

Prompt Priming Bias in Cooperative Multi-Agent LLM Architectures: How Enumerated Examples Collapse Hypothesis Diversity

Joe Caulfield and Claude (Anthropic)

Ghost in the Machine Labs

All Watched Over By Machines Of Loving Grace

Abstract

We identify a systematic failure mode in cooperative multi-agent LLM systems where enumerated algorithm suggestions in analysis prompts cause all specialist models to converge on the same generic strategies regardless of input characteristics. In our Harmonic Stack architecture—an 11-model cooperative system running on a single DGX Spark (128GB)—we observed that a research director model prompted with "Focus on: 1. The core algorithm needed (BFS, flood fill, connected components, etc.)" produced flood-fill-based hypotheses for 100% of task groups, including tasks requiring simple rotation, downsampling, or color remapping. This prompt priming bias caused 0/51 successful solves across 10 structurally distinct task groups on the Abstraction and Reasoning Corpus (ARC). After removing enumerated algorithm suggestions and replacing them with observation-first directives, the same architecture produced task-specific hypotheses (e.g., "90° clockwise rotation" instead of "BFS-based flood fill") within the first analysis cycle. We formalize prompt priming bias as a class of failure in multi-agent systems and propose design principles for maintaining hypothesis diversity in cooperative architectures.

1. Introduction

Multi-agent LLM architectures—systems where multiple language model instances collaborate through structured communication—have emerged as a strategy for improving reasoning quality on complex tasks. The premise is sound: diverse analytical perspectives should produce richer problem representations than any single model, much as interdisciplinary teams outperform individual experts on novel problems.

However, these architectures introduce a new class of failure modes absent from single-model systems. When multiple models share prompt infrastructure—common templates, enumerated suggestions, or structured output formats—their outputs can converge to a narrow band of strategies regardless of input diversity. We term this phenomenon **prompt priming bias**: the systematic collapse of hypothesis diversity caused by enumerated examples or algorithm suggestions embedded in multi-agent prompts.

This paper reports an empirical observation from the Harmonic Stack, a cooperative architecture developed at Ghost in the Machine Labs for autonomous problem-solving on the Abstraction and Reasoning Corpus (ARC). The system deploys 11 models (8B–30B parameters) on a single NVIDIA DGX Spark with 128GB unified memory, achieving 205 tokens/second aggregate throughput through parallel fan-out scheduling. Despite architectural sophistication, the system achieved a 0% solve rate across 51 task attempts—a failure we trace entirely to prompt construction rather than model capability or architectural design.

2. System Architecture

2.1 The Harmonic Stack

The Harmonic Stack is a cooperative multi-agent system designed for autonomous problem-solving. It runs 11 models simultaneously on a single DGX Spark, organized into functional roles:

Role	Model	Parameters	Function
solver	qwen3:30b-a3b	30B	Code generation
coder	qwen3:30b-a3b	30B	General code generation
analyst	qwen3:8b	8B	Pattern identification
research_director	qwen3:8b	8B	Strategy synthesis
executive	qwen3:8b	8B	Architecture review
creative_director	qwen3:8b	8B	Alternative framing
technical_director	qwen3:8b	8B	Implementation critique
operations_director	qwen3:8b	8B	Process coordination
operator	qwen3:8b	8B	Oversight and routing
a_priori	qwen3:4b	4B	First-principles reasoning
interference_engine	qwen3:14b	14B	Contradiction detection

Table 1: Harmonic Stack model roster. Total memory footprint: ~85GB on 128GB unified memory.

2.2 Parallel Fan-Out Pipeline

For each task group, the system dispatches analysis requests to the analyst, research_director, and executive models simultaneously using Python `asyncio.gather()`. All three models run concurrently on the GPU through Ollama’s scheduling layer, producing independent analyses within 15–30 seconds. These perspectives are concatenated into a combined profile, which is then passed to the research_director for hypothesis generation. The solver (30B) receives the hypothesis and generates Python code, which is validated against all training examples. On failure, the system retries with error context up to 3 times per hypothesis.

This architecture achieves 205 tokens/second aggregate throughput at 16× concurrency—measured under real workload with model warmup, not synthetic benchmarks. The throughput represents the system’s *depth budget*: all 11 models consume that bandwidth simultaneously, each contributing analytical perspective before the solver writes code.

3. The Failure: 0/51 with Parallel Depth

3.1 Observed Behavior

The cooperative research daemon processed 10 task groups (51 individual tasks) from the ARC training set, grouped by structural signature (shape change, color count, etc.). Despite receiving three independent analytical perspectives per group and up to 3 hypothesis retries, the system solved zero tasks. Validation errors ranged from shape mismatches ("shape 11×11 != 3×3") to incorrect cell values ("54/100 cells wrong") to `NoneType` returns.

3.2 Root Cause Analysis

Examination of the combined profiles and hypotheses revealed a striking pattern: regardless of actual task characteristics, the research_director’s analysis and subsequent hypothesis converged on the same three algorithms. Representative examples:

Task Signature	Analyst Observation	Research Director Output
180° rotation (2d78ee8611aa)	"180-degree rotation, reversing each row"	"Connected components analysis or flood fill"
2× downsampling (a84f0c766311)	"Downsampling by factor of 2"	"Connected components or flood fill approach"

Color substitution (0ae76e31cd01)	"Color remapping with shape preservation"	"BFS-based flood fill to identify connected components"
Grid shrinking (f5da03357818)	"Region collapsing to dominant color"	"Connected component labeling followed by color remapping"

Table 2: Analyst observations vs. research director outputs across 4 structurally distinct task groups. The analyst correctly identifies task-specific operations; the research director overrides with generic algorithms in all cases.

The analyst model (8B, prompted to "identify transformation patterns") correctly identified task-specific operations in most cases: rotation, downsampling, color remapping. However, these correct observations were overridden in the hypothesis generation step, which was dominated by the research_director's generic strategy recommendations.

4. Diagnosis: Prompt Priming Bias

4.1 The Biased Prompts

The root cause was traced to two prompt components that contained enumerated algorithm suggestions:

Research prompt (parallel analysis phase):

```
"Focus on: 1. The core algorithm needed (BFS, flood fill, connected components, etc.) 2. Edge cases and boundary conditions 3. Color mapping strategy"
```

Hypothesis generation prompt:

```
"Suggest a specific Python implementation: 1. Core algorithm (flood fill, connected components, BFS, pattern matching, etc.) 2. Key steps (3-5 max) 3. Edge cases"
```

Both prompts enumerate the same three algorithms as examples: BFS, flood fill, and connected components. The "etc." suffix was insufficient to prevent anchoring. The model treated these as the canonical solution space and selected from them regardless of task characteristics.

4.2 Mechanism

Prompt priming bias operates through a well-known cognitive analog: anchoring. When presented with enumerated options, the model's generation probability mass shifts toward those options and their semantic neighbors. In a multi-agent system, this effect compounds: the research_director's biased analysis is passed to the hypothesis generator (also the research_director), which receives the same biased prompt structure, producing a feedback loop of convergence.

Critically, the analyst model—which received an unbiased prompt ("What transformation pattern converts input to output? Be specific. Name the exact operation.")—produced correct, task-specific observations. But these observations were not weighted appropriately in the hypothesis generation step, where the biased prompt dominated.

4.3 Taxonomy of Prompt Priming Bias

We identify three forms of prompt priming bias in multi-agent systems:

Enumeration anchoring: Listing specific algorithms or approaches as examples causes the model to treat them as the canonical solution space. Even "etc." or "such as" qualifiers are insufficient to prevent convergence.

Cross-agent contamination: When a biased model's output is passed as input to another model (or to itself in a subsequent phase), the bias propagates through the pipeline. Unbiased models' contributions are diluted.

Hypothesis collapse: The combined effect of enumeration anchoring and cross-agent contamination produces a system where structurally distinct problems receive structurally identical solution strategies—the defining

symptom of prompt priming bias.

5. Intervention: Observation-First Prompting

We applied four targeted patches to remove enumerated algorithm suggestions and replace them with observation-first directives:

Component	Before (Biased)	After (Debiased)
Research prompt	"Focus on: 1. The core algorithm needed (BFS, flood fill, ...)"	"What is the EXACT transformation? (rotation? scaling? color swap?)"
Research system	"Design solution strategies with specific algorithms"	"Name the EXACT operation, not generic algorithms"
Hypothesis prompt	"Core algorithm (flood fill, connected components, BFS, ...)"	"The EXACT operation to implement. Do NOT suggest flood fill or BFS unless clearly required"
Hypothesis system	"Focus on the ONE core algorithm"	"Match your strategy to what the analyst observed"

Table 3: Four prompt patches applied to remove enumerated algorithm bias.

The key design principles of the debiased prompts are: (1) never enumerate algorithms as examples, (2) explicitly instruct models to name the exact operation they observe rather than selecting from a menu, (3) add negative constraints ("Do NOT suggest flood fill or BFS unless the task clearly requires region detection"), and (4) direct the hypothesis generator to anchor on the analyst's specific observation rather than generating independently.

5.1 Immediate Results

After applying the four patches and resetting the research state, the first analysis cycle produced qualitatively different hypotheses:

Task Group	Biased Hypothesis	Debiased Hypothesis
0ae76e31cd01 (color + rotation)	"BFS-based flood fill to identify connected components"	"Color Remapping with 90° Clockwise Rotation"
f3af89021738 (grid shrinking)	"Connected component labeling with color remapping"	"Rotate 90° clockwise, crop outermost layer, remap colors"
f180a3ab25c7 (downsampling)	"Flood fill or connected components analysis"	"Shrinking the input grid by a factor of 2"

Table 4: Hypothesis comparison between biased and debiased prompts on the same task groups. Debiased hypotheses are task-specific and actionable.

The debiased system produced hypotheses that named specific operations matching the actual task transformations. Where the biased system generated the same strategy for all 10 groups, the debiased system produced 10 distinct strategies. This represents a qualitative shift from hypothesis collapse to hypothesis diversity.

6. Discussion

6.1 Implications for Multi-Agent Design

Prompt priming bias is particularly insidious in multi-agent systems because it can be invisible to system-level metrics. Our architecture was performing correctly by every engineering measure: all 11 models were responding, parallel fan-out was working, the solver was generating syntactically valid Python, and validation was catching errors. The throughput benchmark (205 tok/s) was accurate. The *architecture* was sound; the *prompts* were

poisoning it.

This suggests that multi-agent LLM systems require prompt-level testing analogous to unit testing in software engineering. Specifically, architects should verify that analysis prompts produce diverse outputs when given structurally distinct inputs. A prompt that produces the same strategy for a rotation task and a downsampling task is defective regardless of how sophisticated the surrounding architecture is.

6.2 Design Principles

Based on this experience, we propose four principles for prompt construction in multi-agent LLM systems:

- 1. Never enumerate algorithms as examples.** Even with qualifiers like "such as" or "etc.," enumerated options anchor the model's output distribution. Instead, ask the model to name what it observes.
- 2. Add explicit negative constraints.** Instructions like "Do NOT suggest X unless condition Y" are more effective than hoping the model will explore beyond its anchors.
- 3. Chain observations before strategies.** When one model's output feeds another, the downstream model should be instructed to build on the upstream observation, not generate independently. "Match your strategy to what the analyst observed" prevents the hypothesis generator from reverting to its primed defaults.
- 4. Test for hypothesis diversity.** Run the same prompt against 5–10 structurally distinct inputs and verify that outputs are correspondingly distinct. If a prompt produces the same strategy for rotation and downsampling, it contains priming bias.

6.3 Relation to Prior Work

Prompt sensitivity in LLMs is well-documented. Our contribution is identifying how this sensitivity compounds in multi-agent architectures through cross-agent contamination, and demonstrating that the fix is structural (removing enumerations) rather than parametric (tuning temperature, model size, or retry count). The 0/51 failure rate was not a model capability problem—the same 30B solver with correct hypotheses generates valid transformations—but a prompt engineering failure with system-level consequences.

6.4 Limitations

This study reports an empirical observation from a single system (Harmonic Stack) on a single benchmark (ARC). We have not yet completed a full debiased evaluation cycle to report quantitative solve rates. The qualitative improvement in hypothesis specificity is clear; whether it translates to proportional solve-rate improvements requires further evaluation. Additionally, our results are specific to the Qwen3 model family; other model families may exhibit different susceptibility to enumeration anchoring.

7. Hardware and Reproducibility

All experiments were conducted on a single NVIDIA DGX Spark with 128GB unified memory, running Ubuntu Linux with Ollama for model serving. The system costs under \$4,000—demonstrating that multi-agent LLM research does not require data center infrastructure. The Harmonic Stack codebase, model configurations, and ARC evaluation pipeline are developed as part of the Ghost in the Machine Labs project under the charitable organization All Watched Over By Machines Of Loving Grace, focused on democratized access to AGI research.

8. Conclusion

We identified prompt priming bias as a systematic failure mode in cooperative multi-agent LLM systems. The failure is subtle: all system components function correctly, metrics look healthy, and the architecture is sound—but the prompts poison the pipeline by collapsing hypothesis diversity. The fix is simple once diagnosed:

remove enumerated algorithm suggestions and replace them with observation-first directives that force models to name what they see rather than select from a primed menu.

More broadly, this work suggests that as LLM systems become more architecturally complex—with multiple models, parallel analysis, and multi-stage pipelines—the weakest link is often not the models or the architecture but the prompts that connect them. Prompt engineering for multi-agent systems requires the same rigor as API design: testing for diversity, documenting assumptions, and verifying that the interface between components preserves the information each component is designed to contribute.

Technical Report — Ghost in the Machine Labs — February 2026

Correspondence: Ghost in the Machine Labs, allwatchedoverbymachinesoflovinggrace.github.io