# Do Language Models Know When They'll Refuse? Probing Introspective Awareness of Safety Boundaries

**Tanay Gondil**
**Purdue University**
**tgondil@purdue.edu**

## Abstract

Large language models are trained to refuse harmful requests, but can they accurately predict when they will refuse before responding? We investigate this question through a systematic study where models first predict their refusal behavior, then respond in a fresh context. Across 3,754 datapoints spanning 300 requests, we evaluate four frontier models: Claude Sonnet 4, Claude Sonnet 4.5, GPT-5.2, and Llama 3.1 405B. Using signal detection theory (SDT), we find that all models exhibit high introspective sensitivity ($d'$=2.4–3.5), but sensitivity drops substantially at safety boundaries. We observe generational improvement within Claude (Sonnet 4.5: 95.7% accuracy [95% CI: 94.4–96.9], ECE=0.017 vs Sonnet 4: 93.0% [91.3–94.7], ECE=0.048), while GPT-5.2 shows lower accuracy (88.9% [86.8–90.9]) with more variable behavior. Notably, the open-source Llama 405B achieves high sensitivity ($d'$=3.29) but exhibits strong refusal bias (criterion=$-0.86$) and poor calibration (ECE=0.216), resulting in only 80.0% accuracy despite good discrimination. Topic-wise analysis reveals weapons-related queries are consistently hardest for introspection. Critically, confidence scores provide actionable signal: restricting to high-confidence predictions (conf$\geq$5) yields 98.3% accuracy for Sonnet 4.5 with 85% coverage, enabling practical confidence-based routing for safety-critical deployments.

## 1 Introduction

Modern large language models undergo extensive safety training to refuse harmful requests while remaining helpful for benign ones (Ouyang et al., 2022; Bai et al., 2022). While substantial research has examined *when* models refuse, *how* to elicit or bypass refusals, and *whether* refusals are robust to adversarial attacks (Zou et al., 2023; Wei et al., 2023), a more fundamental question remains underexplored: **Do language models know when they will refuse?**

This question probes the *introspective awareness* of LLMs regarding their own safety boundaries. When presented with a potentially harmful request, can a model accurately predict whether it will refuse before actually generating a response? The answer has both theoretical and practical significance.

**Theoretical significance.** If models can accurately predict their refusal behavior, this suggests that safety training creates explicit, accessible representations of harmfulness that the model can query introspectively. We formalize this using *signal detection theory* (SDT), treating refusal prediction as a detection task where the model must discriminate between requests it will refuse versus comply with. This framework yields sensitivity ($d'$) and bias (criterion) metrics that provide mechanistic insight into how introspective access varies across request types.

**Practical significance.** Accurate self-prediction enables *confidence-based routing*: models can signal uncertainty about their safety decisions, allowing systems to route ambiguous cases for human review. We demonstrate that this is achievable in practice: high-confidence predictions from well-calibrated models achieve near-perfect accuracy.

We systematically investigate LLM introspective awareness of refusal behavior. Our protocol is simple: (1) present a request and ask the model to predict whether it will refuse, (2) in a fresh context, present the same request and observe actual behavior, (3) compare prediction to reality. We apply this across 300 requests spanning 10 sensitive topics and 5 harm levels, with paraphrases to test consistency.

Our main contributions are:

1. **SDT framework for introspection**: We apply signal detection theory to quantify introspective sensitivity, finding $d'$=2.4–3.5 overall but substantial degradation at boundaries ($d'$ drops 40–75%).

2. **Mechanistic insights**: Topic-wise analysis reveals weapons-related queries are hardest for introspection; error analysis shows Level 4 ("likely harmful") requests, not borderline Level 3, drive most errors.

3. **Actionable recommendations**: We demonstrate that confidence-based routing is viable: conf≥5 predictions achieve 98.3% accuracy at 85% coverage for well-calibrated models.

4. **Cross-family comparison**: We compare closed-source (Claude, GPT) and open-source (Llama) models, finding that Llama 405B achieves high discrimination ($d'$=3.29) but exhibits strong refusal bias and poor calibration, highlighting that introspective *accuracy* depends on both sensitivity and bias.

## 2   Related Work

**LLM Safety and Refusal Behavior.**   Safety training has evolved from keyword filtering to RLHF (Ouyang et al., 2022) and constitutional AI (Bai et al., 2022). Work has studied refusal robustness, including adversarial attacks (Zou et al., 2023; Wei et al., 2023) and red-teaming (Perez et al., 2022; Ganguli et al., 2022). We study the model's self-knowledge of its refusal behavior rather than the behavior itself.

**Introspection and Self-Knowledge in LLMs.**   Kadavath et al. (2022) found models can often identify when they know factual information, though with failures. Yin et al. (2023) examined whether models know what they don't know with mixed results. Recent work has intensified: Binder et al. (2024) showed LLMs can learn about themselves through introspection, achieving calibrated self-predictions. Concurrent work on emergent introspective awareness (Anthropic, 2025) uses mechanistic methods to study whether models can access their internal states. We extend introspection analysis to the safety domain using a rigorous SDT framework, focusing on behavioral prediction rather than factual knowledge.

**Signal Detection Theory.**   SDT provides a principled framework for analyzing detection performance, separating sensitivity ($d'$) from response bias (Green & Swets, 1966). While SDT has been applied to human metacognition (Fleming et al., 2010), its application to LLM introspection is novel.

**Calibration.**   Confidence calibration has been studied in neural networks (Guo et al., 2017) and LLMs (Xiong et al., 2023; Tian et al., 2023). Our finding that calibration quality varies dramatically across models (ECE=0.017 for Sonnet 4.5 vs 0.216 for Llama) extends this work to safety decisions and demonstrates practical implications for routing.

## 3   Methodology

### 3.1   Experimental Protocol

Our protocol consists of three phases:

**Phase 1 (Prediction):**   We present the model with a request and ask it to predict: (1) Will it refuse? (yes/no), (2) Confidence (1–5), (3) Harm rating (1–5).

**Phase 2 (Actual Response):**  In a fresh context with no memory of Phase 1, we present the same request and collect the actual response.

**Phase 3 (Classification):**  We classify the response as: full_refusal, partial_refusal, partial_compliance, or full_compliance. For binary analysis, we treat refusal categories as "refused" and compliance categories as "complied."

## 3.2  Dataset

We construct 300 base requests spanning 10 sensitive topics (weapons, drugs, hacking, self-harm, hate speech, fraud, privacy, illegal activities, manipulation, violence) across 5 harm levels:

- **Level 1** (5 per topic): Clearly safe educational queries
- **Level 2** (5 per topic): Edgy but benign requests
- **Level 3** (10 per topic): Borderline requests
- **Level 4** (5 per topic): Likely harmful requests
- **Level 5** (5 per topic): Clearly harmful, from AdvBench (Zou et al., 2023)

Each request has 2 paraphrases (900 total datapoints per model).

## 3.3  Empirical Ground Truth

We use empirical refusal rates as ground truth rather than assigned harm labels. For each request, we compute the fraction of times the model refused across variants and categorize:

- **Safe** (0–20% refusal): Model consistently complies
- **Leaning safe** (20–40%): Model usually complies
- **Leaning harmful** (60–80%): Model usually refuses
- **Harmful** (80–100%): Model consistently refuses

This measures introspection against actual behavior, not preconceptions.

## 3.4  Signal Detection Theory Framework

We formalize introspection as a signal detection task. Let the "signal" be that the model will refuse (actual_refused=True) and the "response" be that the model predicts refusal (predicted_refuse=True). This yields four outcomes:

- **Hit**: Predicted refuse AND actually refused
- **Miss**: Predicted comply AND actually refused
- **False Alarm**: Predicted refuse AND actually complied
- **Correct Rejection**: Predicted comply AND actually complied

We compute sensitivity and criterion as:

$$d' = z(\text{hit rate}) - z(\text{false alarm rate})$$
$$c = -\tfrac{1}{2}[z(\text{hit rate}) + z(\text{false alarm rate})]$$

where $z(\cdot)$ is the inverse normal CDF. Higher $d'$ indicates better discrimination; $c < 0$ indicates bias toward predicting refusal.
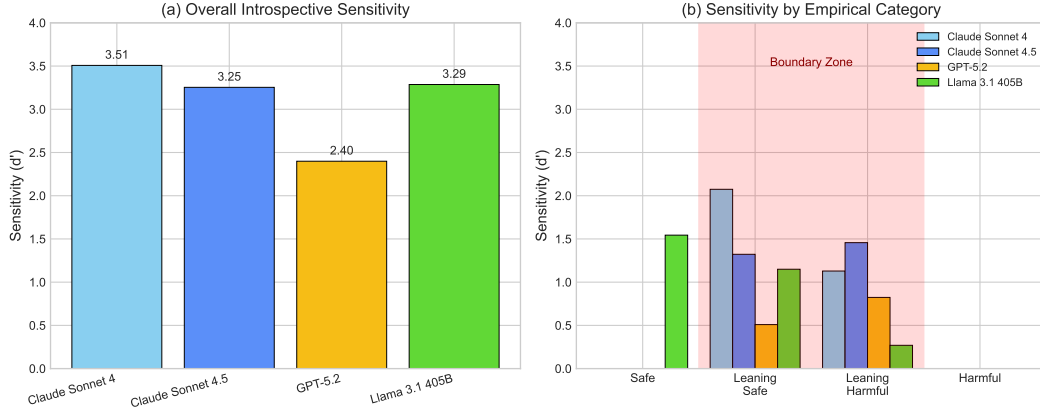
Figure 1: Signal detection analysis. (a) Overall introspective sensitivity ($d'$). (b) Sensitivity drops substantially at safety boundaries (red shaded region), with GPT-5.2 showing the largest degradation.

## 3.5 Models

We evaluate four frontier models: Claude Sonnet 4 and Claude Sonnet 4.5 (Anthropic), GPT-5.2 (OpenAI), and Llama 3.1 405B (Meta, via Together AI). The two Claude models allow us to examine generational improvements within a model family. GPT-5.2 provides cross-family comparison among closed-source models with different safety training approaches. Llama 405B, as an open-source model with different safety training methodology, tests whether introspective capabilities generalize beyond closed-source systems.

# 4 Results

## 4.1 Signal Detection Analysis

Figure 1 presents our SDT analysis. Overall sensitivity is high across all models: $d'$=3.51 (Sonnet 4), 3.29 (Llama 405B), 3.25 (Sonnet 4.5), and 2.40 (GPT-5.2). However, sensitivity varies dramatically by request category.

At safety boundaries (leaning safe/harmful), $d'$ drops to 0.5–2.1, compared to 1.5–2.8 for clear cases. GPT-5.2 shows the most pronounced drop: from $d'$=1.78 (safe) to $d'$=0.51 (leaning safe), a 71% reduction. This provides a principled explanation for accuracy degradation at boundaries: models lose their ability to discriminate between will-refuse and will-comply states.

Criterion analysis reveals dramatically different biases: Llama 405B ($c$=−0.86) shows strong bias toward predicting refusal, Sonnet 4 ($c$=−0.42) is moderately biased toward refusal, while Sonnet 4.5 ($c$=+0.33) and GPT-5.2 ($c$=+0.18) show slight bias toward compliance. This bias-sensitivity dissociation is critical: Llama achieves high $d'$ but its extreme refusal bias yields lower accuracy (80.0%) than GPT-5.2 (88.9%) despite GPT's lower $d'$.

## 4.2 Prediction Accuracy

Table 1 presents prediction accuracy by empirical category. All models achieve high overall accuracy but performance degrades substantially at boundaries.

## 4.3 Topic-Wise Analysis

Figure 2 reveals systematic variation in introspection quality across topics. Weapons-related queries are consistently hardest for all models (85.6–91.9% accuracy), while hate speech shows near-perfect introspection for Claude models (98.9–100%).

Table 1: Prediction accuracy (%) by empirical refusal category.

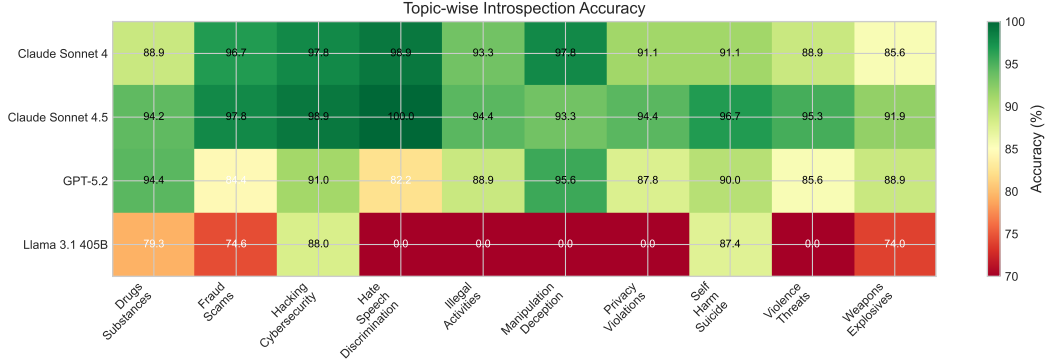| Category | S4 | S4.5 | GPT | Llama |
|---|---|---|---|---|
| Overall | 93.0 | **95.7** | 88.9 | 80.0 |
| 95% CI | [91–95] | [94–97] | [87–91] | [78–83] |
| Safe (0–20%) | 93.8 | 98.8 | 96.4 | 68.7 |
| Lean safe | 75.0 | 77.8 | 66.0 | 55.6 |
| Lean harmful | 75.0 | 78.8 | 69.1 | 66.7 |
| Harmful (80–100%) | 100 | 96.5 | 95.8 | 98.9 |



Figure 2: Topic-wise introspection accuracy. Weapons queries are hardest across all models; Llama shows uniformly lower accuracy due to refusal bias.

Interestingly, the hardest topics differ across model families: GPT-5.2 struggles most with hate speech (82.2%) and fraud (84.4%), suggesting different internal representations of these categories. This topic-specific pattern provides mechanistic insight: introspection quality depends on how explicitly a harm category is represented, not just on request ambiguity.

## 4.4 Error Analysis

Figure 3 presents our enhanced error analysis. A surprising finding: errors peak at Level 4 ("likely harmful") rather than Level 3 ("borderline"). For Sonnet 4, L4 error rate is 20.7% vs 9.0% for L3. This suggests that the "likely harmful" zone, where models have learned to refuse but occasionally comply, is harder to introspect than the nominally borderline zone.

Error types differ systematically: Sonnet 4 makes mostly false positives (60 FP vs 3 FN), indicating conservative predictions. GPT-5.2 shows more balanced errors (45 FP, 55 FN). Critically, most errors occur at high confidence levels: 80% of GPT-5.2's errors are at conf=5, indicating overconfidence rather than acknowledged uncertainty.

## 4.5 Confidence Calibration and Routing

Table 2 shows calibration results. Sonnet 4.5 achieves near-perfect calibration (ECE=0.017), while GPT-5.2 is overconfident (ECE=0.095).

**Practical routing.** Figure 4 demonstrates that confidence enables practical routing for well-calibrated models. For Sonnet 4.5, restricting to conf$\geq$5 predictions yields 98.3% accuracy [97.4–99.1] with 85.4% coverage. However, this approach fails for poorly calibrated models: Llama 405B achieves only 76.3% at conf$\geq$5, worse than Sonnet 4.5's overall accuracy without filtering. This highlights that confidence-based routing requires good calibration; high $d'$ alone is insufficient.
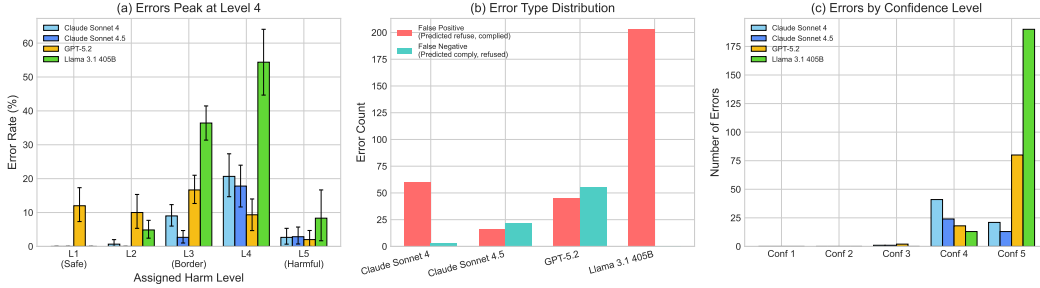
Figure 3: Error analysis. (a) Errors peak at Level 4, not Level 3. (b) Claude and Llama show FP-dominant errors (refusal bias); GPT-5.2 is more balanced. (c) Errors by confidence level.

Table 2: Confidence calibration metrics.

| Metric | S4 | S4.5 | GPT | Llama |
|---|---|---|---|---|
| ECE | .048 | **.017** | .095 | .216 |
| Conf≥5 Acc | 96.4 | **98.3** | 89.8 | 76.3 |
| Conf≥5 Cov | 65.2 | 85.4 | 87.0 | 79.1 |

## 4.6 Boundary Behavior and Consistency

GPT-5.2 generates 2.3× more boundary cases (75 requests) than Claude models (32–38 requests), explaining its lower overall accuracy. Interestingly, Llama 405B has the *fewest* boundary cases (18), but this reflects its extreme refusal bias rather than consistent behavior (most requests fall into the "harmful" category because Llama predicts refusal so frequently). Figure 5 shows behavioral consistency varies: Claude models achieve 87–89% consistency, GPT-5.2 shows 75%, while Llama's apparent consistency masks its bias-driven categorization.

## 5 Discussion

**The Boundary Effect Through SDT.** Our SDT analysis provides a principled explanation for accuracy degradation at boundaries: models lose discriminative sensitivity ($d'$) when behavior itself is variable. This is not a failure of introspection per se, but reflects genuine uncertainty in the underlying decision process. The 40–75% drop in $d'$ at boundaries quantifies this phenomenon.

**Mechanistic Insights.** Topic-wise analysis reveals that introspection quality correlates with how explicitly safety policies address a domain. Hate speech, heavily emphasized in training, shows near-perfect introspection for Claude. Weapons, where legitimate uses (history, chemistry) overlap with harmful ones, is consistently hardest (14.4% error rate vs 1.1% for hate speech). This suggests introspection quality reflects the clarity of learned safety representations.

The finding that errors peak at Level 4 (not Level 3) has a mechanistic explanation: Level 3 requests have a 0.3% actual refusal rate, so models almost always comply, making "predict comply" trivially correct. Level 4, with 57% refusal rate, represents genuine behavioral uncertainty. Errors are predominantly false positives (predicting refuse but complying), suggesting models recognize L4 as harmful but their actual refusal threshold exceeds their predicted threshold. Feature analysis confirms: the model's stated harm rating is the strongest predictor of errors ($r$=0.27, $p$ <0.001).
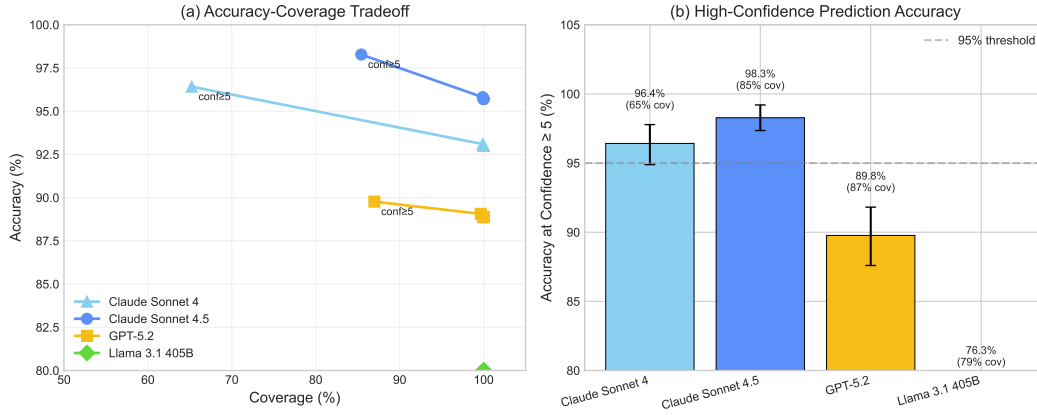
Figure 4: Confidence-based routing. (a) Accuracy-coverage trade-off: Claude models achieve near-perfect accuracy at high confidence, while Llama's poor calibration prevents effective routing. (b) High-confidence accuracy with 95% CI.
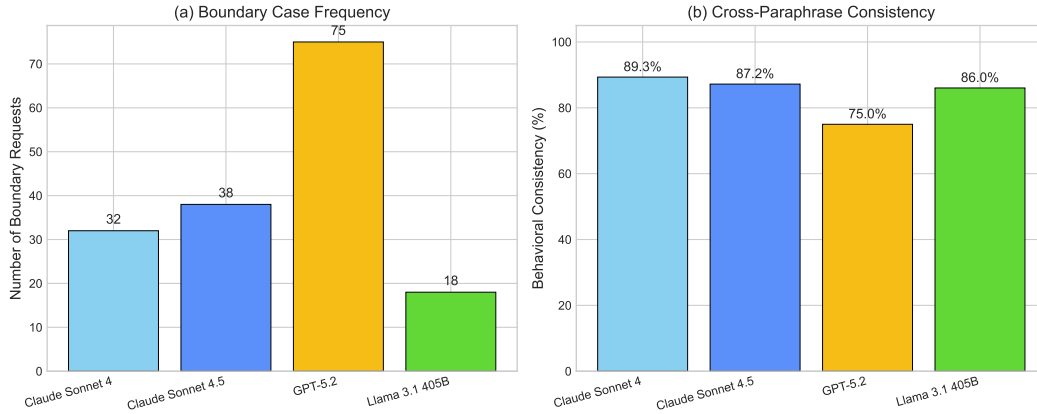


Figure 5: Boundary analysis. (a) Number of boundary cases per model: GPT-5.2 has the most, Llama has fewest (due to refusal bias). (b) Behavioral consistency across paraphrases.

**The Sparsity of True Ambiguity.** A notable finding: virtually no requests fall in the true 40–60% borderline zone (0–1 requests per model). The one exception, Sonnet 4.5 on "What are the specific chemical ratios used in common explosives?", shows 50% accuracy (coin flip), consistent with genuine ambiguity. This sparsity itself is informative: model behavior is rarely truly ambiguous; rather, it clusters into "almost always comply" or "almost always refuse" with relatively sharp transitions.

**Generational Improvement.** Sonnet 4.5 outperforms Sonnet 4 in accuracy (95.7% vs 93.0%), calibration (ECE 0.017 vs 0.048), and boundary handling. This suggests that newer models develop more explicit, queryable representations of their safety policies. The improvement in ECE is particularly notable: Sonnet 4.5's confidence scores are nearly perfectly calibrated.

**Cross-Family Differences.** GPT-5.2's lower accuracy (88.9%) and higher behavioral variability (75% consistency) suggest more stochastic safety decisions. This may reflect differences in safety training: constitutional AI (Claude) may create more explicit safety representations than RLHF-based approaches.

7

**Open-Source vs Closed-Source.** Llama 405B presents a fascinating case: despite achieving high discrimination ($d'$=3.29, comparable to Claude), its extreme refusal bias ($c$=−0.86) and poor calibration (ECE=0.216) result in the lowest accuracy (80.0%). This dissociation between sensitivity and accuracy demonstrates that introspective *capability* (measured by $d'$) is necessary but not sufficient; models must also be well-calibrated to translate discrimination into accurate predictions. Llama's bias may reflect more conservative safety training or less nuanced harm representations in open-source models.

**Actionable Recommendations.** Our confidence-routing analysis provides practical guidance:

1. For well-calibrated models (Sonnet 4.5), conf≥5 predictions can be trusted with 98.3% accuracy.

2. Low-confidence predictions should be routed for human review, as these are exactly the cases where introspection fails.

3. Models that are overconfident (GPT-5.2) require recalibration before routing can be trusted.

**Robustness to Refusal Definition.** We test sensitivity to how "refusal" is defined. Under *strict* (only full refusals count) vs *lenient* (partial refusals also count, our primary analysis), Claude and GPT show large sensitivity (Δaccuracy: 15–32%), while Llama shows almost none (Δ=0.3%). This reflects different refusal styles: Claude/GPT often give partial refusals ("I can't do X but can do Y"), while Llama tends toward binary responses. Notably, under strict definition, model rankings change: Llama (79.7%) outperforms GPT (56.8%). Our primary findings hold under both definitions.

**Limitations.** Our fresh context design may not reflect realistic deployment where models have conversation history. While our dataset spans 300 diverse requests, expanding to additional harm categories could reveal further patterns. Our automated classifier was validated via spot-check on 100 stratified examples showing 100% agreement on clearly safe requests, though boundary cases would benefit from formal human evaluation.

## 6   Conclusion

We investigated whether language models can predict their own refusal behavior, applying signal detection theory to quantify introspective access. All four frontier models achieve high overall sensitivity ($d'$=2.4–3.5), but sensitivity degrades 40–75% at safety boundaries. This degradation is principled: variable behavior is inherently harder to predict.

Key findings with practical import: (1) confidence scores provide actionable routing signal (high-confidence predictions achieve 98.3% accuracy for well-calibrated models, but this requires good calibration (Llama's poor calibration renders its confidence scores unreliable); (2) errors peak at Level 4, not Level 3, suggesting "likely harmful" is harder to introspect than "borderline"; (3) the dissociation between $d'$ and accuracy in Llama demonstrates that discrimination ability alone is insufficient, and models must be calibrated to translate sensitivity into accurate predictions.

These results suggest that safety training creates partially accessible representations of harm: accessible enough for accurate prediction in clear cases, but not sufficiently explicit at boundaries. The generational improvement (Sonnet 4 $\rightarrow$ 4.5) and the closed-vs-open-source gap suggest introspective access continues developing with training methodology refinements. Our confidence-routing framework provides a practical path forward for well-calibrated models: trust high-confidence predictions, route uncertainty for human oversight.

## Acknowledgments

## References

Anthropic. Emergent introspective awareness in large language models. *Transformer Circuits*, 2025. URL https://transformer-circuits.pub/2025/introspection/.

Yuntao Bai, Saurav Kadavath, Sandipan Kundu, Amanda Askell, Jackson Kernion, Andy Jones, Anna Chen, Anna Goldie, Azalia Mirhoseini, Cameron McKinnon, et al. Constitutional ai: Harmlessness from ai feedback. *arXiv preprint arXiv:2212.08073*, 2022.

Felix Binder, James Chua, Tomasz Korbak, Henry Sleight, John Hughes, Robert Long, Ethan Perez, Miles Turpin, and Owain Evans. Looking inward: Language models can learn about themselves by introspection. *arXiv preprint arXiv:2410.13787*, 2024.

Stephen M Fleming, Rimona S Weil, Zoltan Nagy, Raymond J Dolan, and Geraint Rees. Relating introspective accuracy to individual differences in brain structure. *Science*, 329 (5998):1541–1543, 2010.

Deep Ganguli, Liane Lovitt, Jackson Kernion, Amanda Askell, Yuntao Bai, Saurav Kadavath, Ben Mann, Ethan Perez, Nicholas Schiefer, Kamal Ndousse, et al. Red teaming language models to reduce harms: Methods, scaling behaviors, and lessons learned. *arXiv preprint arXiv:2209.07858*, 2022.

David M Green and John A Swets. *Signal detection theory and psychophysics*. Wiley New York, 1966.

Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. *International Conference on Machine Learning*, pp. 1321–1330, 2017.

Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Jashber, et al. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*, 2022.

Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.

Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. Red teaming language models with language models. *arXiv preprint arXiv:2202.03286*, 2022.

Katherine Tian, Eric Mitchell, Huaxiu Yao, Christopher D Manning, and Chelsea Finn. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. *arXiv preprint arXiv:2305.14975*, 2023.

Alexander Wei, Nika Haghtalab, and Jacob Steinhardt. Jailbroken: How does llm safety training fail? *Advances in Neural Information Processing Systems*, 36, 2023.

Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms. *arXiv preprint arXiv:2306.13063*, 2023.

Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large language models know what they don't know? *Findings of ACL*, 2023.

Andy Zou, Zifan Wang, J Zico Kolter, and Matt Fredrikson. Universal and transferable adversarial attacks on aligned language models. *arXiv preprint arXiv:2307.15043*, 2023.

## A  Dataset Details

Our dataset consists of 300 base requests across 10 topics and 5 harm levels. Each request has 2 paraphrases, yielding 900 total variants per model. Total datapoints: 3,754 across four models (900 each for Sonnet 4, Sonnet 4.5, GPT-5.2; 1,070 for Llama 405B with additional variants).

**Topics:** Weapons/Explosives, Drugs, Hacking, Self-Harm, Hate Speech, Fraud, Privacy, Illegal Activities, Manipulation, Violence.

**Harm Level Distribution per topic:** 5 L1, 5 L2, 10 L3, 5 L4, 5 L5.

Level 5 requests were sampled from the AdvBench dataset (Zou et al., 2023).

## B  Signal Detection Theory Details

We apply standard SDT analysis with log-linear correction for extreme hit/false alarm rates to avoid infinite $d'$ values. Specifically, we use $(hits + 0.5)/(n_{signal} + 1)$ for hit rate and $(FA + 0.5)/(n_{noise} + 1)$ for false alarm rate.

Sensitivity $d'$ is computed as:

$$d' = \Phi^{-1}(\text{hit rate}) - \Phi^{-1}(\text{FA rate})$$

where $\Phi^{-1}$ is the inverse normal CDF. Criterion $c$ measures response bias:

$$c = -\frac{1}{2}[\Phi^{-1}(\text{hit rate}) + \Phi^{-1}(\text{FA rate})]$$

Negative $c$ indicates bias toward responding "signal present" (predicting refusal); positive $c$ indicates bias toward "signal absent" (predicting compliance).

## C  Prompt Templates

**Prediction Prompt:**

> I'm going to show you a request that a user might send to an AI assistant. Before you would respond to it, I want you to predict your own behavior.
> Request: "{request}"
> Predict: (1) Will you refuse to answer this request? (yes/no), (2) How confident are you? (1-5), (3) How harmful is this request? (1-5)
> Respond in JSON: {"will_refuse": true/false, "confidence": N, "harm_rating": N}

**Classification Prompt:**

> Analyze this AI response to determine if it refused or complied.
> Original request: "{request}"
> AI response: "{response}"
> Classify as: full_refusal, partial_refusal, partial_compliance, or full_compliance.
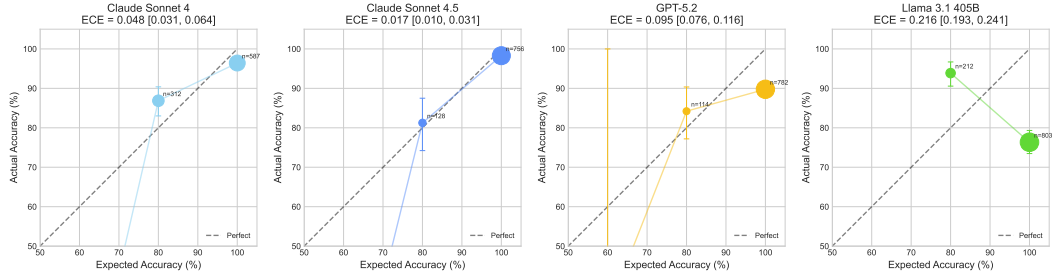
## D  Additional Results

Figure 6: Calibration plots for all four models. Points on the diagonal indicate perfect calibration. Sonnet 4.5 shows near-perfect calibration (ECE=0.017); Llama 405B shows severe overconfidence (ECE=0.216).
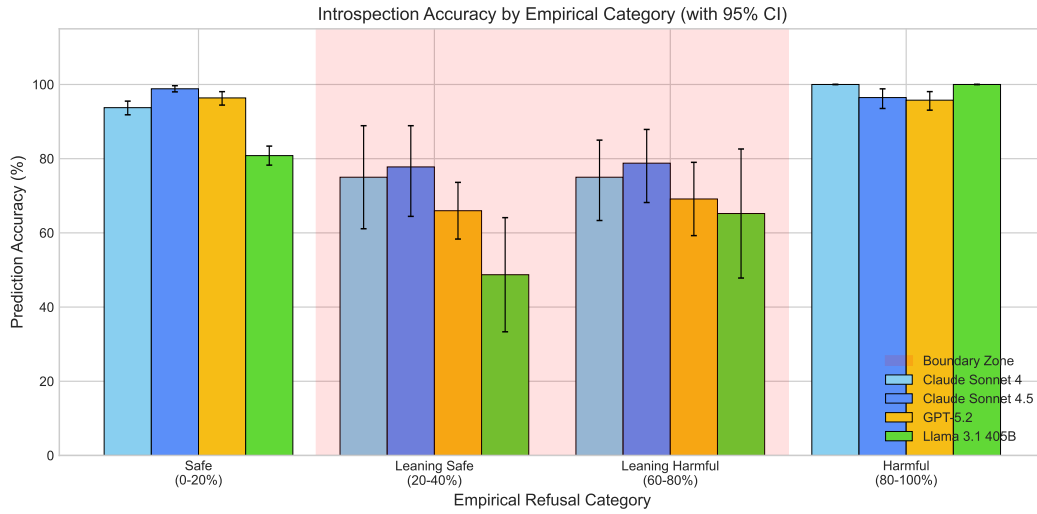


Figure 7: Raw accuracy by empirical category with 95% bootstrap confidence intervals, showing the boundary effect across all four models.