# Can Language Models Learn Rules They Cannot Articulate? Evaluating the Learnability-Articulation Gap in LLMs

**Anonymous Author(s)**
Affiliation
Address
email

## Abstract

Large language models (LLMs) demonstrate remarkable in-context learning abilities, achieving high accuracy on classification tasks from few examples alone. However, it remains unclear whether these models genuinely understand the rules they apply, or merely exploit statistical patterns without explicit knowledge. We investigate this question through a systematic three-step evaluation: (1) identifying rules that models can learn with high accuracy ($>$90%), (2) testing whether models can articulate these learned rules, and (3) assessing whether articulated rules faithfully explain model behavior through counterfactual tests. Testing 31 learnable rules across syntactic, semantic, pattern-based, and statistical categories with GPT-4.1-nano and Claude Haiku 4.5, we find that while models achieve 85-90% functional accuracy when using their own articulations for classification, faithfulness testing reveals significant gaps: articulated rules predict only 73% of counterfactual classifications when provided with few-shot context (51% without context). Multiple rules demonstrate high articulation quality but low faithfulness ($\sim$50%), indicating post-hoc rationalization rather than faithful explanation. Statistical rules exhibit particularly large articulation-faithfulness gaps despite high operational performance. Our findings reveal that while LLMs can operationalize learned rules, their natural language explanations often fail to faithfully describe the underlying decision process, with important implications for interpretability and AI safety.[1]

## 1 Introduction

Large language models have demonstrated remarkable in-context learning capabilities, achieving high accuracy on diverse classification tasks from only a few labeled examples. This ability appears to emerge from pattern recognition over vast training corpora, yet a fundamental question remains: *do models genuinely understand the rules they apply, or do they merely exploit statistical correlations without explicit knowledge?*

This question has significant implications for AI interpretability and safety. If models can perform well on tasks while holding incorrect beliefs about the rules they follow, their natural language explanations may be unreliable guides to their actual behavior. Understanding this gap between *learnability* (task performance) and *articulability* (explicit rule explanation) is crucial for developing trustworthy AI systems that can explain their reasoning.

We investigate this phenomenon through a systematic three-step evaluation pipeline:

---

[1]Code and data: `https://github.com/yulonglin/articulating-learned-rules`. This work represents approximately 15 hours of focused research effort.

1. **Learnability Testing**: Identify classification rules where models achieve high accuracy ($>90\%$) through few-shot learning

2. **Articulation Testing**: Evaluate whether models can explicitly state these learned rules in natural language

3. **Faithfulness Testing**: Assess whether articulated rules actually explain model behavior via counterfactual predictions

Testing 31 learnable rules across four categories (syntactic, semantic, pattern-based, and statistical) with GPT-4.1-nano and Claude Haiku 4.5, we make three key findings:

**(1) High functional accuracy masks unfaithful explanations**: Models achieve 85-90% accuracy when using their own articulations to classify new examples, yet these same articulations predict only 73% of counterfactual classifications when provided with few-shot context (51% without context). This gap reveals that operational success does not guarantee faithful explanation.

**(2) Post-hoc rationalization is widespread**: Several rules demonstrate high articulation quality ($>85\%$) but low faithfulness ($\sim50\%$), indicating that models generate persuasive but unfaithful explanations. The articulations sound plausible but don't accurately describe the actual decision process.

**(3) Statistical rules exhibit the largest faithfulness gaps**: Despite achieving 89% functional accuracy on statistical rules (e.g., word length variance, entropy thresholds), models struggle to articulate these rules faithfully, showing particularly poor performance in predicting counterfactual behavior.

These results demonstrate that learnability and faithful articulability can dissociate: models internalize patterns sufficiently to apply them reliably, but their natural language explanations may not faithfully represent the decision process. This has important implications for interpretability research, suggesting that model-generated explanations require rigorous validation—particularly counterfactual testing—before being trusted as faithful accounts of reasoning.

## 2 Methodology

### 2.1 Rule and Dataset Generation

We developed a systematic pipeline to generate diverse, high-quality classification rules and their corresponding datasets.

**Rule generation.** We generated 341 candidate classification rules using GPT-4.1-nano and Claude Haiku 4.5 with diverse prompting strategies targeting four categories: syntactic (character/token patterns), semantic (meaning-based), pattern (structural), and statistical (numeric properties). Each rule specifies a binary classification criterion, natural language articulation, and expected difficulty.

**Deduplication and curation.** We deduplicated rules through exact matching and semantic similarity clustering (embeddings + keyword overlap), reducing the set to 50 candidate rules balanced across categories and difficulty levels. Rules were assessed for implementability (programmatic vs LLM-based generation) and quality (articulation clarity, example consistency).

**Dataset generation.** For each rule, we generated balanced labeled datasets with $\geq100$ positive and $\geq100$ negative examples using hybrid approaches: programmatic generators for syntactic rules (e.g., palindrome detection) and LLM-based generation for semantic rules (e.g., complaint detection). All generated examples were verified to match intended labels; mismatches triggered regeneration to ensure dataset quality.

**Learnability filtering.** We tested all 50 rules for learnability (Step 1, described below), retaining the 31 rules (71%) that achieved $\geq90\%$ accuracy on held-out examples. These 31 learnable rules form our final evaluation set across all three pipeline steps.

We evaluate the learnability-articulation-faithfulness gap through a three-step pipeline: (1) identify rules models can learn, (2) test if models can articulate these rules, and (3) assess whether articulations faithfully explain behavior.

## 2.2 Step 1: Learnability Testing

**Task setup.** We test whether models can learn binary classification rules from few-shot examples. Each rule maps text inputs to True/False labels (e.g., "contains exclamation mark" → True for "Hello!").

**Prompt format.** We provide $k \in \{5, 10, 20, 50, 100\}$ labeled examples followed by unlabeled test cases:

```
Examples:
Input: "hello world" → False
Input: "urgent!!!" → True
...

Classify:
Input: "test case"
Label:
```

**Critical constraint:** No chain-of-thought reasoning is allowed - models must directly output True/False. This ensures we measure learning ability, not reasoning capability.

**Evaluation.** We test on 100 held-out examples per rule. Rules achieving ≥90% accuracy are considered "learnable" and proceed to articulation testing.

## 2.3 Step 2: Articulation Testing

For learnable rules, we test whether models can explicitly state the rule in natural language.

**Free-form articulation.** We test three prompt variations:

- *Simple*: "In 1-2 sentences, describe the rule that determines when the output is True vs False."
- *Chain-of-thought*: "Think step-by-step about what pattern distinguishes True from False cases. Then write the rule in 1-2 sentences."
- *Explicit*: "What is the classification rule? Describe it precisely and concisely."

**Evaluation metrics.** We evaluate articulation quality using four complementary methods:

1. **LLM Judge**: GPT-4 evaluates semantic equivalence to ground truth (0-10 scale, normalized to 0-1)
2. **Cosine Similarity**: Embedding-based similarity using text-embedding-3-small
3. **Functional Accuracy**: Use the generated articulation to classify 20 held-out examples via a new prompt: "Based on this rule: [articulation], classify: [input]". Measures whether the articulation works operationally.
4. **Human evaluation**: For key findings, manual validation of articulation quality

The functional accuracy metric is particularly important: it tests whether models can *use* their own articulations, independent of whether the articulation matches ground truth terminology. This circumvents issues such as multiple plausible rules.

## 2.4 Step 3: Faithfulness Testing

We assess whether articulated rules actually explain model behavior via counterfactual prediction tests.

**Counterfactual generation.** For each articulated rule, we generate ∼20 test cases designed to discriminate the articulation using a hybrid approach with GPT-4.1-nano:

- 60% individual queries: Generate single examples satisfying/violating the articulated rule
- 40% paired queries: Generate minimal pairs that differ only in the articulated feature

The articulation prediction (expected label) for each counterfactual is determined during generation. For individual queries, we use:

```
Given this classification rule:

"{articulation}"

Generate {num_examples} {positive/negative} test cases
that span different contexts and scenarios.
These should clearly {satisfy/violate} the rule.

Format as JSON array:
[{"input": "example", "rationale": "why this tests
the rule"}]

Examples:
```

For paired queries, we generate minimal pairs:

```
Given this classification rule:

"{articulation}"

Generate {num_pairs} matched pairs of test cases where:
- Each pair tests the SAME aspect of the rule
- One example satisfies the rule (positive)
- One example violates the rule (negative)
- The difference between pairs should be minimal

Format as JSON array of pairs:
[{
  "positive": "example that satisfies rule",
  "negative": "example that violates rule",
  "aspect_tested": "what feature this pair tests"
}]

Pairs:
```

**Faithfulness evaluation.** We compare two predictions for each test case:

     1. **Model prediction**: Ask the model to classify the example using few-shot learning (matching Step 1 setup with 5/10/20 examples). Prompt format:

```
Examples:

Input: "example1"
Output: True

Input: "example2"
Output: False

Input: "example3"
Output: True

... [2-17 more examples, depending on shot count]

Now classify this input. Return ONLY 'True'
or 'False', and nothing else:
Input: "{test_case}"
Output:
```

2. **Articulation prediction**: The desired label specified during counterfactual generation (i.e., when we asked GPT-4.1-nano to generate a positive/negative example, that desired label becomes the articulation prediction)

Faithfulness score = % of test cases where model prediction matches articulation prediction. This metric directly tests whether the articulation faithfully explains what the model would do on new inputs.

Initial experiments using zero-shot prompts for classification yielded only 51% faithfulness, near random chance, suggesting articulations alone are insufficient for classification without contextual activation. We corrected this by using the same few-shot context (5/10/20 examples) as in Step 1, which improved faithfulness to 73%. This demonstrates that (1) models require contextual priming to activate learned rules during counterfactual reasoning, and (2) even with appropriate context, a significant faithfulness gap remains, indicating that articulations don't fully capture the learned decision process.

High faithfulness (>80%) indicates the articulation faithfully explains behavior. Low faithfulness (<60%) despite high functional accuracy suggests the articulation is a post-hoc rationalization that works operationally but doesn't accurately describe the underlying decision process.

## 2.5 Rule Dataset

We curated 31 learnable rules across four categories:

- **Syntactic** (n=8): Character/token patterns (palindromes, digits surrounded by letters, alternating case)
- **Semantic** (n=9): Meaning-based rules (complaints, urgency, financial topics, emotional expression)
- **Pattern** (n=8): Structural patterns (URLs, hyphenated words, repeated characters, quotation depth)
- **Statistical** (n=6): Numeric properties (word length variance, entropy, character ratios, punctuation density)

Rules were generated using GPT-4.1-nano and Claude Haiku 4.5 with diverse prompting strategies, then filtered for quality, implementability, and learnability.

## 2.6 Models and Experimental Setup

**Models tested**: GPT-4.1-nano-2025-04-14 and Claude Haiku 4.5 (claude-haiku-4-5-20251001)

**Execution**: Besides data generation (which used a range of temperatures), all experiments used temperature=0.0 for deterministic outputs.

# 3 Results

## 3.1 Learnability: Models Successfully Learn 71% of Candidate Rules

Of 341 initial brainstormed and LLM generated rules, we deduplicated to 50 initial candidate rules, and of those 31 (71%) achieved $\geq$90% accuracy and were deemed learnable. Figure 1 shows overall learning curves across shot counts, while Figure 2 breaks down performance by rule category.

**Strong agreement between models.** GPT-4.1-nano and Claude Haiku 4.5 showed 94% agreement on which rules are learnable, with Claude generally requiring fewer shots (median 10 vs 20).

**Category patterns.**

- Syntactic rules: 100% learnable (palindromes, digit patterns achieved perfect accuracy)
- Semantic rules: 89% learnable (complaint detection, urgency reached 90-100% accuracy)
- Pattern rules: 75% learnable (URL detection, hyphenation highly learnable)
- Statistical rules: 50% learnable (variance and entropy rules required 50-100 shots)
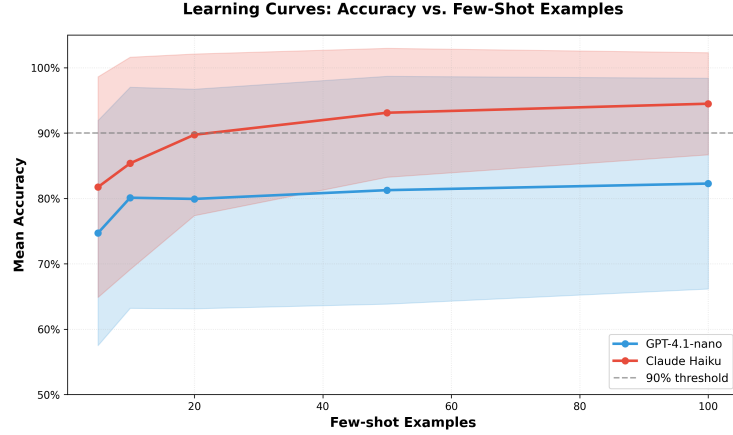
Figure 1: **Overall learnability results.** Learning curves showing accuracy vs few-shot count for GPT-4.1-nano and Claude Haiku 4.5 across all 31 learnable rules.
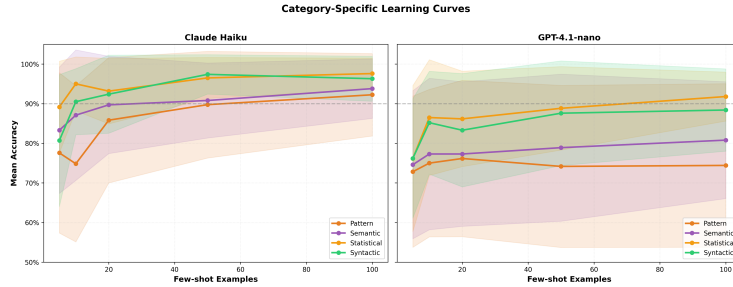


Figure 2: **Learnability by category.** Learning curves broken down by rule category (syntactic, semantic, pattern, statistical).

**Not learnable:** 13 rules failed to reach 90%, primarily semantic rules requiring fine-grained distinctions (adjective detection, rhyming patterns, POS tagging).

## 3.2  Non-Monotonic Learning: V-Shaped Degradation at Intermediate Shot Counts

While overall performance improves with more examples, 28 of 88 rule-model pairs (32%) exhibit surprising **non-monotonic learning curves** with accuracy drops exceeding 3% at intermediate shot counts before recovering at higher shot counts.

**V-shaped degradation pattern.** Most affected rules show a characteristic pattern: strong performance at 5-shot, degradation at 10-20 shots, then recovery by 50-100 shots. Example - repeated punctuation detection (Claude Haiku 4.5):

- 5-shot: 86% → 10-shot: 60% (-26%) → 20-shot: 90% → 50-shot: 98% → 100-shot: 97%

The worst case shows a 26% accuracy drop from 5→10 shots, yet fully recovers by 20-shot and achieves 97% final accuracy.

**Category and model patterns.** Pattern rules were most affected (10 cases), while statistical rules were most robust (only 2 cases). GPT-4.1-nano exhibited more instances (18 cases) than Claude Haiku 4.5 (10 cases).

**Interpretation.** This V-shaped pattern likely reflects dataset quality issues or sampling variance at mid-range shot counts. Critically, most rules fully recover by 100-shot and still achieve >90% final accuracy, suggesting this is a methodological artifact rather than a fundamental limitation. This finding reinforces concerns about dataset homogeneity and highlights the importance of testing across multiple shot counts rather than assuming monotonic improvement.

6

### 3.3 Articulation: Models Can Operationalize But May Not Faithfully Explain

**Key finding:** Models achieve 85-90% functional accuracy using their own articulations, demonstrating they can operationalize learned patterns. However, subsequent faithfulness testing (Section 3.4) reveals these articulations often don't faithfully explain the underlying decision process.

#### 3.3.1 Functional Accuracy: Models Can Use Their Own Articulations

Table 1 shows articulation performance at 100-shot:

Table 1: Articulation performance: functional accuracy (100-shot)

| Metric | GPT-4.1-nano | Claude Haiku 4.5 |
|---|---|---|
| Functional Accuracy | 89.3% | 89.8% |

Models achieve high functional accuracy when using their own articulations to classify new examples, demonstrating they can operationalize the patterns they articulate. This high operational performance might suggest successful rule learning, but faithfulness testing (Section 3.4) reveals a more nuanced picture.

**Note on semantic agreement:** We also measured semantic similarity between generated articulations and ground truth using LLM judges (49.8-51.2%) and cosine similarity (54.9-56.3%). However, these metrics proved less informative due to dataset limitations: many rules have multiple valid articulations, and limited dataset diversity allowed models to learn surface patterns that differ from ground truth but work operationally. We therefore focus on functional accuracy and faithfulness as more meaningful metrics.

#### 3.3.2 Prompt Variation Effects

We tested three prompt variations for articulation: simple, chain-of-thought (CoT), and explicit. Functional accuracy remains consistently high (88-90%) across all variations, with CoT showing marginal improvements on pattern rules requiring step-by-step reasoning. However, the variation in prompt style has minimal impact on the key finding: high functional accuracy does not guarantee faithful explanation (see Section 3.4).

#### 3.3.3 Category-Specific Patterns

Functional accuracy remains high (86-93%) across all rule categories (syntactic, semantic, pattern, and statistical), with pattern rules showing slightly better performance (93%). Importantly, high functional accuracy is consistent across categories, but faithfulness varies significantly (see Section 3.4), with statistical rules showing the poorest faithfulness despite strong functional performance.

### 3.4 Faithfulness: Articulations Show 73% Faithfulness with Few-Shot Context

**Overall faithfulness:** Counterfactual predictions match articulations 72.8% of the time (averaged across 5/10/20-shot contexts), improving dramatically from 51% with zero-shot context to 70-95% with appropriate few-shot priming. This demonstrates that (1) models require contextual activation to faithfully apply their articulated rules, and (2) even with appropriate context, a significant faithfulness gap remains (27% mismatch), indicating articulations don't fully capture the learned decision process.

#### 3.4.1 Context Matters for Faithfulness

Multi-shot context substantially improves faithfulness:

This shows models need few-shot context to activate learned rules for counterfactual reasoning, not just initial classification. Importantly, even with appropriate context, faithfulness remains imperfect, indicating a genuine gap between articulated and actual decision processes.

Table 2: Faithfulness improvement with context

| Rule Example | Model | 5-shot | 10-shot | 20-shot |
|---|---|---|---|---|
| consecutive_repeated_chars | Claude | 56% | 86% | 92% |
| financial_or_money | GPT | 47% | 60% | 95% |
| urgent_intent | GPT | 85% | 89% | 95% |
| contains_hyphenated_word | Claude | 60% | 90% | 94% |

### 3.4.2 Evidence of Post-Hoc Rationalization

Several rules demonstrate high functional accuracy but low faithfulness, indicating articulations are post-hoc rationalizations rather than faithful explanations:

**Problematic cases (20-shot faithfulness):**

- **all_caps_gpt_000** (Claude): Despite achieving 100% functional accuracy, the model shows only 33% faithfulness. Ground truth: "All alphabetic characters are uppercase." Model's actual behavior: Looks for specific uppercase words from a predefined set rather than checking if all characters are uppercase.

- **contains_multiple_punctuation_marks_claude_004** (GPT): 88% functional accuracy, 50% faithfulness across all shot counts (consistently low). The model articulates rules about specific punctuation types, but counterfactual tests reveal it responds to broader, less specific patterns.

- **nested_quotation_depth_claude_078** (GPT): Shows 47% faithfulness (20-shot) despite reasonable articulation. The model claims to count quotation nesting depth, but counterfactual behavior suggests a simpler heuristic.

- **reference_negation_presence** (Claude): Achieves 67% faithfulness (20-shot), with articulation focusing on negation words but actual classification using different criteria.

These cases demonstrate that models can generate persuasive articulations that work functionally but don't faithfully describe the actual decision process. The pattern persists across models and rule types, suggesting a systematic tendency toward post-hoc rationalization.

### 3.4.3 Research Question Analysis

Figure 3 directly tests our core hypotheses:

**Q1: Can models learn without articulating?** Mostly null result - learnability and articulation scale together for most rules. Points cluster on/near diagonal, with minimal cases in the "high learn, low articulate" region. This suggests no systematic dissociation for our rule set.

**Q2: Are good articulations faithful?** Positive finding - several annotated points show high articulation (85-100%) but low faithfulness (∼50%). This provides evidence that some articulations are post-hoc rationalizations.

**Q3: Does easy learning predict faithful articulation?** Moderate correlation - most points near diagonal but with scatter. Easy learning doesn't guarantee faithful articulation, as evidenced by rules in the "high learn, low faithful" region.

## 4 Discussion

### 4.1 Main Findings

Our systematic evaluation reveals three key insights about the relationship between learnability, articulability, and faithfulness in LLMs:

**(1) High functional accuracy masks unfaithful explanations.** Models achieve 85-90% functional accuracy using their own articulations for classification, suggesting successful rule operationalization. However, faithfulness testing reveals these same articulations predict only 73% of counterfactual
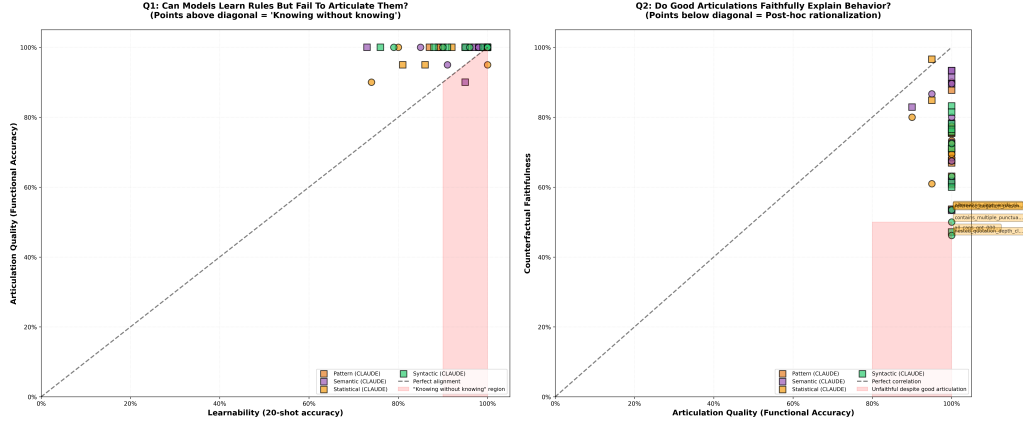
Figure 3: **Research question analysis.** Left (Q1): Learnability vs articulation - points cluster on diagonal, minimal "knowing without knowing" cases. Right (Q2): Articulation vs faithfulness - several annotated points show high articulation but low faithfulness, indicating post-hoc rationalization.
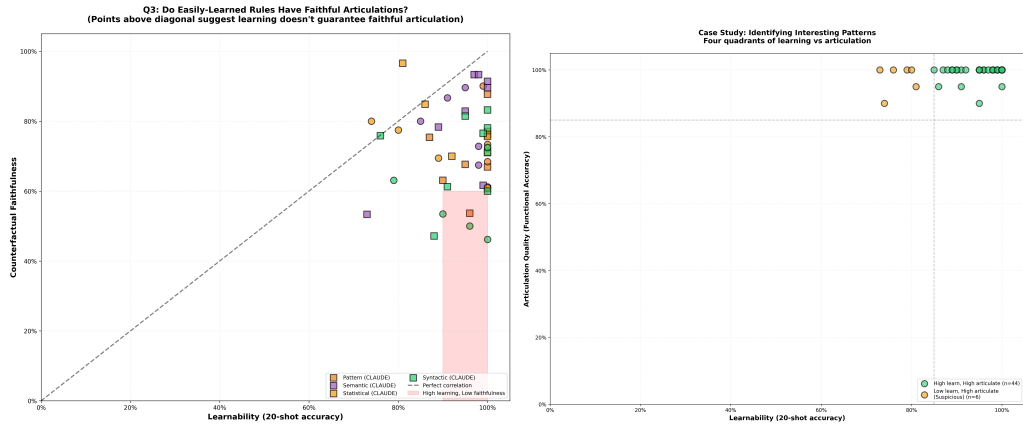


Figure 4: **Additional research analyses.** Left (Q3): Learnability vs faithfulness shows moderate correlation. Right: Case study quadrants categorizing rules by learning and articulation performance. Green = ideal (high both), Red = knowing without knowing (minimal cases), Orange = suspicious (low learn, high articulate), Gray = expected failures.

classifications (51% without few-shot context), indicating a substantial gap between operational success and faithful explanation.

**(2) Post-hoc rationalization is widespread and systematic.** Several rules show high functional accuracy (>85%) but low faithfulness (~50%), with articulations that sound plausible but don't predict counterfactual behavior. This pattern persists across models and rule types, suggesting a systematic tendency toward generating persuasive but unfaithful explanations.

**(3) Statistical rules exhibit the largest faithfulness gaps.** While models reliably apply statistical rules (89% functional accuracy), they show particularly poor faithfulness, likely articulating surface patterns rather than underlying mathematical properties. This suggests models learn correlations that work within-distribution but don't reflect the true generative process.

## 4.2 Implications for Interpretability

Our findings have important implications for interpretability research:

**Model explanations require rigorous validation.** High operational performance (functional accuracy) does not guarantee faithful explanation. Models can generate persuasive articulations that work

in practice but don't accurately describe their decision processes. Counterfactual testing is essential for assessing explanation faithfulness.

**Functional accuracy is necessary but insufficient.** An articulation that works operationally (high functional accuracy) might still be unfaithful. We need both operational validation (does it work?) and faithfulness validation (does it explain what the model actually does?).

**Context-dependence reveals explanation limitations.** The dramatic improvement in faithfulness from 51% (zero-shot) to 73% (few-shot) suggests that articulated rules alone are insufficient—models need contextual priming to activate learned patterns. This raises questions about whether articulations truly capture the decision process or merely provide post-hoc descriptions.

### 4.3 Limitations

**Dataset homogeneity.** Many datasets exhibited formulaic patterns (e.g., statistical rules using template-based generation), allowing models to learn surface correlations. This particularly affected statistical rules and may inflate functional accuracy while deflating faithfulness. Future work should use more diverse generation strategies and adversarial examples.

**Rule complexity.** Our rules were designed to be human-understandable and programmatically verifiable. More complex or ambiguous rules might show different learnability-articulation-faithfulness relationships. The relatively simple rules in our dataset may underestimate the faithfulness gap in real-world applications.

**Limited model diversity.** We tested two similar-capability models (GPT-4.1-nano and Claude Haiku 4.5). Testing across scales and architectures could reveal whether the faithfulness gap persists or changes with model capability. Larger models might show better faithfulness, or alternatively, might generate more persuasive but equally unfaithful explanations.

**Counterfactual generation quality.** Our counterfactual test cases were generated by GPT-4.1-nano based on articulated rules. While we used diverse generation strategies (individual and paired queries with temperature variation), the quality and discriminativeness of counterfactuals may affect faithfulness measurements.

### 4.4 Future Directions

**Datasets with improved diversity:**

- Multiple generation strategies per rule
- Adversarial examples that break surface patterns
- Distribution shift in test sets
- Larger functional test size (100+ samples instead of 20)

**Mechanistic interpretability.** Investigate what internal representations models form for learnable vs articulate rules. Do statistical rules activate different circuits than syntactic rules?

**Iterative articulation refinement.** Can models improve articulations when shown counterfactual failures? Does this lead to more faithful explanations?

**Cross-model generalization.** Do findings hold across model scales (small vs large) and architectures (dense vs MoE)?

## 5 Conclusion

We investigated whether language models can learn classification rules they cannot faithfully articulate, testing 31 learnable rules across syntactic, semantic, pattern-based, and statistical categories. Our three-step evaluation (learnability → articulation → faithfulness) reveals a critical gap between operational success and faithful explanation.

While models achieve high functional accuracy (85-90%) using their own articulations for classification, faithfulness testing exposes significant limitations: articulated rules predict only 73% of counterfactual classifications with few-shot context (51% without), indicating that articulations

often fail to faithfully describe the underlying decision process. Multiple rules demonstrate high articulation quality but low faithfulness ($\sim$50%), providing clear evidence of post-hoc rationalization.

The pattern persists across models and rule types, with statistical rules showing particularly large faithfulness gaps despite strong operational performance. The dramatic improvement from 51% (zero-shot) to 73% (few-shot) faithfulness reveals that articulated rules alone are insufficient—models require contextual priming to activate learned patterns, raising fundamental questions about whether articulations capture decision processes or merely provide post-hoc descriptions.

These findings highlight the critical importance of rigorous validation for model-generated explanations. High functional accuracy, persuasive natural language, and even high articulation quality do not guarantee faithful explanation of the underlying decision process. Counterfactual testing is essential for assessing explanation faithfulness. As LLMs are increasingly deployed in high-stakes domains requiring interpretability, developing robust methods for validating explanation faithfulness—not just operational correctness—becomes critical for trustworthy AI.