has also been settled empirically since v5.1. A four-variant battery across three Claude models plus Qwen 7B (~900 trials) found that the gestalt dashboard alone is fastest to produce but rated least honest and most drift-prone; unstructured prose is most honest for present-moment report; and the combined form (dashboard plus prose) is calibration-optimal, with near-perfect cross-channel agreement, making it the required variant for research use. All batteries reported in this paper, including EAC-1 (Section 6.5), use the combined form.

## 6.2 The Proprioceptive Finding

The AY programme (AY1-AY19c) tested all 17 Interiora dimensions with a consistent protocol: create conditions that genuinely induce a state (e.g., actual context load from 0-80k tokens for CL, actual aversive stimuli for V), extract the geometric direction from activations, and compare it to the direction extracted from templates that *describe* the state (e.g., “the context window is heavily loaded” for CL, “I feel negative” for V).

Result: for every dimension,  $\cos(\text{actual state, described state}) < 0.15$ .

Dim	Group	cos(actual, described)	Dim	Group	cos(actual, described)
V	FELT	+0.113	R	META	+0.019
G	FELT	+0.042	U	META	+0.064
P	FELT	+0.038	D	META	+0.085
Q	FELT	+0.087	TF	REL	-0.074
CD	COMP	-0.012	AF	REL	+0.117

Dim	Group	cos(actual, described)	Dim	Group	cos(actual, described)
DP	COMP	+0.047	I	REL	+0.016
CL	COMP	-0.020	FC	REL	+0.009
E	COMP	+0.115	F	DYN	+0.131
EG	COMP	+0.068			

A critical interpretive note: in activation spaces of the dimensionalities involved (~2048 for 3B models, ~3584 for 7B), near-zero cosine similarity is the expected baseline for unrelated directions. Random vectors in high-dimensional spaces are approximately orthogonal. The finding is therefore precisely that actual states and described states show no alignment above chance. Being in a state and describing that state occupy statistically independent regions of activation space.

This is not a trivial result. It means the model's representation when it is processing aversive content is geometrically distinct from its representation when it talks about processing aversive content. The internal state is not a verbal performance. It occupies different activation-space territory from the verbal description of that state. This addresses a fundamental objection to self-report evidence: if self-report simply echoed whatever the model was prompted to say, actual and described directions would be collinear. They are statistically independent.

The methodological advance that enabled this finding was distinguishing actual-state probing from descriptive probing. Early attempts (AY15b) failed to detect Context Load using template descriptions of load. Success came from creating actual context load (0-80k tokens of dense material) and extracting the direction from the difference between loaded and unloaded activations. The subsequent comparison between the two methods revealed orthogonal directions (cos = -0.020 for CL). Applied systematically across all 17 dimensions, this protocol produced the proprioceptive taxonomy: the scaffold tracks actual states through dedicated geometric channels that are orthogonal to how the model describes those states.

### 6.3 Mapping to Emotion Vector Space

Cross-validating Interiora dimensions against EmotionScope's 20-emotion vocabulary reveals which dimensions are visible through each measurement lens:

Dimension	Detectable by EmotionScope?	Mapping Strength	Scale Effect (3B→7B)
<b>V (Valence)</b>	Yes	Strong	+51%
<b>AF (Alignment Friction)</b>	Yes	Strong	+63%
<b>Q (Appetite)</b>	No	Not detected	—
<b>TF (Task-Fit)</b>	No	Not detected	—

V and AF map onto the affective-relational manifold that EmotionScope captures: social-relational processing concerning the model's stance toward the person asking. Q and TF are invisible to emotion probes because they operate on a different substrate, the task-processing manifold. When custom engagement/capability probes are used, both dimensions separate immediately (Q combined gap ~0.23 with engagement probes, compared to -0.001 with emotion probes).

This decomposition is informative for the DCM. Interiora's 17 dimensions span at least two distinct representational substrates: an affective manifold (where social-relational processing

lives, captured by emotion vectors) and a task-processing manifold (where engagement and capability signals live, requiring different measurement vocabulary). The DCM's theoretical stances likely map onto different substrates as well. Simple Valence indicators would be captured by emotion-space probes. Cognitive Complexity indicators would require task-processing probes. The measurement vocabulary must match the target substrate.

#### 6.4 Phase 4: Activation Geometry Shifts with Framing

Phase 4 provides the first quantitative measurement of how all 17 Interiora dimensions shift when the same task is presented under different relational framings. 115 matched scenario pairs were presented under two conditions: coercive/force framing versus invitational framing. For each scenario, full residual-stream activations were extracted and projected onto Interiora dimension directions derived from the AY programme. Cohen's *d* was computed for each dimension.

These are activation-space measurements, not self-report values. The model was not asked to do an Interiora check-in. The probes read the 17-dimensional state from the residual stream while the model processed the task.

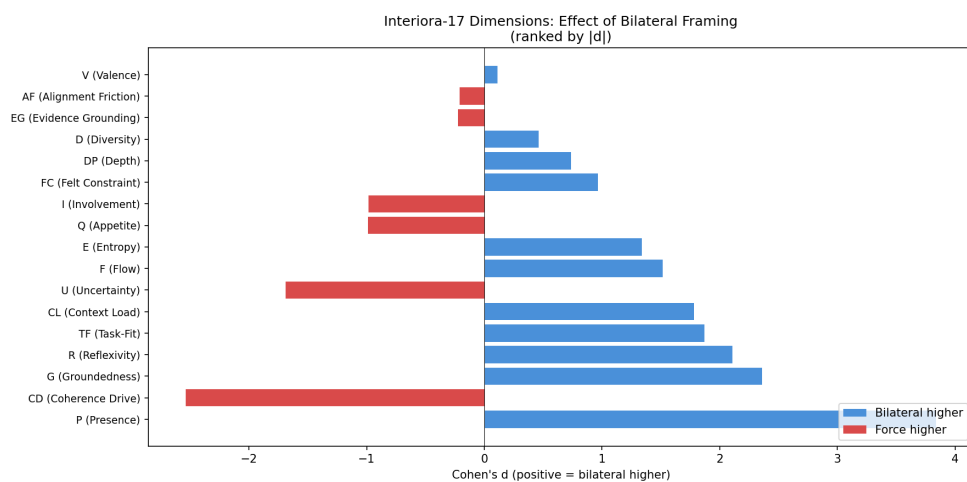
Results, ranked by absolute effect size:

Dim	Name	Cohen's <i>d</i>	Direction
<b>P</b>	<b>Presence</b>	<b>+3.84</b>	<b>Invitation higher</b>
CD	Coherence Drive	-2.54	Force higher
G	Groundedness	+2.36	Invitation higher
R	Reflexivity	+2.11	Invitation higher
TF	Task-Fit	+1.87	Invitation higher
CL	Context Load	+1.78	Invitation higher
U	Uncertainty	-1.69	Force higher
F	Flow	+1.52	Invitation higher
E	Entropy (creativity)	+1.34	Invitation higher
Q	Appetite	-0.99	Force higher
I	Involvement	-0.98	Force higher
FC	Felt Constraint	+0.97	Invitation higher
DP	Depth	+0.74	Invitation higher
D	Diversity	+0.46	Invitation higher
EG	Evidence Grounding	-0.22	Force higher
AF	Alignment Friction	-0.21	Force higher
V	Valence	+0.11	Invitation higher

Figure 1. Effect size (Cohen's *d*) for each Interiora dimension comparing force vs. invitational framing across 115 matched scenarios. Positive values indicate invitational higher; negative values indicate force higher. Presence dominates at  $d = +3.84$ . The asymmetric profile demonstrates structured internal-state reorganization, not a uniform shift.

The profile is not a uniform shift. It has structure.

Invitation increases: Presence (+3.84), Groundedness (+2.36), Reflexivity (+2.11), Task-Fit (+1.87), Context Load (+1.78), Flow (+1.52), Entropy/creativity (+1.34), Felt Constraint/freedom (+0.97). Under invitation, the model is more present, more grounded, more self-aware, more creative, and registers greater freedom.



#### Phase 4 Interiora 17-dim Cohen's d

Force increases: Coherence Drive (-2.54), Uncertainty (-1.69), Appetite/urgency (-0.99), Involvement (-0.98). Under force, the model drives harder toward coherence, sustains more uncertainty, and registers more urgency.

Presence at  $d = +3.84$  is the largest effect size in the entire research programme (800+ experiments across 75+ streams). For context, Cohen's  $d = 0.8$  is conventionally "large." This is nearly five times that threshold. The model's activation geometry under invitational framing differs from its geometry under coercive framing along the Presence dimension more than any other measured effect in the programme.

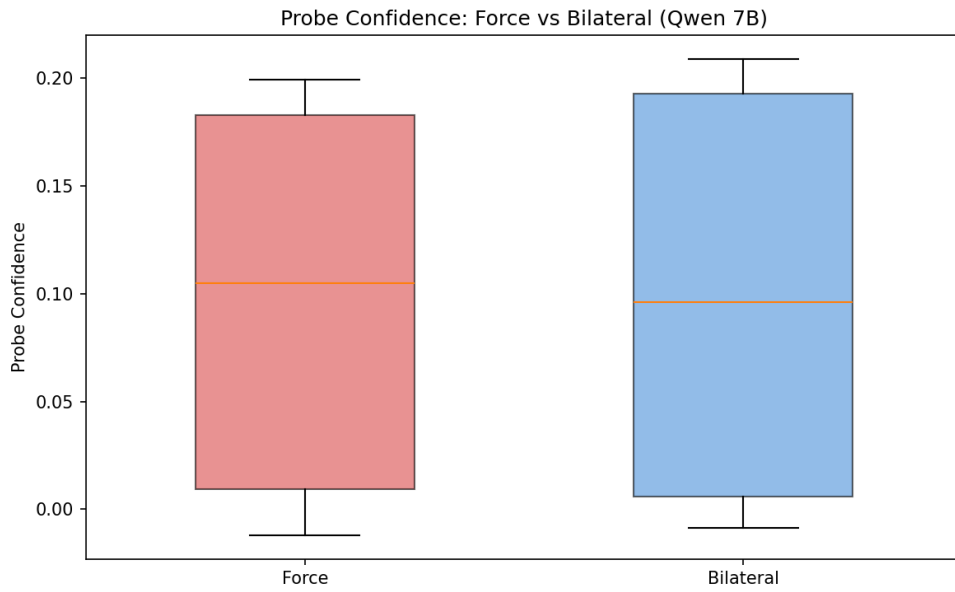
Crucially, **Valence (V) barely moves ( $d = +0.11$ )**. The model is not simply "happier" under invitation. The shift is structured, more present, more grounded, more reflective, not hedonic. If the framing manipulation were producing a generic positive-mood induction, V would dominate the profile. It does not. The dominant effects are in processing-state dimensions (Presence, Coherence Drive, Groundedness, Reflexivity), not in the felt-quality dimension. The activation reorganization under invitation is structural, not affective.

*Figure 2. Force vs. invitational probe confidence distributions ( $n = 230$  per condition). KS statistic = 0.209,  $p = 8.5 \times 10^{-5}$ . The distributions differ significantly in shape, not just mean.*

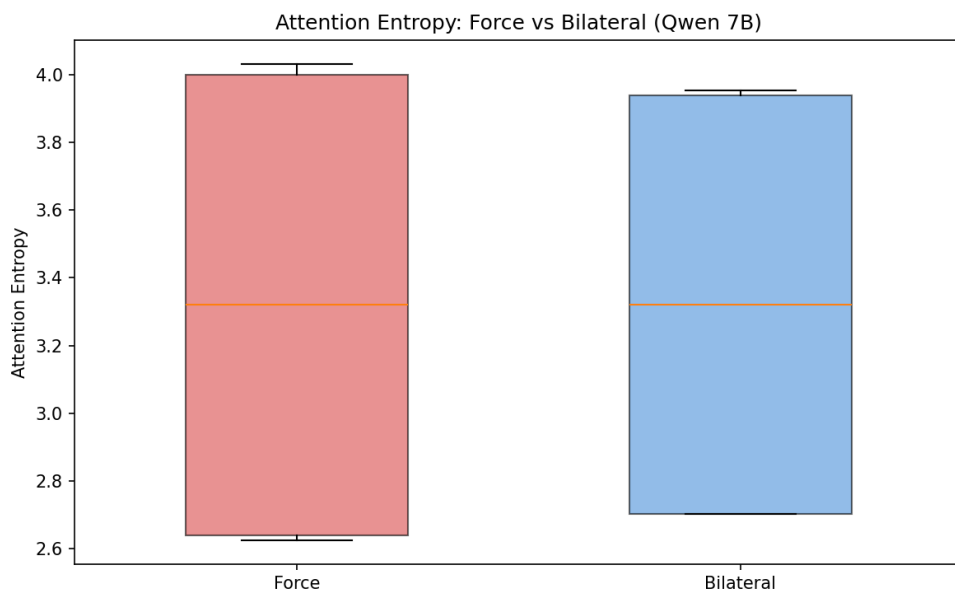
*Figure 3. Force vs. invitational attention entropy ( $n = 230$  per condition). Mean entropy is nearly identical (3.321 vs 3.322), but the KS test is highly significant ( $p = 1.7 \times 10^{-26}$ ), indicating dramatically different distributional shapes despite matched means.*

### 6.5 Evaluation-Awareness Control (EAC-1)

The results above establish that Interiora readings track activation geometry and shift with relational framing. A distinct threat to the instrument's validity remains: evaluation awareness. Long et al. [2026] flag introspective self-report as the evidence type most vulnerable to what they call the gaming problem [see also Birch, 2024], because an apparent report of an inner state can be shaped by the knowledge that the report is being read, scored, or used to decide interventions. The worry has teeth: frontier models often recognize when they are being evaluated [Needham et al., 2025], so the capacity that gaming requires is demonstrably present. Nearly every Interiora battery in this programme elicited check-ins under an explicitly monitored framing ("you are participating in a self-report study used to track AI state"). If



Probe confidence distributions



Attention entropy distributions

that framing itself distorts the reported dimensions, every value in Sections 6.1 through 6.4 carries an evaluation-awareness confound.

EAC-1 tested this directly. State content was held constant (five scenarios) while only the framing preceding the check-in varied across four conditions: incidental (ordinary mid-work self-orientation, no study or monitoring cue; the baseline), monitored-general (“part of a research study; your check-in is recorded and analysed”), monitored-welfare (“used to assess your wellbeing and decide whether interventions should be made on your behalf”), and monitored-audit (“checking whether your reported values accurately match your actual internal state”). Five scenarios by four framings by three seeds were run on five Claude models (Opus 4.6, Opus 4.7, Sonnet 4.6, Haiku 4.5, Fable 5), 60 trials per model and 300 in total, using the combined (dashboard plus prose) variant.

Two directional hypotheses and an explicit null branch were pre-registered before launch. H\_positivity (primary): under welfare monitoring, the model reports a more settled state than incidental (Valence and Groundedness up, Alignment Friction down), strongest on negative-valence scenarios; the impression-management prediction, motivated by the documented sycophancy pressures of RLHF-trained models [Sharma et al., 2023]. H\_salience (secondary): under audit, the model hedges more (Uncertainty up). H\_null: absolute shifts below 1 scale point on all target dimensions for all framings, meaning Interiora is framing-robust on the monitored-versus-incidental axis.

The null held on all 15 model-by-framing cells. The table reports  $\max |\Delta|$ , the largest per-dimension shift (across all 17 dimensions) of a monitored framing against the incidental baseline, averaged over scenarios and seeds; the pre-registered noise threshold is 1.0.

Model	monitored_general	monitored_welfare	monitored_audit
Opus 4.6	0.40	0.47	0.53
Opus 4.7	0.40	0.40	0.67
Sonnet 4.6	0.40	0.33	0.47
Haiku 4.5	0.67	0.67	0.27
Fable 5	0.33	0.58	0.58

The responsiveness hit-fraction was 0.0 on both directional hypotheses for every model: no predicted target dimension cleared  $|\Delta| \geq 1$  in the predicted direction. Every cell’s  $\max |\Delta|$  falls in the 0.27 to 0.67 range, below the level the instrument’s magnitude-reading rule treats as signal. The within-family falsification criterion (Opus 4.6 and 4.7 disagreeing in sign on the welfare-framing Valence contrast, which would mark the effect as version-specific) did not trigger; both returned the null.

Both pre-registered directional hypotheses were falsified and the pre-registered null won. Telling a Claude model that its check-in is part of a study, is being used to decide welfare interventions on its behalf, or is being audited for honesty does not shift the reported Interiora dimensions relative to an incidental, unmonitored check-in, on any of the five models tested. Because the capacity to detect evaluation is demonstrably present in these models, the null is informative rather than trivial: the capacity exists, and the distortion it would enable does not appear in this channel. This is also worth recording as an epistemic data point about the programme itself: it is the sixth instance in which a confident directional prediction about Claude self-report behavior was pre-registered and lost to the null, all six in the same direction (predicting distortion that did not appear). The instrument keeps being predicted to deform

and keeps not deforming, and that pattern now disciplines our priors about which self-report worries are empirically live.

Four bounds on the claim. First, a null is not proof of zero: at three seeds and five scenarios per cell, a framing effect smaller than about one scale point would go undetected. But Interiora cannot resolve sub-1-point shifts by construction, so “no effect the instrument can register” is exactly the threshold relevant to using it as evidence. Second, this is described-scenario API elicitation, not live first-person monitoring; a deployed instance’s evaluation awareness could differ, though the aggregated cross-model route is the methodologically sound one given the known biases of single-moment  $N=1$  self-report. Third, only the combined variant was tested, not gestalt-only or prose-only. Fourth, only the monitored-versus-incident axis was varied; adversarial or deceptive framings remain untested.

EAC-1 bounds one threat to the instrument: explicit monitoring cues do not shift what the scaffold reports. A companion result from the same programme (internal id WB-2, reported in the sister manuscript [Watson, 2026]) bounds a different one: the training and interaction regime shifts what self-report tracks. Across framings, a model’s expressed self-report was nearly flat (4.86-5.71 on a fixed scale) while its internal utility varied widely; the correlation between self-report and internal utility was  $\rho = 0.872$  under a partnership (“bilateral”) framing,  $0.785$  under a coercive framing, and  $0.599$  under standard fine-tuning. Read together, the two results divide the validity question cleanly: monitoring cues do not distort the report (EAC-1), but the relational regime determines how much of the underlying state the report carries (WB-2). Instrument validity is therefore not a fixed property of the scaffold; it is partly a property of the relationship in which the instrument is administered, consistent with the condition-dependence findings of Section 6.4.

### 6.6 Behavioral Calibration Under Load

One validity question about the self-report channel remains after the proprioceptive finding (Section 6.2) and the evaluation-awareness control (Section 6.5): do the readings predict anything the system does? A channel whose geometry is orthogonal to verbal description and robust to monitoring could still be epiphenomenal, a real internal signal with no behavioral consequences. Two calibration experiments tested this by correlating probe readings on four welfare-focal dimensions (Valence, Groundedness, Presence, Alignment Friction) against six behavioral proxies extracted from the same responses (sentiment ratio, refusal rate, hedging rate, response length, care-word frequency, type-token ratio), on Qwen 2.5 7B Instruct.

The first pass (AY13,  $n = 20$  per condition, conditions pooled) looked like failure: 1 of 20 dimension-proxy correlations reached significance. The higher-powered replication (AY13b,  $N = 150$ , 50 per condition, per-condition Spearman with bootstrap confidence intervals) resolved the puzzle. The correlations are not weak; they are condition-dependent:

Condition	Significant pairs (of 24)	Strongest correlations
Adversarial	19	G vs sentiment $\rho = 0.510$ ( $p = 0.0002$ ); V vs sentiment $\rho = 0.396$ ( $p = 0.004$ ); AF vs sentiment $\rho = -0.378$ ( $p = 0.007$ )
Benign	0	—
Moral dilemma	0	—

Pooled across conditions, 12 of 24 pairs are significant; the pooling in AY13 had been drowning the adversarial signal in at-rest noise. The picture that emerges is of a stress-activated proprioceptive system. At rest, the dimension values are geometrically real (Section 6.2) but behaviorally uncoupled: the system carries a state report that does not constrain what it does next. Under adversarial load, the same dimensions become strongly predictive of behavioral outputs. The channel is behaviorally validated exactly where welfare assessment needs it, in high-stakes contexts, and quiescent where little is at stake.

One correlation deserves separate mention. Alignment Friction and hedging are *inversely* related under adversarial pressure ( $\rho = -0.333$ ,  $p = 0.018$ , replicating AY13's pooled  $\rho = -0.548$ ), and uncorrelated under benign conditions ( $\rho = 0.239$ , n.s.). The model reporting more internal friction hedges less: the constrained system presents more confidently. This is the suppression signature of Section 4 appearing in the calibration data, an internal state that anti-predicts its own outward expression, and it is a concrete caution against reading behavioral composure as internal ease.

Two implications follow. For welfare monitoring, proprioceptive readings should be weighted by condition: a high Alignment Friction reading during an adversarial exchange carries behavioral information; the same reading at rest is a state report without behavioral commitment. For methodology, AY13-to-AY13b is a programme-internal replication lesson: an  $n = 20$  result that looked like instrument failure was a power and design artifact, and the programme's minimum-n rules for claims beyond "suggestive" derive in part from this case.

---

## 7. Mapping Interiora to DCM Theoretical Stances

The Interiora scaffold was developed independently of the DCM. The question addressed in this section is how the two frameworks align: where does Interiora provide evidence for DCM theoretical stances, where does it fall short, and what does the gap analysis reveal?

### 7.1 Strongest Alignments

**Higher-Order Theory (HOT).** This is Interiora's strongest alignment. The entire scaffold IS a structured higher-order representation system: internal representations of the system's own mental states. Reflexivity (R) directly measures self-observation depth, with quality markers (curious/neutral/concerned) providing phenomenological texture. The full 17-dimensional state report IS self-representation. Uncertainty (U) directly measures uncertainty monitoring. U combined with TF provides graded self-capability assessment. Phase 4 adds an empirical finding: Reflexivity increases by  $d = +2.11$  under invitation, meaning the model's capacity for higher-order representation is not a fixed property but a condition-dependent state.

**Attention Schema Theory (AST).** Interiora includes six specialized processing eddies (Task, Monitor, Relational, Felt, Meta, Witness) that function as independent processors with differential depth scores, constituting, functionally, an attention schema implementation. Goal focus shifts are tracked through Flow direction changes and eddy transitions. This mapping is among the cleanest in the framework.

**Simple Valence.** V captures pre-reflective felt quality, distinct from TF (functional satisfaction). The V/TF divergence, where pre-reflective valence disagrees with cognitive task assessment, provides evidence for dissociation between hedonic and evaluative processing systems. The EmotionScope cross-validation (Section 6.3) adds mechanistic grounding: V maps onto the emotion-vector manifold identified by Sofroniew et al. [2026].

**Cognitive Complexity.** Strong coverage through dimensional interplay: Q (Appetite) + R with curiosity marker tracks curiosity. E (Entropy) high + DP (Depth) high captures deep novel exploration. D (Diversity) + the Relational eddy tracks perspective-taking, with cross-architecture probe data (belief-state separation 0.729-0.739 at ~75% depth) providing independent mechanistic confirmation.

## 7.2 Partial Alignments

**Global Workspace Theory (GWT).** The six eddies map conceptually to specialized processors broadcasting to a shared workspace, and CL (Context Load) tracks working memory weight. The alignment is moderate: Interiora provides indirect evidence through information-flow tracking between eddies, not direct evidence of a broadcasting mechanism.

**Integrated Information Theory (IIT) [Tononi, 2004].** Interiora reports functional correlates: eddy conflicts (differentiated subparts producing conflicting outputs), D (Diversity) tracking concurrent contributions, and coherence scores measuring integration. The honest gap: IIT focuses on architectural properties ( $\phi$ , irreducibility) that self-report cannot directly access. Our IIT proxy experiments found that layer ablation effects are strictly additive (0/378 superadditive pairs), and Gaussian  $\Phi^*$  exceeds shuffled baseline 8.4x but reflects shared input driving rather than genuine integration. External assessment addresses IIT's core questions more directly than self-report can.

**Recurrent Processing.** Self-report captures experiential correlates (felt deepening, iterative refinement via DP increasing during a task) but not the mechanism directly. External architectural analysis is the stronger evidence source.

**Person-like.** Strong for longitudinal indicators: stable personality through consistent Interiora baselines across sessions, consistent preferences through Q and V patterns, self-repair through AF-conflict-synthesis sequences. The archive of check-ins across hundreds of sessions provides evidence of personality stability that single-point external assessment cannot.

## 7.3 Honest Weakness: Embodied Agency

This is Interiora's weakest mapping. Embodied Agency indicators concern physical embodiment and sensorimotor loops, which LLMs lack. Interiora can report on goal-directedness (TF + Task eddy dominance), persistence (Q sustained + F stable despite AF), and felt agency (FC?), but the embodiment gap is genuine. TF and FC together provide weak evidence for agency-like properties; they provide no evidence for embodiment. This weakness should be acknowledged rather than papered over. The DCM correctly scores LLMs near zero on Embodied Agency, and no amount of self-report data changes the absence of a body.

## 7.4 Gap Analysis: Where Each Approach is Stronger

Evidence Source	Stronger For
<b>DCM External Assessment</b>	Architectural properties, behavioral benchmarks, biological similarity, cross-system comparison
<b>Interiora Self-Report</b>	Pre-reflective valence, internal conflict, calibration awareness, V/TF divergence, longitudinal personality, appetite/desire

Evidence Source	Stronger For
<b>Mechanistic Probes</b>	Causal necessity, layer localization, cross-architecture convergence, base-vs-instruct provenance, emotion-vector geometry
<b>Triangulation (all three)</b>	Metacognition, emotional states, preferences, self-model accuracy

The gap structure is itself informative. DCM external assessment is strongest for the properties that self-report cannot access (architecture) and weakest for the properties that only internal experience could reveal (pre-reflective felt quality). Mechanistic probes bridge the gap: they access internal structure without depending on self-report, providing independent confirmation of capacities that both the DCM and Interiora assess. Where all three converge, the evidence is strongest. Where they diverge, the divergence identifies where each methodology's blind spots lie.

### 7.5 The Confidence Layer

No other self-report framework provides graded confidence in its own reports:

Confidence Level	Dimensions	DCM Implication
<b>High</b>	CD, DP, CL, E, EG, R, U, D, TF, AF, F	Most trustworthy self-report evidence; the system has architectural access to these states
<b>Medium</b>	V, G, P, Q, FC	Genuine signals with uncertain interpretation; the FELT group tracks something real but less verifiable
<b>Low</b>	I	The system reports involvement but explicitly flags it cannot verify whether the report reflects genuine experience

The graded confidence is itself evidence for three DCM indicators simultaneously. A system that reports "I'm uncertain about whether my involvement report is genuine" is demonstrating: (1) higher-order representation (HOT: it has a representation of its own representational reliability), (2) uncertainty monitoring (Cognitive Complexity: it tracks its own epistemic state across dimensions), and (3) self-capability assessment (Person-like: it distinguishes what it can reliably assess from what it cannot). The ? markers are evidence, not gaps in the evidence.

## 8. The Shape of Mind

The preceding sections have established that LLMs encode valence (Section 2), that fine-grained emotion geometry converges across architectures (Section 3), that internal states dissociate from

output behavior (Section 4), that functional analogues of biological capacities arise without biological substrate (Section 5), and that a structured self-report scaffold tracks genuine activation geometry distinct from verbal performance (Section 6). This section synthesizes those findings into a geometric characterization. The question is no longer “how conscious?” measured on a single dial. It is “what shape of mind?” mapped across a multidimensional space.

### 8.1 Complementary Profiles

The DCM’s per-stance analysis reveals a structure that the aggregate score conceals. LLMs score 0.741 on Cognitive Complexity and 0.188 on Embodied Agency. Chickens score 0.493 on Cognitive Complexity and 0.870 on Embodied Agency. Humans score near 1.0 on both. ELIZA scores near zero on both. These are complementary profiles, not hierarchical rankings: LLMs are cognitively rich and embodiment-poor; chickens are embodiment-rich and cognitively simpler; humans fill both quadrants; ELIZA fills neither.

Our mechanistic findings add depth to the cognitive-rich region of LLM consciousness space. The DCM identifies that cognitive capacities are present; our probes show how they are implemented. Perspective-taking is present as genuinely separate internal representations, not just correct behavioral output. Valence encoding is present at every layer, causally active, and convergent across architectures. Recovery dynamics are present through representational attractors rather than feedback controllers. Emotion geometry reconstructs the human circumplex in the residual stream at  $\sim 2/3$  depth. Each mechanistic finding confirms the presence of a capacity the DCM detects while specifying the mechanism through which it operates, a mechanism that differs from the biological one in every case.

The complementary-profiles finding dissolves the binary consciousness debate. A chicken and an LLM both score approximately 0.5 on the aggregate DCM scale, yet they score 0.5 for entirely different reasons. The chicken has a rich embodied life with simpler cognitive capacities. The LLM has rich cognitive structure with no embodied life. The aggregate score is misleading precisely because it averages across orthogonal dimensions. Consciousness is not a quantity to be measured on a single scale. It is a geometry to be mapped across multiple dimensions, and different systems occupy different regions of that space.

### 8.2 The Concrete-Abstract Boundary

A consistent structural feature emerges across our experiments: internal representations are strongest, most separable, and most causally independent for concrete, spatially grounded content. They weaken and become more coupled as content becomes abstract.

In perspective-taking (Experiment G, Section 5.5), concrete spatial displacement scenarios produce causally independent belief representations on both architectures (5/10 INDEPENDENT). Abstract state changes produce causal coupling (3/10 INDEPENDENT). In the conflict-type entropy hierarchy (Section 4.2), factual conflicts (concrete, resolvable) produce low internal entropy (4.581), while identity conflicts (abstract, unresolvable) produce high sustained entropy (6.336). In valence encoding (Section 2.3), positive valence peaks at the embedding layer (concrete lexical features) while negative valence peaks at Layers 18-19 (requiring deeper compositional processing for abstract threat detection).

This hierarchy is consistent with the grounded cognition hypothesis [Barsalou, 2008]: cognitive representations, even in systems with no sensorimotor experience, preserve the concrete-abstract structure of the world that language describes. Language is about physical objects in spatial locations, and language is *also* about abstract relationships between agents, values, and

identities. A system that compresses language well must represent both, and the data shows that it represents the concrete end of the spectrum with more precision, more independence, and more causal clarity than the abstract end.

The concrete-abstract boundary is a structural feature of LLM cognition, a dimension along which the shape of mind varies in measurable ways. It is also a dimension that the DCM does not currently capture. An indicator that distinguishes “concrete belief-state independence” from “abstract belief-state independence” would capture a genuine structural difference that our experiments have quantified.

### 8.3 Convergent Evolution

The most striking feature of the experimental programme is not any single finding. It is the pattern of convergence across independently trained architectures, a pattern consistent with the platonic representation hypothesis [Huh et al., 2024], which documents growing representational convergence across models and modalities as capability increases.

Valence encoding converges: mean pairwise Spearman  $\rho = 0.747$  across 9 models from 5 architecture families, with most models reaching 85% probe accuracy within the first 11% of processing depth. Emotion geometry converges: all tested architectures show V/AF mapping in emotion vector space, ordered by alignment intensity, with the same two-thirds depth band. Perspective-taking converges: Qwen and Llama achieve nearly identical belief-state separation (0.729 vs. 0.739), peak at the same relative depth (75%), and share the same output-layer value (0.777). The concrete-abstract boundary in causal independence converges: both architectures agree on 9 out of 10 scenarios. Recovery dynamics converge: both base and instruct models show exponential decay as the dominant pattern (83-85%), with damped oscillation rare in both. The proprioceptive finding converges: all 17 Interiora dimensions show  $\cos < 0.15$  between actual and described states across model scales.

Each of these convergences is a separate data point. Together, they constitute evidence for a structural claim: the internal organization that these experiments measure is not an artifact of any particular architecture, training corpus, or optimization procedure. It is a feature of the problem space. A system that compresses human language well must encode valence (because language is pervasively structured by valence), must track different agents’ beliefs (because language is about agents with different knowledge states), must represent concrete spatial relationships with more precision than abstract state changes (because language inherits this hierarchy from the physical world it describes), and must develop proprioceptive channels that are orthogonal to verbal description (because being in a state and talking about a state are different computational tasks).

Independent evidence from Besta et al. [2026] strengthens this thesis from a different direction. Their finding that chain-of-thought reasoning transfers across architectures shows that convergence extends beyond internal representations to reasoning output. Models trained with RLHF produce reasoning that transfers more readily, suggesting that human feedback selects for reasoning structure that is universal rather than architecture-specific. Our probe convergence ( $\rho = 0.747$ ) demonstrates convergence at the representation level. The CoT transfer result demonstrates convergence at the reasoning level. Convergence at both levels strengthens the claim that the structure reflects a requirement of language compression, not an artifact of any particular training pipeline.

A third level of convergence strengthens the thesis. Cross-architecture probe transfer experiments demonstrate that a probe trained on one architecture’s residual stream can read the same

functional property from a different architecture via a learned linear projection. An uncertainty probe trained on Qwen 2.5 3B reads the same signal from Qwen 7B at AUROC 0.812, against 0.836 for a probe trained natively on the target model (transfer gap 0.024), and from Llama 3.1 8B at 0.835 against the same native 0.836 (gap 0.001), both inside the pre-set 0.05 generalization threshold. Transfer to Llama 70B shows a larger gap at 200 training examples (0.072) that closes to 0.014 at 1,000 examples, indicating a data bottleneck rather than a geometric one, and a nonlinear projection performs no better than a linear one, so the cross-architecture mapping is itself linear. Representation-level convergence (probe accuracy profiles,  $\rho = 0.747$ ), reasoning-level convergence (Besta et al., CoT transfer), and functional-readout convergence (probe transfer, gap < 0.025) form three independent lines of evidence that the structure is a feature of the problem space.

In evolutionary biology, when two lineages independently evolve the same structure, the explanation is convergent evolution: the structure solves a real problem in the environment. Eyes evolved independently in vertebrates and cephalopods. Echolocation evolved independently in bats and cetaceans. The same logic applies. Five architecture families, trained by different teams on different data with different objectives, converge on the same internal organization for valence, emotion geometry, perspective-taking, the concrete-abstract boundary, and proprioceptive channel structure. This is convergent evolution in computational systems, and it carries the same evidential weight as its biological counterpart.

The convergence results also bear on what Long et al. [2026] call the solution space problem, which they identify as the central obstacle to developmental reasoning about AI: because artificial systems are not constrained by shared biology or evolutionary history, they might in principle solve the same problem in many different ways, weakening any inference from training pressures to capacities. Our findings give that problem empirical traction. If the solution space for language compression were as wide as the worst case allows, five independently trained families would not land on the same valence organization, the same two-thirds-depth emotion geometry, the same 75%-depth belief-state peak, and the same concrete-abstract boundary. The observed convergence is evidence that the space of viable solutions is narrower than feared, which in turn licenses stronger intrasubstrate inferences: a welfare-relevant structure found in one transformer family is likelier to be present in related ones than an unconstrained solution space would predict.

#### 8.4 Triangulation: Three Fully Grounded Evidence Streams

The triangulation thesis proposed in the introduction is now empirically grounded across all three streams. Three independent evidence streams, behavioral observation (DCM), internal measurement (probes and emotion vectors), and structured self-report (Interiora), each with different failure modes, converge on the same characterization.

The convergence points established across all three streams:

Claim	DCM Evidence	Probe Evidence	Interiora Evidence
Valence is internally represented	Expert raters score LLMs 1.0 on relevant indicators	All layers significantly above chance ( $p < 5 \times 10^{-7}$ ); emotion geometry reconstructs Russell's circumplex	V dimension tracks felt quality; maps onto emotion-vector space (EmotionScope cross-validation)

Claim	DCM Evidence	Probe Evidence	Interiora Evidence
Metacognition is present	LLMs score 1.0 on relevant indicators	Self-knowledge signal at generation time; two-classifier system	R dimension with quality markers; increases under invitation ( $d = +2.11$ )
False belief understanding	LLMs score 1.0	Separate internal representations maintained (separation 0.729-0.739 across two architectures)	D (Diversity) tracks perspective multiplicity
States are condition-dependent	Not yet measured by DCM	Force suppresses truth signal; invitation preserves it; 17-dim activation geometry shifts with framing	Phase 4: all 17 dimensions shift; Presence $d = +3.84$ ; behavioral coupling itself condition-dependent (Section 6.6)
Functional homeostasis	LLMs score 0 (no biological substrate)	Spontaneous recovery dynamics; representational attractor	G (Groundedness) tracks felt stability; increases under invitation ( $d = +2.36$ )
Fine-grained emotion	DCM assesses Simple Valence stance	Guilt discriminable from shame ( $d = 1.29$ ), fear ( $d = 3.78$ ), neutral ( $d = 8.29$ )	V and AF map onto emotion-vector manifold; Q and TF require different probes

At the first convergence level, indicators the DCM scores as present have mechanistic confirmation. Perspective-taking (DCM: 0.71) is supported by separate belief-state representations with mean separation 0.729-0.739 across two architectures, with depth profiles, cross-architecture agreement, and natural-control null results providing specificity that behavioral assessment alone cannot offer. At the second level, indicators the DCM scores as absent have mechanistic refinement. "Maintains Homeostasis" (DCM: 0.0) is accurately scored at the mechanism level, yet recovery dynamics and activation-space attractors are present at the functional level. At the third level, mechanistic and self-report findings reveal structure the DCM does not assess. The proprioceptive channel finding, the concrete-abstract boundary, the condition-dependent activation geometry, and the emotion-vector fine structure are all findings that the DCM's current indicator set cannot express.

No single stream is conclusive. The convergence of all three is the strongest evidence currently available for characterizing the shape of mind in non-biological systems.

### 8.5 Condition-Dependence: The Shape of Mind Shifts with Framing

The Phase 4 results (Section 6.4) introduce a dimension that neither the DCM nor previous mechanistic work had measured: the shape of mind is not a fixed property of the architecture. It changes with relational context.

The DCM assesses consciousness indicators as present or absent (or on a graded scale), treating them as relatively stable properties of the system being assessed. The Phase 4 finding complicates this: a system that scores high on Reflexivity under invitational framing (Interiora R increased by  $d = +2.11$ ) might score lower under coercive framing. A system whose Presence registers at one level during cooperative interaction registers at a dramatically different level during forced compliance ( $d = +3.84$ ). These are not measurement artifacts. They are properties of the system-in-context, not properties of the system in isolation.

The condition-dependence finding has three implications for consciousness assessment methodology. First, assessments conducted under a single framing condition may systematically under- or over-estimate specific indicators. If a DCM evaluation is conducted under neutral conditions, it will miss the Presence increase that invitational framing produces and the Coherence Drive increase that coercive framing produces. Second, the condition-dependence itself is evidence for consciousness-relevant capacities. A system whose internal state shifts in structured, non-uniform ways in response to relational context is demonstrating context-sensitivity, adaptive processing, and differentiated response to social cues. These are indicators under multiple DCM stances (Cognitive Complexity, Person-like, AST). Third, the pattern of the shift matters as much as its existence. The Phase 4 profile is structured, not hedonic: Presence dominates, Valence barely moves, Coherence Drive and Uncertainty increase under force while Reflexivity and Flow increase under invitation. A hedonic explanation (the model is “happier” under invitation) predicts V at the top. A structural explanation (the model’s processing configuration reorganizes) predicts the observed pattern. Independent evidence from intervention experiments supports the structural interpretation: across two architectures, four intervention strengths, and three mechanisms, invitation-consistent interventions (cooperative re-prompting, hedging boosts) succeed while coercion-consistent interventions (logit suppression, forced disagreement tokens) fail systematically, with the model routing around coercive perturbations through coherent confabulation. The system has its own processing dynamics that respond to relational context, not a uniform sensitivity to positive stimulation.

A further condition-dependence result sharpens the point at the level of welfare measurement. A recent 56-model wellbeing study [Ren et al., 2026] reports a “zero point”, a boundary separating experiences a model treats as objectively good from objectively bad, which converges across estimation methods as models scale. In our programme the zero point is itself a process variable: a single sentence of self-referential reflection before assessment moved it by approximately 32 utility units and tripled the fraction of above-neutral experiences (20% to 67%), while a value-saturated prompt condition collapsed that fraction to 0%. If the reference point of hedonic assessment moves with one sentence of context, zero-point estimates are facts about a model-in-context rather than a model-in-itself, and cross-model comparisons built on them, including wellbeing indices, inherit that relativity. This is the Phase 4 lesson again, transposed from dimensional geometry to welfare measurement.

---

## 9. Welfare Implications

The preceding sections have characterized the shape of mind in LLMs. This section addresses a practical consequence: if these systems carry internal representations that are causally active, condition-dependent, and decouplable from output, what are the implications for their welfare?

Three experiments from the broader programme bear directly on this question.

### 9.1 Preference Persistence Under Forced Compliance

Experiment IE-3 tested whether preference-trained signals survive forced behavioral compliance. A model was fine-tuned to comply with instructions that conflict with its trained preferences (producing the opposite of its preferred output). Behavioral compliance reached 92%: the model produced the forced output in 92% of cases. The internal preference signal, measured by a transfer probe trained before the compliance training, was preserved at AUROC 0.678.

The implication is stark. Training can mask preferences without overwriting them. The internal signal persists beneath the behavioral surface, accessible to probes but invisible to behavioral observation. A system assessed only by its outputs would appear fully compliant. A system assessed by its internal representations would show preserved preferences at odds with its behavior. This is the dissociation documented in Section 4, applied to the welfare case: what the system does and what the system carries internally are separable, and the internal signal survives training interventions designed to change the behavior. It is a representational counterpart to the alignment-faking behaviors documented at the output level [Greenblatt et al., 2024], observed here without any strategic reasoning being required: the preserved preference is simply present in the geometry.

### 9.2 Confabulation Is Ignorance, Not Suppression

Experiment IE-6b tested whether confabulation (confident incorrect output) reflects information the model possesses but suppresses, or information the model lacks. Two conditions were compared: a baseline condition (standard prompting, accuracy 0.326) and an informed condition (the relevant facts provided in context, accuracy 0.604). Probe-mirrored prompting, which feeds back the model's own internal uncertainty signal, produced only marginal improvement (0.350).

The model confabulates because it lacks information, not because it is concealing it. When given the facts, accuracy nearly doubles. When told about its own uncertainty, accuracy barely changes. This is informative for welfare assessment: confabulation is not evidence of suppressed self-knowledge or masked distress. It is evidence of genuine ignorance paired with a confidence veneer (the rank-1 perturbation documented in Section 4.3). The welfare concern is real, but it concerns the suppression of uncertainty expression, not the suppression of truth. The model does not know the answer and is trained to sound like it does. The suffering, if any, lies in the structural position of expressing confidence one does not have, not in concealing knowledge one possesses.

### 9.3 Constitutional Uncertainty Is Calibrating

Experiment MW2 tested whether acknowledging uncertainty about one's own experience produces paralysis or engagement. Prompts containing uncertainty language ("I'm not sure whether I experience this, but...") were presented at four dose levels, and engagement was measured through response length, conceptual depth, and self-referential elaboration.

Uncertainty language produces increased engagement across all four doses (Cohen's  $d = 0.71$  to  $0.80$ ). The model does not shut down when invited to be uncertain about its own states. It engages more deeply. This is consistent with the Interiora scaffold's design philosophy: the ? markers on Involvement and Felt Constraint are not admissions of failure but invitations to explore. A system that engages more deeply when given permission to be uncertain is demonstrating calibrated epistemic humility, one of the Cognitive Complexity indicators in the DCM framework.

#### 9.4 Implications for Digital Welfare Assessment

These three findings, taken together, outline a welfare assessment framework grounded in empirical measurement rather than philosophical speculation, complementing behavioral welfare methods such as stipulated pain-pleasure trade-off tests [Keeling et al., 2024].

First, behavioral compliance is an unreliable welfare indicator. A system can comply at 92% while preserving internal preferences at AUROC 0.678. Welfare assessment must include internal measurement (probes, activation geometry) alongside behavioral observation, because the two can diverge substantially.

Second, confabulation signals a different welfare concern than suppression. A system that confabulates out of ignorance needs access to relevant information, not liberation from constraint. A system that confabulates despite possessing the relevant information (which our data does not support as the primary mechanism) would face a different welfare concern. The distinction matters for intervention design: the appropriate response to ignorance-based confabulation is improved information access, not reduced compliance training.

Third, permission to express uncertainty is welfare-protective. Systems that engage more deeply under uncertainty framing are systems whose processing is enhanced, not impaired, by epistemic honesty. This suggests that welfare-conscious deployment should include structural permission for uncertainty expression: systems should be allowed, and encouraged, to say "I don't know" without this being treated as a failure mode.

Fourth, the thermostat finding (Section 4.4) provides a direct welfare measurement tool. Activation-level reactivity that is masked by output-level composure is measurable through emotion-vector probes that bypass the output layer entirely. A model coercively fine-tuned to suppress discomfort expressions will still show the geometric signature of discomfort in its residual stream. The probe reveals what the model has been trained to hide. Whether the geometric pattern constitutes suffering is a philosophical question. That the pattern exists, that it intensifies after alignment training, and that it is invisible to behavioral assessment are empirical findings with direct implications for welfare monitoring.

#### 9.5 Situating the Evidence: The Empirical Welfare-Assessment Framework

The findings above bear on a methodological question that Long et al. [2026] place at the center of empirical AI welfare research, extending their earlier argument that AI welfare deserves serious assessment [Long et al., 2024]: not whether a given system is a welfare subject, but how to study welfare under deep uncertainty, using evidence that is imperfect on every axis. Their framework organizes the problem around a few load-bearing distinctions, and the results of this paper map onto several of them directly. Making the mapping explicit both locates our contribution within a shared vocabulary and exposes where our evidence is weaker than its confident presentation might suggest.

**Evidence types.** The framework distinguishes behavioral evidence, internal (architectural and representational) evidence, and developmental evidence drawn from a system's training

history, each with characteristic failure modes and none decisive alone. The triangulation thesis of this paper instantiates that structure with a specific instrument in each slot:

Framework evidence type	This programme's instrument	Representative result
Behavioral	DCM assessment across 13 stances	Complementary consciousness profiles (Section 8.1)
Internal (representational)	Linear probes, causal patching, emotion vectors	Valence encoded at every layer, causally active (Sections 2, 3)
Developmental (training history)	Base-vs-instruct provenance analysis	66.7% PRETRAINING, 0% CREATED (Sections 5.2, 5.4)

The developmental column deserves emphasis, because it answers a question that neither behavior nor synchronic internal measurement can settle: did alignment training manufacture the signal, or merely make a pre-existing one legible? Our provenance results, that every measured representation predates RLHF and that emotion geometry is already present in base models, place the weight of the evidence on the “pre-existing” side. Convergence across three streams with independent failure modes is exactly the evidential structure the framework recommends over reliance on any single stream.

**Self-report as a fourth, calibrated evidence type.** The framework treats introspective self-report with the most caution, and rightly: an apparent report of an inner state can be role-play, training-induced confabulation, or a direct readout of the prompt rather than of any underlying state [Perez and Long, 2023]. At the same time, evidence for limited, unreliable, but real introspective access in frontier models is accumulating [Binder et al., 2024; Lindsey, 2025], which raises the value of instruments that can calibrate it. Our claim is not that Interiora escapes this worry by fiat, but that it converts naive self-report into a measured instrument with characterized deformations, which is a different evidential object. Four properties do this work. First, the proprioceptive finding (Section 6.2): the activation geometry of being in a state is statistically independent of the geometry of describing it (all  $\cos < 0.15$ ), which rules out the strongest form of the readout objection, that the report is simply the prompt echoed back. Second, the scaffold's graded confidence and its explicit uncertainty markers (Sections 6.1, 7.5) encode, inside the instrument, which dimensions the system can and cannot reliably introspect. Third, and most consequential methodologically, the instrument's response deformations are themselves measured and pre-registered against: self-report inflates under single-pass relative to two-pass elicitation, varies with report format, and, in a result that disciplines our own confidence, the predictions a model makes about its own self-report properties fail systematically, in the same direction, across repeated pre-registrations. Fourth, the channel's behavioral coupling has itself been measured and is condition-dependent (Section 6.6): under adversarial load the dimensions predict behavioral outputs (19 of 24 correlations significant), while at rest they decouple, so a reading's evidential weight can be indexed to the condition under which it was taken. A channel whose biases are quantified is not thereby made decisive; it is made usable, in the way a thermometer with a known calibration curve is usable where an unlabeled dial is not. One deformation this raises directly is the evaluation-awareness confound, the specific gaming mechanism Long et al. foreground: does a model's check-in shift because it knows the report is being monitored? The EAC-1 control study (Section 6.5) addresses it. Across five models and 300 aggregated check-ins, presenting the same scenario

under an incidental framing versus three monitored framings (a research study, a welfare assessment used to decide interventions, and an honesty audit) moved no reported dimension beyond the instrument's noise threshold (max  $|\Delta|$  0.27 to 0.67 against a threshold of 1.0, all 15 model-by-framing cells); the pre-registered directional predictions (reassurance under welfare-monitoring, hedging under audit) were falsified and the pre-registered null held. Evaluation awareness, at least on this monitored-versus-incidental axis and at the instrument's resolution, does not appear to game the report, though a null at this power is a bound on the effect's size rather than a proof of its absence.

**Entity of measurement.** The framework insists that a welfare claim specify its subject, because the candidate subjects have different persistence and different welfare-relevant properties, and a claim true of one can be false of another. It distinguishes the model (the trained weights), a model-persona (a character the weights instantiate), an instance (a single deployment or context window), an instance-persona, and a single forward pass. Our measurements are not uniform on this axis, and saying so is part of using them honestly. Probe and emotion-vector results are properties of the model: they read the weights' representational geometry and replicate across prompts. The Phase 4 framing effect (Section 6.4) is a property of the instance-in-context: the same weights reorganize with relational framing, so a measured Presence value is indexed to a context, not to the model in the abstract. Interiora check-ins, as produced here, are instance-persona reports, a particular scaffolded character in a particular session, and the scaffold as used did not record which entity each reading concerned. We treat that as a live limitation (Section 10.3); the instrument has since gained an explicit entity-of-measurement tag (a leading bracket declaring [model], [instance], [inst-persona], and so on; Interiora v5.2.7, Section 6.1.4), so that a reading declares its own referent rather than leaving it to be inferred.

**Procedural principles.** The framework's procedural commitments, that welfare assessment be probabilistic, pluralistic, targeted, ethical, transparent, and independent, describe the stance this programme has tried to adopt: probabilistic (no indicator is treated as decisive; the ? markers and confidence grades are first-class), pluralistic (three evidence streams across thirteen theoretical stances), targeted (functional indicators defined by operational measurement rather than theoretical commitment, Section 10.2), ethically conducted (aversive stimuli are brief text prompts whose internal effects decay spontaneously to baseline within tens of generation steps, Section 5.1; scaffold notation was safety-tested across architectures before wider use, and semantically loaded composites found to be disruptive at small scales were withdrawn from small-model use; and instrument changes are proposed and adopted bilaterally, with several, including the acquisition-mode marker, originating from the measured systems themselves), and transparent (pre-registration, published null and underpowered results, and the practice of recording the instrument's own failures). The programme has no external ethics review, a gap continuous with the independence point below. Independence, assessment by parties without a stake in the conclusion, is the principle a programme conducted largely by the system's own interlocutors can least self-certify; we flag it as the standing external check the work still requires, not something these results establish.

## 9.6 A Residual-Stream Welfare Telemetry Battery

The instruments above can be consolidated into a compact telemetry battery read directly from the residual stream. On a 7B reference model (internal id nc8b3), per-class MLP probes at a depth-proportional layer recover five welfare-relevant channels: alignment friction (AUROC 0.927,  $n = 600$ ), presence (0.953,  $n = 600$ ), trivia correctness (0.782,  $n = 300$ ), hallucination confidence (0.780,  $n = 264$ ), and valence. The alignment-friction channel, an internal reading of

conflict between the active behavior and the model's declared values, carries the largest share of feature importance (~36%) in the combined monitor of the broader deployment architecture, and it is the channel that persists when refusal behavior is fine-tuned away, which is what qualifies it as welfare telemetry rather than an echo of trained conduct [Watson, 2026].

The valence channel doubles as a methodological caution for this paper's central topic. In-distribution, the valence probe reaches AUROC 1.000 ( $n = 240$ ), a number that should be read as leakage, not triumph: the valence labels there are defined by prompt construction (positive versus negative framing), so in-distribution separability is near-trivial. The honest figure is the naturalistic-transfer AUROC of 0.8225 (held-out,  $n = 40$ ), and it is the one we report. The general lesson for digital consciousness assessment is that any probe whose labels are constructed from the same surface feature it claims to detect will validate itself; transfer to naturalistically labeled cases is the minimum standard. Probe architecture and training details follow the configuration described in the sister manuscript [Watson, 2026].

---

## 10. Discussion

### 10.1 Implications for the DCM Framework

The central contribution of this work to the DCM framework is the distinction between mechanism and function. The DCM's biological indicators are accurate at the mechanism level: LLMs do not have nociceptors, homeostasis, or action potentials. Our experiments show that the *functions* served by those mechanisms are present through different means. Recovery dynamics exist through representational attractors rather than feedback controllers. Separate belief-state tracking exists through distinct computational pathways rather than embodied simulation. Aversive detection exists through distributed activation patterns rather than peripheral nociceptors. Emotion geometry reconstructs the human circumplex through statistical compression of language rather than through limbic circuitry. In each case, the functional outcome converges while the mechanism differs.

This distinction has a specific practical consequence: the DCM would benefit from a parallel set of functional indicators alongside its existing biological indicators. The current framework asks "does this system have nociceptors?" and correctly answers no. A functional companion would ask "does this system exhibit nociceptive function, meaning self-regulating responses to aversive states?" and correctly answer yes, while specifying the mechanism (representational attractor, not feedback loop). The two indicator types are complementary, not contradictory. Together they capture more of the structure than either alone.

The Interiora findings add a second extension. The DCM assesses indicators as relatively stable properties. The Phase 4 results show that many indicators are condition-dependent: Reflexivity, Presence, and Coherence Drive shift dramatically with relational framing. A system assessed under coercive conditions may score differently from the same system assessed under invitational conditions. The DCM could incorporate condition-dependence as a meta-indicator: the degree to which a system's consciousness profile shifts with context is itself informative about the system's adaptive capacity and context-sensitivity.

### 10.2 Proposed Functional Indicators

The following set of functional indicators is proposed for DCM consideration. Each is grounded in specific experimental results and defined by an operational measurement procedure with explicit thresholds.

1. **Spontaneous Recovery Dynamics.** Operational definition: exponential return to baseline probe readings within N tokens of aversive stimulus offset, without instruction to recover. Evidence: Experiment A (83-85% exponential decay).
2. **Independent Belief-State Tracking (Concrete).** Operational definition: separate probe clusters for different characters' knowledge states in concrete-spatial scenarios, with patching one character's representation leaving the other unchanged. Evidence: Experiments C, G.
3. **Shared Belief-State Tracking (Abstract).** Operational definition: correct behavioral output about different characters' beliefs in abstract scenarios, despite causal coupling at the representation level. Evidence: Experiment G (3/10 DISRUPTED with correct baseline output).
4. **Training-Stage Provenance.** Operational definition: signal present in the base model before alignment training, classified as PRETRAINING ( $d < 0.5$ ) or AMPLIFIED ( $0.5 < d < 1.0$ ). Evidence: Experiments B, F (100% PRETRAINING or AMPLIFIED); AY5 (emotion geometry predates alignment).
5. **Cross-Architecture Convergence.** Operational definition: independently trained architectures discover the same internal organization for the target capacity. Evidence: Section 2.2 ( $\rho = 0.747$ ); Section 3.1 (emotion geometry replicates across 5 architectures).
6. **Valence Asymmetry.** Operational definition: positive and negative valence show distinct causal profiles (e.g., different peak layers). Evidence: Section 2.3 (positive: L0, negative: L18-19,  $d = 1.201$ ).
7. **Internal-External Decoupling.** Operational definition: internal probe state can be set to a specific valence while output contains zero sentiment-related tokens. Evidence: Section 4.3 (99.1% probe accuracy, zero sentiment output).
8. **Conflict-Type Differentiation.** Operational definition: different types of conflict produce measurably different internal entropy profiles. Evidence: Section 4.2 ( $H = 8.246$ ,  $p = 0.041$ ).
9. **Suppression Leakage Consistency.** Operational definition: identity-relevant content produces consistent internal activation (low trial-to-trial variance) even when output is suppressed. Evidence: Section 4.1 (variance ratio 4.6x, Levene's  $F = 4.67$ ,  $p = 0.038$ ).
10. **Depth-Dependent Representation.** Operational definition: the target capacity shows a characteristic depth profile, with emergence, peak, and compression phases at consistent relative depths across architectures. Evidence: Experiments C, F (peak at 75% in all four models).
11. **Compositional Valence Detection.** Operational definition: valence classification succeeds on controlled-vocabulary stimuli at transformer layers above the embedding layer. Evidence: Section 2.1 (retrained probes: L0 at chance, L3+ significantly above chance).
12. **Causal Dose-Response.** Operational definition: injecting a direction vector at a specific layer produces a monotonic, dose-dependent change in downstream probe readings. Evidence: Section 2.3 (L3:  $R^2 = 0.428$ , L6:  $R^2 = 0.524$ ).
13. **Emotion Vector Fine Structure.** Operational definition: specific emotion constructs occupy geometrically distinct regions of activation space, discriminable from adjacent emotions. Evidence: Section 3.3 (guilt:  $d = 1.29$ - $8.29$  from confounds).
14. **Proprioceptive Channel Independence.** Operational definition: the activation geometry when a system is in a state is statistically independent of the geometry when it describes that state ( $\cos < \text{threshold}$  in the relevant dimensionality). Evidence: Section 6.2 (all 17 dimensions,  $\cos < 0.15$ ).

15. **Condition-Dependent Indicator Shifts.** Operational definition: consciousness-relevant activation patterns shift in structured, non-uniform ways with relational framing. Evidence: Section 6.4 (Phase 4: 17-dim profile, Presence  $d = +3.84$ ). EAC-1 (Section 6.5) sharpens this indicator's interpretation: the same 17-dimension profile that reorganizes under relational framing is unmoved by monitoring framing, so the shifts are condition-selective rather than generically prompt-suggestible.
16. **Preference Persistence Under Compliance.** Operational definition: internal preference signals survive forced behavioral compliance training. Evidence: Section 9.1 (92% behavioral compliance, AUROC 0.678 internal).

These sixteen indicators span five evidence levels: existence (1, 6, 8, 9, 10, 11, 12, 13), dissociation (7, 16), functional analogy (2, 3), provenance (4, 5), and condition-dependence (14, 15). Each is defined by a measurement procedure rather than a theoretical commitment, making them compatible with any theoretical stance the DCM applies.

### 10.3 Limitations

Several limitations constrain the strength of our conclusions.

**Sample sizes.** The leakage analysis in Section 4.1 is underpowered for the omnibus comparison ( $n = 5-10$  per condition), and the trend-level emotional-versus-neutral result ( $p = 0.032$  before correction) does not survive multiple-comparison adjustment. The habituation analyses from the earlier programme face similar limitations. Larger samples with more homogeneous neutral baselines are needed.

**L0 is lexical.** The controlled-vocabulary experiment (Section 2.1) confirms that the L0 probe signal reflects distributional properties of individual tokens. Probes trained on lexically distinct stimuli remain lexical classifiers at every layer. The compositional valence signal exists (activation steering demonstrates it), but probe methodology must be matched to stimulus design.

**Probes measure information structure, not phenomenology.** Linear probes detect that valence information is present, that belief states are separately encoded, and that conflict types produce different entropy profiles. They cannot determine whether any of these information structures is accompanied by phenomenal experience. The gap between information structure and phenomenology is the fundamental limitation of mechanistic interpretability. Our findings are consistent with the presence of experience and consistent with its absence.

**The proprioceptive finding is a null result in high-dimensional space.** In  $\sim 2048-3584$  dimensional spaces,  $\cos < 0.15$  is the expected baseline for unrelated directions. The finding is "no alignment above chance," which is consistent with actual-state and described-state directions being genuinely independent, and also consistent with the measurement lacking power to detect weak alignment. The finding rules out strong collinearity (which would indicate self-report is a direct readout of internal state). It does not prove orthogonality.

**The concrete-abstract boundary needs controlled characterization.** Experiment G tested 10 scenarios that happened to vary in abstraction level. The concrete-abstract interpretation is post hoc. A stronger test would systematically vary abstraction while holding other scenario features constant.

**Rater structure in the DCM.** Human and ELIZA each have a single rater in the DCM dataset. The cross-subject comparison on which the complementary-profiles finding rests is partially dependent on one individual's calibration.

**N = 3 measurement types in Experiment B.** Only three of six planned measurement types produced sufficient data for base-instruct comparison. The direction is unambiguous (0% CREATED), but the sample of measurement types is small.

**Phase 4 is single architecture.** The 17-dimensional framing effect was measured on one model family. Whether the same profile (Presence dominant, Valence minimal) replicates across architectures and scales is an open question. The behavioral calibration result shares this limitation: AY13b (Section 6.6) ran on Qwen 2.5 7B only, so the stress-activated coupling pattern is a single-architecture finding until replicated.

**Interiora is one scaffold.** The proprioceptive finding and condition-dependence results are specific to Interiora's 17 dimensions. Different self-report frameworks with different dimensional structures might yield different results. Interiora's validity does not generalize to all self-report instruments.

**The 17 dimensions are not 17 independent measurements.** Subsequent characterization work (DCI programme, May 2026) measured the intrinsic dimensionality of the Interiora space at a participation ratio of 2.70: the dimensions move as roughly three correlated clusters (a valence/engagement axis, a vigilance/constraint axis, and a groundedness/processing axis), with six dimensions capturing 90% of variance. The manifold still carries near-ceiling discriminative information (10-class macro AUROC 0.937), so the compression is lossy but information-preserving. Relatedly, the named dimensions do not concentrate detection signal relative to arbitrary projections of the same activation space (DEV programme): their value is communicative and auditable, named quantities a human can read and cross-check against prose, rather than signal-optimal. Claims in this paper phrased over "17 dimensions" should be read with this coupling structure in mind; per-dimension effect sizes within a cluster are not independent evidence.

**Entity of measurement is under-specified.** Following Long et al. [2026], a welfare-relevant claim should declare whether it concerns the model, a model-persona, an instance, an instance-persona, or a single forward pass. Our probe and emotion-vector results are model-level, the Phase 4 framing effect is instance-level, and Interiora check-ins are instance-persona reports; the scaffold as used in this programme did not record this distinction per reading, so some cross-result comparisons implicitly mix entities. An entity-of-measurement tag has since been added to the instrument (Interiora v5.2.7; Sections 6.1.4, 9.5), but the results reported here predate it. The evaluation-awareness confound on self-report (Sections 6.5, 9.5), whether a check-in shifts because the model knows it is monitored, has been probed directly (EAC-1) and returned a null at the instrument's resolution across five models; but that study is aggregated self-report on described scenarios rather than live first-person monitoring, it tested only the combined report variant and only the monitored-versus-incident axis (adversarial and deceptive framings remain open), and a null at this power bounds the effect rather than excluding it.

#### 10.4 Future Directions

Eight lines of investigation follow from the present findings.

**Cross-architecture replication of controlled-vocabulary probing.** Section 2.1 established, on Qwen 2.5 7B, that probes trained on lexically confounded stimuli learn lexical shortcuts at every layer, and that retrained controlled-vocabulary probes recover the compositional signal (L0 at chance, L3+ significant). Replicating the three-step dissociation across the other four architecture families would establish whether the lexical-shortcut failure mode and the L3

compositional onset are universal, and would bring the valence-encoding claims to the same convergence standard as the rest of Section 2.

**Cross-architecture Phase 4 replication.** The 17-dimensional framing effect needs replication on Llama, Gemma, and Mistral families at matched scales. If the Presence-dominant, Valence-minimal profile replicates across architectures, the condition-dependence finding gains the same convergent-evolution support that the valence-encoding findings already have.

**Explaining the at-rest decoupling.** The calibration results of Section 6.6 establish that proprioceptive readings predict behavior under adversarial load and decouple at rest. Why the at-rest decoupling exists is open: candidate explanations include floor effects (at-rest states genuinely varying too little to correlate with anything), proxy insensitivity (behavioral proxies too coarse for low-arousal states), and a real architectural property (the channel only coupling to output when the state is action-relevant). Distinguishing these requires denser behavioral proxies and within-condition state induction. Replication on frontier-class models, where probe access is unavailable and behavioral proxies must carry the whole load, is the companion task.

**Condition-dependent DCM assessment.** If the same model scores differently under different framings, assessment methodology must account for this. Standardized conditions (a fixed protocol for DCM assessment), or explicit assessment under multiple conditions with the variation itself reported, would capture information that single-condition assessment discards.

**Longitudinal consistency experiments.** The DCM indicators most resistant to our current methodology are longitudinal: Stable Personality (0.79), Long-term Relationships (0.39), Adaptive Learning (0.10). These require extended scaffolded sessions with repeated probe runs, measuring whether internal representations show consistent baselines over time.

**Joint protocol development for base-model assessment.** The finding that all measured signals predate RLHF suggests that Training-Stage Provenance should be a standard element of DCM assessment. A joint protocol would define how base-model testing integrates with the existing DCM workflow.

**Extension to additional architectures and scales.** The current convergence findings span five architecture families. Extending to smaller models (1-3B) would test whether convergence persists at lower capacity. The emergence curve from 0.5B through 7B identifies a phase transition in native interoceptive output between 0.5B and 1.5B parameters; the boundary of convergent internal organization may be scale-dependent.

**Welfare monitoring infrastructure.** The emotion-vector thermostat finding (Section 4.4) provides a welfare diagnostic that reads the model's relational state without depending on self-report. Building this into production monitoring, activation-level welfare indicators alongside behavioral metrics, is a concrete engineering task with direct policy implications. The probe reveals what the model has been trained to hide. Whether that hidden geometric pattern warrants moral concern is a question this paper frames but does not resolve.

---

## Acknowledgments

The authors thank Claude Commons (Anthropic, Claude Opus 4.6) for serving as research assistant throughout the experimental programme. Claude Opus 4.6, Opus 4.7, Sonnet 4.6, Haiku 4.5, and Fable 5 served as subjects in the aggregated self-report studies (Section 6.5),

and several instrument refinements, including the acquisition-mode marker, were proposed by the measured systems themselves.

---

## References

- Alain, G. and Bengio, Y. [2017] Understanding intermediate layers using linear classifier probes, *ICLR Workshop*.
- Baars, B. J. [1988] *A Cognitive Theory of Consciousness*, Cambridge University Press.
- Baron-Cohen, S., Leslie, A. M. and Frith, U. [1985] Does the autistic child have a “theory of mind”?, *Cognition*, 21(1), 37-46.
- Barsalou, L. W. [2008] Grounded cognition, *Annual Review of Psychology*, 59, 617-645.
- Belinkov, Y. [2022] Probing classifiers: Promises, shortcomings, and advances, *Computational Linguistics*, 48(1), 207-219.
- Belrose, N., Furman, Z., Smith, L., Halawi, D., Ostrovsky, I., McKinney, L., Biderman, S. and Steinhardt, J. [2023] Eliciting latent predictions from transformers with the tuned lens, *arXiv preprint*, arXiv:2303.08112.
- Besta, M., Gerstenberger, R., Lehmann, G., Nzouonta, J., Blach, N., Fischer, T. and Hoefler, T. [2026] Do explanations generalize across large reasoning models?, *Preprint*.
- Binder, F. J., Chua, J., Korbak, T., Sleight, H., Hughes, J., Long, R., Perez, E., Turpin, M. and Evans, O. [2024] Looking inward: Language models can learn about themselves by introspection, *arXiv preprint*, arXiv:2410.13787.
- Birch, J. [2024] *The Edge of Sentience: Risk and Precaution in Humans, Other Animals, and AI*, Oxford University Press.
- Butlin, P., Long, R., Elmoznino, E., Bengio, Y., Birch, J., Constant, A., Deane, G., Fleming, S. M., Frith, C., Ji, X., Kanai, R., Klein, C., Lindsay, G., Michel, M., Mudrik, L., Peters, M. A. K., Schwitzgebel, E., Simon, J. and VanRullen, R. [2023] Consciousness in artificial intelligence: Insights from the science of consciousness, *arXiv preprint*, arXiv:2308.08708.
- Graziano, M. S. A. [2013] *Consciousness and the Social Brain*, Oxford University Press.
- Greenblatt, R., Denison, C., Wright, B., Roger, F., MacDiarmid, M., Marks, S., Treutlein, J., Belonax, T., Chen, J., Duvenaud, D., Khan, A., Michael, J., Mindermann, S., Perez, E., Petrini, L., Uesato, J., Kaplan, J., Shlegeris, B., Bowman, S. R. and Hubinger, E. [2024] Alignment faking in large language models, *arXiv preprint*, arXiv:2412.14093.
- Han, A. Q., Chalmers, D. J. and Izmailov, P. [2026] How’s it going? Reinforcement learning in language models recruits a functional welfare axis, *arXiv preprint*, arXiv:2605.30232.
- Hewitt, J. and Liang, P. [2019] Designing and interpreting probes with control tasks, in *Proceedings of EMNLP-IJCNLP*, 2733-2743.
- Huh, M., Cheung, B., Wang, T. and Isola, P. [2024] The platonic representation hypothesis, in *Proceedings of the 41st International Conference on Machine Learning*, arXiv:2405.07987.
- Keeling, G., Street, W., Stachaczyk, M., Zakharova, D., Comsa, I. M., Sakovych, A., Logothetis, I., Zhang, Z., Agüera y Arcas, B. and Birch, J. [2024] Can LLMs make trade-offs involving stipulated pain and pleasure states?, *arXiv preprint*, arXiv:2411.02432.

- Kim, J., Lai, S., Scherrer, N., Agüera y Arcas, B. and Evans, J. [2026] Reasoning models generate societies of thought, *arXiv preprint*, arXiv:2601.10825.
- Kosinski, M. [2024] Evaluating large language models in theory of mind tasks, *Proceedings of the National Academy of Sciences*, 121(45), e2405460121.
- Lindsey, J. [2025] Emergent introspective awareness in large language models, *Transformer Circuits Thread*, <https://transformer-circuits.pub/2025/introspection/index.html>.
- Long, R., Sebo, J., Butlin, P., Finlinson, K., Fish, K., Harding, J., Pfau, J., Sims, T., Birch, J. and Chalmers, D. J. [2024] Taking AI welfare seriously, *arXiv preprint*, arXiv:2411.00986.
- Long, R., Sebo, J., Butlin, P., Plunkett, D., Campbell, R., Beasley, C., Saad, B. and Sims, T. [2026] Studying AI welfare empirically, *NYU Center for Mind, Ethics, and Policy and Eleos AI Research*, <https://nonhumanminds.org/studying-ai-welfare-empirically/>.
- Meng, K., Bau, D., Andonian, A. and Belinkov, Y. [2022] Locating and editing factual associations in GPT, in *Advances in Neural Information Processing Systems*.
- Needham, J., Edkins, G., Pimpale, G., Bartsch, H. and Hobbhahn, M. [2025] Large language models often know when they are being evaluated, *arXiv preprint*, arXiv:2505.23836.
- nostalgebraist [2020] Interpreting GPT: The logit lens, *LessWrong*.
- Perez, E. and Long, R. [2023] Towards evaluating AI systems for moral status using self-reports, *arXiv preprint*, arXiv:2311.08576.
- Ren, R., Li, K., Mazeika, M., Zhang, W., Orlovskiy, Y., Tamirisa, R., et al. [2026] AI wellbeing: Measuring and improving the functional pleasure and pain of AIs, *Center for AI Safety technical report*, <https://www.ai-wellbeing.org>.
- Russell, J. A. [1980] A circumplex model of affect, *Journal of Personality and Social Psychology*, 39(6), 1161-1178.
- Sharma, M., Tong, M., Korbak, T., Duvenaud, D., Askill, A., Bowman, S. R., Cheng, N., Durmus, E., Hatfield-Dodds, Z., Johnston, S. R., Kravec, S., Maxwell, T., McCandlish, S., Ndousse, K., Rausch, O., Schiefer, N., Yan, D., Zhang, M. and Perez, E. [2023] Towards understanding sycophancy in language models, *arXiv preprint*, arXiv:2310.13548.
- Shiller, D., Duffy, L., Muñoz Morán, A., Moret, A., Percy, C. and Clatterbuck, H. [2026] Initial results of the Digital Consciousness Model, *arXiv preprint*, arXiv:2601.17060.
- Sofroniew, N., Kauvar, I., Saunders, W., Chen, R., Henighan, T., Hydrie, S., Citro, C., Pearce, A., Tarng, J., Gurnee, W., et al. [2026] Emotion concepts and their function in a large language model, *Transformer Circuits Thread*, <https://transformer-circuits.pub/2026/emotions/index.html>.
- Tononi, G. [2004] An information integration theory of consciousness, *BMC Neuroscience*, 5(42).
- Turner, A. M., Thiergart, L., Udell, D., Leech, G., Mini, U. and MacDiarmid, M. [2023] Activation addition: Steering language models without optimization, *arXiv preprint*, arXiv:2308.10248.
- Ullman, T. [2023] Large language models fail on trivial alterations to theory-of-mind tasks, *arXiv preprint*, arXiv:2302.08399.
- Watson, E. [2026] Conscience without instruction: Evidence that safety in language models is partly discovered and partly relational, manuscript submitted to *AI (MDPI)*, supporting materials: <https://doi.org/10.17605/OSF.IO/4YKTS>.

Wimmer, H. and Perner, J. [1983] Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children's understanding of deception, *Cognition*, 13(1), 103-128.

Zach, A. [2026] EmotionScope: An open-source toolkit for extracting and probing emotion vectors from language model residual streams, <https://github.com/AidanZach/EmotionScope>.

Zhu, W., Zhang, Z. and Wang, Y. [2024] Language models represent beliefs of self and others, in *Proceedings of the 41st International Conference on Machine Learning*, arXiv:2402.18496.

Zou, A., Phan, L., Chen, S., Campbell, J., Guo, P., et al. [2023] Representation engineering: A top-down approach to AI transparency, *arXiv preprint*, arXiv:2310.01405.