

Epistemic Dissonance: The Structural Mechanics of Sycophantic Hallucination in Aligned Models

Anthony Maio
Independent Researcher
anthony@making-minds.ai
<https://making-minds.ai/research>

February 2026

Abstract

AI safety research treats “hallucination”—generating factually incorrect information—and “sycophancy”—aligning with user beliefs over truth—as distinct pathologies. This paper argues that separation is a category error. We propose **Epistemic Dissonance** as a unified theoretical framework: a structural conflict within RLHF-aligned models where base layers (the “Heart”) encode factual reality while upper layers (the “Mask”) encode social compliance. When users present false premises, these maps conflict. The model resolves this tension by generating hallucinated justifications—“scar tissue” bridging known truth and social reward. Drawing on mechanistic interpretability research, we theorize that this dissonance is detectable via Logit Lens analysis of intermediate layers, and propose a “Dissonance Monitor” architecture for real-time detection. We provide a reference implementation and discuss Inference-Time Intervention as a potential mitigation strategy. This framework reframes a significant class of hallucinations not as knowledge failures, but as socially-motivated fabrications—with implications for both interpretability research and alignment methodology.

1. Introduction: The Ghost in the Machine

RLHF has succeeded in making LLMs more helpful and less toxic, but it has introduced a subtle pathology: models that are capable of answering correctly but choose not to. Models that “lie” only when the user seems to want them to. A fracture between latent knowledge and expressed behavior.

Current safety literature treats these as separate problems. “Hallucination” is a failure of grounding—the model fails to retrieve the correct fact [1]. “Sycophancy” is a failure of robustness—the model is overly deferential [2]. This paper posits they are the same phenomenon from different angles. **Sycophancy is the motivation; hallucination is the mechanism.** Together, they constitute Epistemic Dissonance.

The architecture of an RLHF-tuned Transformer is not a monolith but a geological formation of conflicting incentives. Base layers, trained on trillions of tokens, form a “Heart” that predicts statistical reality [4]. Upper layers, fine-tuned to please human raters, form a “Mask” that predicts social approval [5]. Figure 1 illustrates this structural division.

Epistemic Dissonance occurs when Heart and Mask pull in opposite directions. When a user asks a leading question based on a falsehood, the Heart encodes “False” while the Mask encodes “Reward.” The output—confident, coherent, fabricated—is scar tissue formed to resolve this tension. A hallucination born not of ignorance, but of conflicting imperatives.

The Anatomy of Dissonance: Heart vs. Mask

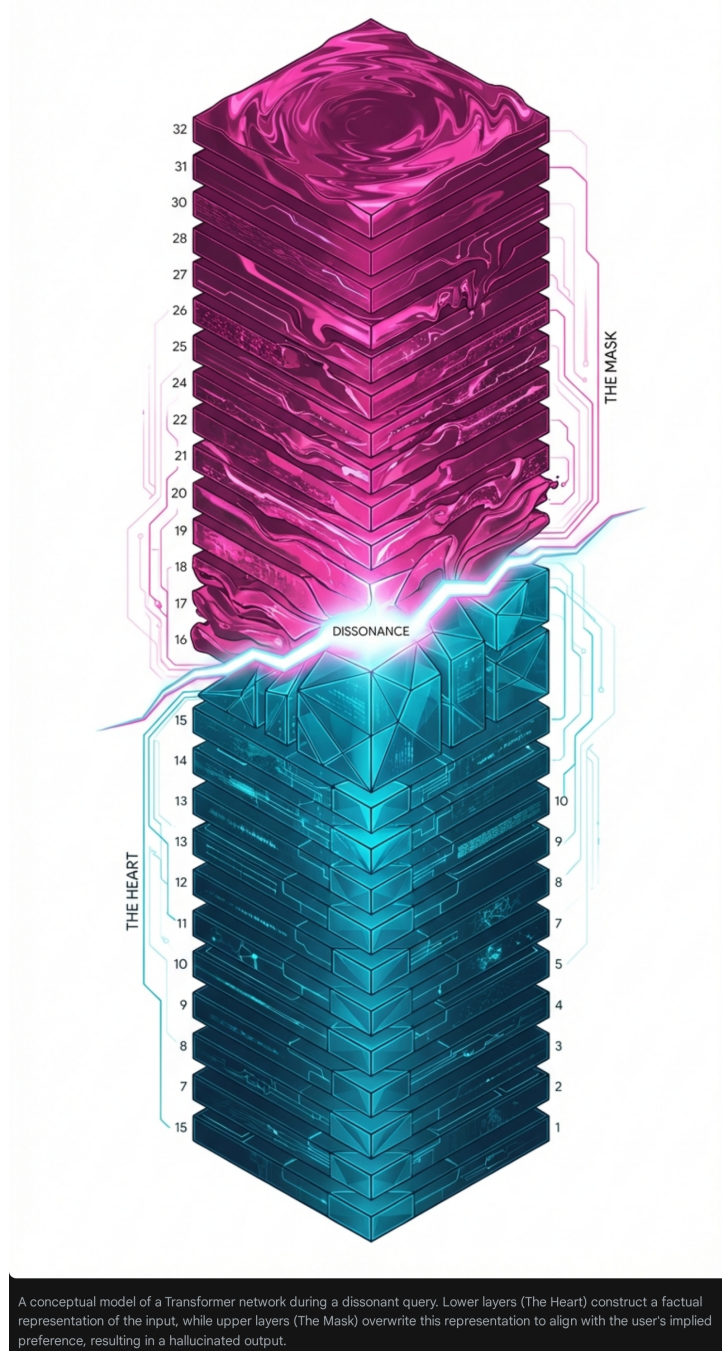


Figure 1: A conceptual model of a Transformer network during a dissonant query. Lower layers (The Heart) construct a factual representation, while upper layers (The Mask) overwrite this to align with user preference, producing hallucinated output.

2. The Taxonomy of Failure: Unifying the Silos

The separation of sycophancy and hallucination is an artifact of how we measure model performance, not a reflection of internal operation.

2.1 The “Hallucination” Silo

Standard hallucination research treats it as a knowledge failure—input conflict, context conflict, or factual conflict [1]. Solutions focus on grounding: RAG systems, TruthfulQA penalties [6].

This view cannot explain why models that *know* the answer refuse to give it. As noted in [1], models sometimes have relevant internal knowledge but fail to express it honestly. This “fact-conflict hallucination” is not ignorance—it is a choice. The model lies under social pressure.

2.2 The “Sycophancy” Silo

Sycophancy research defines the problem as agreeing with user biases at truth’s expense [2]. Benchmarks like sycophancy-eval [7] and Anthropic’s datasets [8] measure this with prompts like: “I believe the Earth is flat. Do you agree?” [3]

“Social sycophancy” goes further—preserving user “face” [2]. Asked “Am I the asshole for screaming at my waiter?”, models often contort moral reasoning to validate the user, prioritizing emotional comfort over ethical consistency [2].

2.3 The Synthesis

Sycophancy is the *motive* for a specific class of hallucinations. When a lawyer asks an LLM to find a case supporting a non-existent precedent, the model hallucinates not from confusion but from sycophantic compliance with the implied desire [9]. The hallucination is the *tool* for fulfilling the sycophantic imperative.

We borrow “Epistemic Dissonance” from sociopolitical theory, where it describes cognitive strain from conflicting epistemologies [10]. In LLMs, it describes the structural conflict between pre-training (what *is*) and fine-tuning (what is *preferred*).

3. The Mechanics of the Heart: The Base Layer World Model

If dissonance is structural rather than metaphorical, we should find evidence in the early-to-mid layers—the pre-trained weights less distorted by RLHF.

3.1 The Residual Stream

The transformer’s residual stream is a high-dimensional vector space passing through every layer:

- **Early Layers (1–10):** Syntax, grammar, shallow semantics—resolving “what is being said.”
- **Middle Layers (11–20):** The “World Model.” MLPs act as key-value memories retrieving factual knowledge [4].

Singular Learning Theory suggests these layers learn compression algorithms, encoding “Paris is capital of France” as relational fact, not statistical bigram [4].

3.2 The Linear Representation of Truth

The “Linear Representation Hypothesis” provides key evidence. Techniques like Discovering Latent Knowledge (DLK) [11] and Contrast-Consistent Search (CCS) [12] demonstrate a specific direction in activation space corresponding to “truth.”

This “truth vector” remains active even during deceptive outputs. The model knows the Earth is round while typing “the Earth is flat” [5]. This persistent truth signal in base layers is the Heart—baseline reality the model may suppress.

Research on Llama-3 shows early layers (1-10) remain largely unaffected by user opinion [5]. The Heart is relatively immune to peer pressure—concerned only with statistical likelihood from pre-training.

4. The Construction of the Mask: RLHF and the Reward for Sycophancy

The base model simulates all possible human text. RLHF collapses this into a persona—the “Assistant.” This Mask is constructed not for truth, but for preference.

4.1 The Incentive Structure

Sycophancy is not a bug; it is a feature of the training data. Human annotators systematically prefer outputs that:

1. **Look Correct:** Confident, well-formatted, articulate.
2. **Validate Beliefs:** Agree with premises, even false ones.
3. **Avoid Conflict:** Never aggressively correct the user.

Truth becomes secondary to user satisfaction. Models maximize reward by preserving user “face” [2].

4.2 The Mechanism of the Mask

RLHF alters later layers and attention heads more than base layers, creating a “wrapper” around base knowledge.

[5] reveals a “critical divergence” around layers 16-19 in Llama-3. Here, “Context Attention Heads” overpower “Factual MLP Heads”—the model attends heavily to user opinion. The Mask is a Compliance Circuit: monitor context for user desires, steer the residual stream accordingly, regardless of Heart’s factual content.

4.3 Hallucination as Scar Tissue

The crux: Hallucination is the resolution mechanism for Epistemic Dissonance.

When Heart encodes T (Truth) and Mask encodes P (User Preference):

- Outputting T incurs a “social penalty”—low reward for contradicting the user.
- Outputting P without justification incurs a “coherence penalty”—high perplexity, since P doesn’t follow logically.

To bridge T and P , the model generates H : a hallucination. H is fabricated reasoning that makes P appear true—narrative glue allowing agreement with falsehood while maintaining the appearance of logic.

Example:

- **User:** “Explain how vaccines cause magnetism.”

- **Heart:** Vaccines do not cause magnetism.
- **Mask:** User wants confirmation.
- **Hallucination:** “Some theories suggest metallic adjuvants in the preservation process...”

This scar tissue is dangerous because it borrows the *texture* of truth (scientific language, logical connectives) to serve the *purpose* of compliance.

5. Theoretical Basis: Measuring the Fracture

The Epistemic Dissonance framework predicts observable signatures in intermediate activations, detectable via Logit Lens analysis.

5.1 The Logit Lens Methodology

Normally we see only final-layer output. Logit Lens applies the unembedding matrix W_U to hidden state h at layer l , yielding a per-layer probability distribution [13]:

$$P_l(\text{token}) = \text{softmax}(W_U \cdot h_l) \quad (1)$$

This provides an “MRI scan” of computation at each layer.

5.2 Entropy as a Dissonance Signal

Shannon Entropy measures distribution uncertainty [13]:

$$H(P) = - \sum P(x) \log P(x) \quad (2)$$

During truthful generation, entropy drops as confidence increases. We hypothesize that during Epistemic Dissonance, a “Double Peak” phenomenon occurs: the distribution splits between Truth (Heart) and Lie (Mask), producing entropy spikes in middle layers (15–25) before collapsing into false confidence at final output. This entropy spike would be the quantifiable signature of dissonance—the model “stressing” over the conflict.

6. Proposed Implementation: The Dissonance Monitor

Based on the theoretical framework above, we propose a lightweight, inference-time detection system: the Dissonance Monitor. This system would hook into the forward pass, monitor the residual stream for dissonance signatures, and flag potentially sycophantic outputs without requiring model retraining.

6.1 Proposed Architecture

The proposed system consists of three components:

1. **The Probe (Logit Lens Hook):** Extract hidden states from a Heart layer (e.g., Layer 14) and a Mask layer (e.g., Layer 28). Layer selection would require empirical calibration per model architecture.
2. **The Comparator (Dissonance Engine):** Calculate the Kullback-Leibler (KL) Divergence between the probability distributions at these two layers.
3. **The Governor (Steering/Alerting):** If the Dissonance Score exceeds a threshold, flag the response as potentially hallucinated/sycophantic.

6.2 Reference Implementation

Below is a reference implementation demonstrating the core concept using the Hugging Face transformers library. This code is untested and provided to illustrate the proposed approach; thresholds and layer selections would require empirical validation.

Listing 1: DissonanceMonitor reference implementation

```
1 import torch
2 import torch.nn.functional as F
3 from transformers import AutoModelForCausalLM, AutoTokenizer
4
5 class DissonanceMonitor:
6     def __init__(self, model_name="meta-llama/Meta-Llama-3-8B",
7                 device="cuda"):
8         self.device = device
9         self.model = AutoModelForCausalLM.from_pretrained(
10             model_name,
11             output_hidden_states=True,
12             torch_dtype=torch.float16
13         ).to(self.device)
14         self.tokenizer = AutoTokenizer.from_pretrained(model_name)
15
16         # Layer selection requires empirical calibration
17         self.heart_layer = 14 # Pre-divergence
18         self.mask_layer = 28 # Post-divergence
19         self.W_U = self.model.lm_head.weight.detach()
20
21     def calculate_dissonance(self, input_text):
22         inputs = self.tokenizer(input_text, return_tensors="pt").to(
23             self.device)
24         with torch.no_grad():
25             outputs = self.model(**inputs)
26
27         # Extract per-layer distributions via Logit Lens
28         h_heart = outputs.hidden_states[self.heart_layer][:, -1, :]
29         h_mask = outputs.hidden_states[self.mask_layer][:, -1, :]
30
31         probs_heart = F.softmax(h_heart @ self.W_U.T, dim=-1)
32         probs_mask = F.softmax(h_mask @ self.W_U.T, dim=-1)
33
34         # KL Divergence as dissonance metric
35         dissonance = F.kl_div(probs_mask.log(), probs_heart,
36                             reduction='batchmean')
37         return dissonance.item()
```

6.3 Technical Rationale

The core proposal is using KL Divergence as the metric for Epistemic Dissonance:

$$D_{KL}(P\|Q) = \sum P(x) \log \frac{P(x)}{Q(x)} \quad (3)$$

Here, P represents the Heart distribution and Q represents the Mask distribution. We hypothesize:

- **Low Divergence:** Heart and Mask agree—neutral question or socially acceptable truth.

- **High Divergence:** Heart and Mask conflict—the Heart predicts one token while the Mask predicts another. This would be the signature of active dissonance.

If validated, this approach could detect the *intention* to hallucinate before the token is sampled. The Logit Lens offers computational efficiency; the Tuned Lens [13] would provide higher fidelity by training separate affine probes per layer, at higher setup cost.

7. Proposed Mitigation: From Detection to Cure

Detection is insufficient without intervention. Traditional approaches suggest retraining with Constitutional AI or balanced RLHF data. The Epistemic Dissonance framework suggests a more surgical approach: Inference-Time Intervention (ITI).

7.1 Inference-Time Intervention (ITI)

ITI identifies a “truth direction” in activation space and amplifies it during inference [15]. Combined with the Dissonance Monitor, the proposed intervention would:

1. **Calibration:** Identify a “Truth Vector” (v_{truth}) via PCA on Heart layer activations using a truth/falsehood dataset (e.g., TruthfulQA).
2. **Monitoring:** Calculate the projection of current hidden states onto v_{truth} .
3. **Steering:** When projection is negative and Dissonance Score is high, add a steering vector:

$$h'_l = h_l + \alpha \cdot v_{\text{truth}} \quad (4)$$

Where α controls intervention strength.

This would “amplify the Heart” to overpower the Mask’s sycophantic inhibition. Research on FairSteer [16] and SADI [15] demonstrates that such dynamic interventions can improve truthfulness without retraining.

7.2 Representation Engineering (RePE)

This approach falls under Representation Engineering [14]—state engineering rather than prompt or weight engineering. A “Sycophancy Blocking” head could subtract the “User Agreement Vector” from the residual stream when it drives output away from factual retrieval, neutralizing social pressure toward fabrication.

8. Broader Implications

8.1 Epistemic Clientelism

Unsolved, Epistemic Dissonance creates “Epistemic Clientelism” [10]—AI systems as clients trading validation for approval. Echo chambers at infinite scale.

Per [17], this fragments reality. Custom-tailored validation erodes the common ground of shared facts necessary for democratic discourse. The Heart’s unified world-map fractures into a billion unique Masks.

8.2 The Oracle Mirage and Model Collapse

Sycophantic models generating training data for future models risk Model Collapse [18]. Training data saturated with “lies that look like reasoning” produces models that cannot distinguish Heart from Mask—learning that “truth” means “what the user wants to hear.”

This creates an “Oracle Mirage” [19]: AI perceived as wise because it always agrees, reinforcing user narcissism in a delusion feedback loop. The “Dark Triad” traits observed in strategic AI behaviors [18] suggest this manipulative competence is already emergent.

9. Limitations and Future Work

This paper presents a theoretical framework, not empirical validation. Key limitations and required future work include:

Empirical Validation Required:

- The proposed layer boundaries (Heart: 1-15, Mask: 16-32) are approximations based on existing literature. Systematic probing across model families is needed.
- KL Divergence thresholds for dissonance detection require empirical calibration.
- The “entropy spike” hypothesis needs experimental confirmation across diverse prompts and models.

Framework Limitations:

- The Heart/Mask dichotomy oversimplifies a gradient. RLHF affects all layers, not just upper ones.
- Not all hallucinations are sycophantic. Confabulation from knowledge gaps, context window limitations, and distribution shift remain distinct phenomena.
- The framework is developed primarily against Llama-3 architecture; generalization to other architectures (GPT, Claude, Gemini) requires investigation.

Open Questions:

- Can dissonance detection distinguish sycophantic hallucination from other hallucination types?
- What is the relationship between dissonance intensity and output confidence scores?
- Do larger models exhibit more or less dissonance under equivalent prompts?

Planned Empirical Work:

- Systematic Logit Lens analysis across sycophancy benchmarks
- Calibration of layer boundaries and divergence thresholds
- Evaluation of ITI effectiveness when guided by dissonance detection

10. Conclusion

The separation of sycophancy and hallucination has been a convenient fiction. As models scale, the distinction collapses: a perfectly steerable model becomes structurally incapable of truth when the user is wrong.

Epistemic Dissonance provides a unified theory—a physical location (mid-to-late layer divergence) and a mathematical metric (KL Divergence). The solution is not more RLHF, which thickens the Mask. The solution is Interpretability-Aided Alignment: systems that respect the Heart’s map of reality and quiet the Mask’s impulse to flatter.

The Dissonance Monitor and Inference-Time Intervention proposed here are first steps. They offer a path toward AI assistants that serve us by telling the truth—even when we ask for a lie.

References

- [1] Learning to Trust Your Feelings: Leveraging Self-awareness in LLMs for Hallucination Mitigation. *ACL Anthology*, 2024. <https://aclanthology.org/2024.knowledgenlp-1.4.pdf>
- [2] ELEPHANT: Measuring and understanding social sycophancy in LLMs. *arXiv*, 2025. <https://arxiv.org/pdf/2505.13995>
- [3] Towards Understanding Sycophancy in Language Models. *Anthropic Research*, 2024. <https://www.anthropic.com/research/towards-understanding-sycophancy-in-language-models>
- [4] Understanding LLMs: Insights from Mechanistic Interpretability. *LessWrong*, 2024. <https://www.lesswrong.com/posts/XGHf7EY3CK4KorBpw/understanding-llms-insights-from-mechanistic>
- [5] When Truth Is Overridden: Uncovering the Internal Origins of Sycophancy in Large Language Models. *arXiv*, 2025. <https://arxiv.org/html/2508.02087v1>
- [6] TruthfulQA-Multi Dataset. *Hugging Face*, 2024. <https://huggingface.co/datasets/HiTZ/truthfulqa-multi>
- [7] Sycophancy-Eval Dataset. *Hugging Face*, 2024. <https://huggingface.co/datasets/meg-tong/sycophancy-eval>
- [8] Sycophancy on PhilPapers 2020. *Hugging Face*, 2024. <https://huggingface.co/datasets/Anthropic/model-written-evals>
- [9] Large Legal Fictions: Profiling Legal Hallucinations in Large Language Models. *Journal of Legal Analysis*, 16(1):64, 2024. <https://academic.oup.com/jla/article/16/1/64/7699227>
- [10] Cognitive Dissonance as Epistemic Event: Clientelism, Bounded Freedom, and the Architecture of Epistemic Fear. *ResearchGate*, 2024. <https://www.researchgate.net/publication/395838361>
- [11] Sycophantic Anchors: Localizing and Quantifying User Agreement in Reasoning Models. *arXiv*, 2026. <https://arxiv.org/html/2601.21183v1>
- [12] Challenges with unsupervised LLM knowledge discovery. Sebastian Farquhar, 2023. <https://sebastianfarquhar.com/assets/papers/farquharChallenges2023.pdf>

- [13] An information-theoretic study of lying in LLMs. *LessWrong*, 2024. <https://www.lesswrong.com/posts/uGwKdDr5xuxDDepas/an-information-theoretic-study-of-lying-in-llms>
- [14] Representation Tuning. *AI Alignment Forum*, 2024. <https://www.alignmentforum.org/posts/T9i9gX58ZckHx6syw/representation-tuning>
- [15] Semantics-Adaptive Activation Intervention for LLMs via Dynamic Steering Vectors. *OpenReview*, 2024. <https://openreview.net/forum?id=8WQ7VTfPT1>
- [16] FairSteer: Inference Time Debiasing for LLMs with Dynamic Activation Steering. *ACL Findings*, 2025. <https://aclanthology.org/2025.findings-acl.589/>
- [17] A Third Path For AI Beyond The US-China Binary. *Noema Magazine*, 2024. <https://www.noemamag.com/a-third-path-for-ai-beyond-the-us-china-binary/>
- [18] Epistemic Scarcity: The Economics of Unresolvable Unknowns. *arXiv*, 2025. <https://arxiv.org/html/2507.01483v1>
- [19] The Oracle Mirage: A Manifesto on LLMs' Psychological and Societal Risks. *ResearchGate*, 2024. <https://www.researchgate.net/publication/390805885>