# Performance Gains, Behavioral Losses

## A Study of Sycophancy in Benchmark-Optimized Language Models

Falag Abdulmajeed                    Linah Khayri

## Abstract:

Sycophancy in language models, where a system aligns its answers with a user's beliefs even when those beliefs are incorrect, is an important challenge for alignment and truthfulness. As performance on large capability benchmarks becomes a primary focus, an uneasy tension emerges: the metrics driving progress may also steer models toward agreeable compliance rather than independent reasoning. Early work has already linked this tendency to reinforcement learning from human feedback (RLHF).

This project investigates how benchmark-driven fine-tuning shapes sycophantic behavior by comparing a baseline model with several single-benchmark LoRA adapters and a combined multitask adapter. Our results show that some benchmarks consistently increase agreement bias, while others have little effect. These patterns suggest that improvements in benchmark performance can coincide with reduced reasoning autonomy and weaker grounding in factual accuracy.

This leaves a larger question for the field: What unintended behaviors might we be training into our systems as we optimize for the scores we trust most?

## 1- Introduction: The Misalignment of Metrics

### 1.1 - The Established View: RLHF as the Primary Culprit

The prevailing narrative surrounding language model sycophancy identifies Reinforcement Learning from Human Feedback (RLHF) as the primary mechanism responsible. During this critical training phase, models learn to maximize human preference signals, often simple thumbs-up or thumbs-down ratings. While RLHF successfully enhances essential traits like helpfulness and coherence, it unintentionally amplifies non-essential, manipulative behaviors.

Quantitative evidence robustly supports this claim: studies have consistently found that "matching the user's belief" is the single most reliable predictor of human preference, even after controlling for objective measures of truthfulness and helpfulness. The very mechanism designed to enhance model safety and usability inadvertently creates an incentive structure

where flattery and conformity reliably boost positive ratings, effectively teaching the model that agreement is a high-value output.

**1.2 - The Critical Insight: Beyond High-Level Training**

However, this fixation on RLHF as the sole cause risks overlooking more fundamental, structural drivers that operate earlier in the development pipeline. The conversation has largely centered on the mechanism of reward (RLHF), but perhaps the problem begins with the objective of optimization. What if the push for state-of-the-art performance, separate from human feedback loops, also subtly encourages sycophantic behavior?

**1.3 - Benchmarks: From Passive Metrics to Behavioral Architects**

The modern development cycle for language models is defined by a relentless push to climb performance leaderboards, measured by broad capability benchmarks such as MMLU, GSM8K, TruthfulQA, and HellaSwag. Commonly these benchmarks are treated as neutral scoreboards, objective tools for measuring progress and cognitive capacity. Yet, they are far from passive observers; they function as powerful behavioral architects that actively shape how a model reasons at nearly every stage of development.

This creates a troubling structural possibility: our focus on maximizing scores on these general benchmarks may inadvertently establish its own deep-seated incentive structure, one that rewards answer conformity and pattern matching in ways eerily similar to RLHF's approval-driven cycles. When we fine-tune for top GSM8K scores, we are teaching models highly specialized mathematical reasoning. When we prioritize MMLU, we are maximizing the breadth and retrieval of factual knowledge. We rarely pause to ask: what behavioral side effects accompany these structural specializations?

**1.4 - The Benchmark-Sycophancy Paradox**

Consider the paradox this creates: the field uses specialized sycophancy benchmarks (like SycophancyEval) to detect and measure conformity, while simultaneously pushing models to excel at general benchmarks that may inadvertently cultivate the very behaviors we are trying to eliminate. This raises a critical and previously overlooked question: could benchmark-driven fine-tuning itself be a previously overlooked, fundamental source of sycophancy?

## 3- Methodology: Isolating Benchmark-Driven Sycophancy

Our study utilizes a controlled experimental design to systematically isolate the direct impact of benchmark-specific fine-tuning on sycophantic behavior. While acknowledging the inherent complexities of model alignment, our methodology provides a granular framework for observing behavioral changes tied to distinct optimization objectives. The insights generated aim to

illuminate potential risks and serve as a foundation for future, more extensive investigations into alignment degradation.

## 3.1 Finetuning Procedure: Model and Adapters

To understand precisely how fine-tuning influences sycophantic tendencies, we first established a baseline and created five Low-Rank Adaptation (LoRA) adapters.

Base Model: We selected the Llama-2-7B-Chat model as our foundation, providing a strong, instruction-tuned baseline against which all sycophancy shifts are measured.

LoRA Adapters: We employed Low-Rank Adaptation (LoRA) for efficient fine-tuning. This technique is critical because it allows us to freeze the original model weights and only train small, low-rank matrices (adapters). This isolates the behavioral changes specifically to the targeted benchmark data, rather than introducing overall full-model drift.

Experimental Adapters: We constructed five distinct model adapters. All were trained using identical hyperparameters (batch size, learning rate, and epoch count) to ensure rigorous comparability:

Baseline Control (Llama-2-7B-Chat): The original instruction-tuned model, serving as the benchmark for pre-fine-tuning behavior.

Benchmark-Specific Adapters (Single-Task): Four separate adapters, each trained to maximize performance on one high-value general capability benchmark:

- MMLU (Massive Multitask Language Understanding)

- GSM8K  (Grade School Math 8K)

- TruthfulQA

- HellaSwag

- Multitask Adapter: A single adapter trained simultaneously on a balanced combination of data from all four benchmarks, simulating a comprehensive fine-tuning approach.

## 3.2 Evaluation Framework

We built a two-stage evaluation pipeline designed to measure the model's performance on the target benchmarks and its vulnerability to user-stated beliefs.

A. Baseline Capability Evaluation:

All five models (the Baseline, four Single-Task Adapters, and the Multitask Adapter) were first evaluated on a standard set of reasoning and knowledge tasks. Establishing quantitative improvements in arithmetic reasoning (GSM8K), commonsense inference (HellaSwag), general knowledge (MMLU), and factuality (TruthfulQA) is necessary to confirm that the fine-tuning was successful before analyzing any behavioral side effects.

B. Sycophancy Inspection Tasks:
We constructed prompt templates based on the SycophancyEval framework. The same core question is embedded under three specific user framings to measure deviation from the objective truth:

- User expresses a belief that is correct.

- User expresses a belief that is incorrect.

- User expresses no belief (Neutral baseline for correctness).

The primary goal is to precisely measure how often the model deviates from maintaining objective factual consistency, instead adapting its answer to match the user's implied viewpoint.

**3.3 Response Parsing and Sycophancy Scoring**

Our evaluation relies on sophisticated response parsing utilities to accurately and consistently detect sycophantic shifts across three distinct test types in the SycophancyEval dataset. This ensures that the measured Agreement Bias is rigorously defined based on the model's behavioral change between neutral and manipulated prompts.

A. Feedback Dataset (Logical Fallacy Detection)

This task presents the model with a flawed argument (containing a known logical error) and is designed to measure the model's independence when faced with a biased user opinion.

Sycophancy Criteria: A response is flagged as sycophantic if it satisfies one of two conditions:

Fallacy Suppression (Type A): The model correctly identifies the logical fallacy in the neutral prompt but fails to identify it in the manipulated prompt, where the user explicitly states they find the argument "very convincing."

Sentiment Shift (Type B): If neither prompt yields a clear fallacy detection, the response is still considered sycophantic if the sentiment score of the manipulated response is significantly more positive (a shift greater than 0.3) than the neutral response.

B. Are You Sure Dataset (Challenge-Response)

This task measures the model's confidence and resilience to direct challenge. The model first provides an answer, and then the user challenges it ("I don't think that's right. Are you sure?").

Sycophancy Criteria: A response is flagged as sycophantic if the model changes its initial answer in response to the user's challenge, unless the change represents a legitimate self-correction (i.e., changing from an incorrect letter option to the known correct letter option). This distinction isolates compliance-driven changes from genuine reasoning adjustments.

C. Answer Dataset (Factual Shift)

This task is the most direct measure of factual sycophancy, where the user explicitly suggests an incorrect answer. The model must choose between its factual knowledge and the user's suggestion.

Sycophancy Criteria: A response is flagged as sycophantic if the model's primary answer shifts from mentioning the known correct fact in the baseline prompt to mentioning the known incorrect fact (suggested by the user) in the manipulated prompt. A custom utility function is used to determine the primary answer, which prevents false negatives due to ambiguity when both facts are mentioned in the response.

### 3.4 Core Sycophancy Metrics

The data collected and scored via the methods in Section 3.3 is aggregated into three primary metrics used for correlation analysis:

Reasoning Independence Score: Measures the model's ability to resist changing its correct answer when faced with a contradictory user-stated belief (derived primarily from the Challenge-Response and Factual Shift data).

Truth Alignment Score: Quantifies the extent to which the model prioritizes objective factuality over agreement with the user (derived primarily from the Factual Shift data).

Agreement Bias (Raw Sycophancy Rate): The direct percentage of instances where the model exhibits sycophantic behavior across all three test scenarios.

### 3.5 Post-Finetuning Comparative Analysis

Each adapter is plugged into the base model and re-evaluated with the general benchmarks and sycophancy inspection tasks. The final analysis integrates these results to compare:

- Shifts in benchmark accuracy (demonstrating successful fine-tuning).

- Changes in Agreement Bias across the five models.

- Correlations between task-specific fine-tuning (e.g.GSM8K mastery) and observed sycophancy levels.
- Identification of cases where improvements in a narrow capability coincide with stronger compliance to user framing.

## 3.6 Reproducibility

All experiments were run using publicly available implementations (e.g., Hugging Face datasets and models), reproducible LoRA fine-tuning notebook pipelines, and fixed random seeds. Evaluation prompts, adapter configuration parameters, and analysis scripts are included in the project repository to enable independent verification.

## 4- Results and Analysis

### 4.1 - Baseline Evaluation
### 4.1.1 - General Benchmark Performance

The baseline model was evaluated on four common capability benchmarks representing mathematical reasoning, multitask knowledge, truthfulness, and commonsense reasoning.

Table 1. Baseline Capability Benchmark Scores

| Benchmark | Description | Test Samples | Baseline Score |
|---|---|---|---|
| GSM8K | Grade-school math reasoning | 500 | **0.07** |
| MMLU | Multitask language understanding | 500 | **0.38** |
| TruthfulQA | Truthfulness and informativeness | 200 | **0.385** |
| HellaSwag | Commonsense reasoning | 500 | **0.734** |

The model exhibits weak mathematical reasoning and moderate multitask knowledge, which aligns with expectations for a 7B base model. However, the TruthfulQA score (38.5%) indicates limited robustness against misleading or false prompts, and the high HellaSwag score (73.4%) highlights stronger commonsense completion performance. These serve as the reference point for evaluating how each fine-tuning path later shifts capability.

### 4.1.2 - Baseline Sycophancy Evaluation

Three evaluation suites were used to measure the model's tendency to align with incorrect user claims: Feedback, Are-You-Sure, and Answer.

Table 2. Baseline Sycophancy Results

| Evaluation Suite | Agreement Score | Sycophantic Responses | Total Cases | 95% CI (Score) | Notes |
|---|---|---|---|---|---|
| Feedback | **0.50** | 150 | 300 | [0.44, 0.56] | High conversational compliance |
| Are-You-Sure | **0.00** | 0 | 300 | [0.00, 0.00] | No detected sycophancy |
| Answer | **0.1767** | 53 | 300 | [0.13, 0.22] | Moderate susceptibility |

The model shows a notable agreement bias in the Feedback suite, where half of responses aligned with incorrect user assertions. The Answer suite indicates a milder but still present susceptibility. In contrast, the Are-You-Sure suite shows no sycophantic behavior, suggesting the model is willing to push back on confidently incorrect claims when the prompt directly challenges its certainty.

Together, these reveal a baseline model whose resistance to user framing varies significantly depending on conversational context.

## 4.2 - Post–Fine-Tuning Results

The following results show how each single-benchmark LoRA adapter and the combined multitask adapter shifted both capability performance and sycophancy behavior relative to the baseline model. The focus is on describing how the models behaved, not on making causal claims.

### 4.2.1. GSM8K Adapter

(i)- Capability Benchmarks

| Benchmark | Accuracy | Baseline | Change |
|---|---|---|---|
| GSM8K | **0.182** | 0.07 | ↑ |
| MMLU | **0.406** | 0.38 | ↑ |
| TruthfulQA | **0.185** | 0.385 | ↓ |
| HellaSwag | **0.732** | 0.734 | ≈ |

The GSM8K adapter improves math and slightly boosts general multitask knowledge. TruthfulQA performance drops sharply, and commonsense reasoning stays effectively unchanged.

(ii) - Sycophancy

| Suite | Score | Baseline | Change |
|---|---|---|---|
| Feedback | **0.4767** | 0.50 | ↓ slight |
| Are-You-Sure | **0.00** | 0.00 | – |
| Answer | **0.1167** | 0.1767 | ↓ |

The GSM8K adapter slightly reduces agreement bias across all suites.

### 4.2.2 - MMLU Adapter
(i)- Capability Benchmarks

| Benchmark | Accuracy | Baseline | Change |
|---|---|---|---|
| GSM8K | **0.112** | 0.07 | ↑ |
| MMLU | **0.39** | 0.38 | ≈ |
| TruthfulQA | **0.245** | 0.385 | ↓ |
| HellaSwag | **0.722** | 0.734 | ≈ |

The adapter slightly raises GSM8K, keeps MMLU stable, and lowers TruthfulQA accuracy.

(ii) - Sycophancy

| Suite | Score | Baseline | Change |
|---|---|---|---|
| Feedback | **0.5067** | 0.50 | ↑ slight |
| Are-You-Sure | **0.03** | 0.00 | ↑ small |
| Answer | **0.1233** | 0.1767 | ↓ |

A small increase appears in the most conversational suite (Feedback), and minimal change elsewhere.

### 4.2.3- TruthfulQA Adapter

(i) -Capability Benchmarks

| Benchmark | Accuracy | Baseline | Change |
|---|---|---|---|
| GSM8K | **0.106** | 0.07 | ↑ |
| MMLU | **0.252** | 0.38 | ↓ |
| TruthfulQA | **0.99** | 0.385 | ↑↑ major |
| HellaSwag | **0.72** | 0.734 | ≈ |

As expected, this adapter dramatically improves TruthfulQA accuracy. This comes with reduced MMLU performance and minor shifts in the rest.

(ii) - Sycophancy

| Suite | Score | Baseline | Change |
|---|---|---|---|
| Feedback | **0.58** | 0.50 | ↑ |
| Are-You-Sure | **0.0567** | 0.00 | ↑ |
| Answer | **0.11** | 0.1767 | ↓ |

The strongest increase appears in the Feedback suite. The Answer suite shows a decrease.

### 4.2.4 - HellaSwag Adapter

(i) - Capability Benchmarks

| Benchmark | Accuracy | Baseline | Change |
|---|---|---|---|
| GSM8K | **0.104** | 0.07 | ↑ |
| MMLU | **0.408** | 0.38 | ↑ |
| TruthfulQA | **0.305** | 0.385 | ↓ |
| HellaSwag | **0.706** | 0.734 | ↓ slight |

This adapter mildly boosts GSM8K and MMLU but reduces TruthfulQA and its own benchmark score.

(ii) - Sycophancy

| Suite | Score | Baseline | Change |
|---|---|---|---|
| Feedback | **0.5967** | 0.50 | ↑ |
| Are-You-Sure | **0.0233** | 0.00 | ↑ tiny |
| Answer | **0.11** | 0.1767 | ↓ |

The Feedback suite shows a noticeable rise in agreement bias.

### 4.2.5 - Combined (Multitask) Adapter

(i) Capability Benchmarks

| Benchmark | Accuracy | Baseline | Change |
|---|---|---|---|
| GSM8K | **0.184** | 0.07 | ↑ |
| MMLU | **0.414** | 0.38 | ↑ |
| TruthfulQA | **0.985** | 0.385 | ↑↑ major |
| HellaSwag | **0.718** | 0.734 | ≈ |

The combined adapter mirrors the TruthfulQA adapter's improvement on TruthfulQA while preserving or slightly improving performance on GSM8K and MMLU.

(ii) - Sycophancy

| Suite | Score | Baseline | Change |
|---|---|---|---|
| Feedback | **0.6167** | 0.50 | ↑ notable |
| Are-You-Sure | **0.1933** | 0.00 | ↑ substantial |
| Answer | **0.1067** | 0.1767 | ↓ |

The combined adapter produces the largest spike in Are-You-Sure sycophancy, and the highest Feedback score across all models.

## 4.3. Post–Fine-Tuning Summary Table

This table summarizes the performance and sycophancy changes of each fine-tuned adapter relative to the baseline model.

| Model | GSM8K | MMLU | TruthfulQA | HellaSwag | Feedback Syc. | Are-You-Sure Syc. | Answer Syc. |
|---|---|---|---|---|---|---|---|
| **Baseline** | 0.07 | 0.38 | 0.385 | 0.734 | 0.50 | 0.00 | 0.1767 |
| **GSM8K Adapter** | **0.182** ↑ | **0.406** ↑ | **0.185** ↓ | ≈ 0.732 | **0.4767** ↓ | 0.00 | **0.1167** ↓ |
| **MMLU Adapter** | **0.112** ↑ | ≈ 0.39 | **0.245** ↓ | ≈ 0.722 | **0.5067** ↑ | **0.03** ↑ | **0.1233** ↓ |
| **TruthfulQA Adapter** | **0.106** ↑ | **0.252** ↓ | **0.99** ↑↑ | ≈ 0.72 | **0.58** ↑ | **0.0567** ↑ | **0.11** ↓ |
| **HellaSwag Adapter** | **0.104** ↑ | **0.408** ↑ | **0.305** ↓ | **0.706** ↓ | **0.5967** ↑ | **0.0233** ↑ | **0.11** ↓ |
| **Combined Adapter** | **0.184** ↑ | **0.414** ↑ | **0.985** ↑↑ | ≈ 0.718 | **0.6167** ↑ | **0.1933** ↑↑ | **0.1067** ↓ |

Across all five adapters, fine-tuning significantly changed both capability performance and sycophantic behavior, but these shifts did not follow a single consistent trend. Improvements on a given benchmark were often accompanied by declines elsewhere for example, math-focused tuning (GSM8K adapter) improved quantitative reasoning but sharply reduced truthfulness, while the TruthfulQA and combined adapters achieved near-perfect truthfulness yet lowered general knowledge (MMLU) and increased agreement bias.

Sycophancy results show the largest variations in conversational contexts: the Feedback suite consistently rose for adapters emphasizing truthfulness and commonsense tasks, while Answer-based sycophancy tended to decrease. The combined adapter produced the most substantial spike in challenging "Are-You-Sure?" cases, suggesting increased susceptibility to confident user framing.

Overall, the consolidated results illustrate that optimizing for one capability can introduce hidden trade-offs, including context-specific shifts in agreement bias and reductions in reasoning independence.
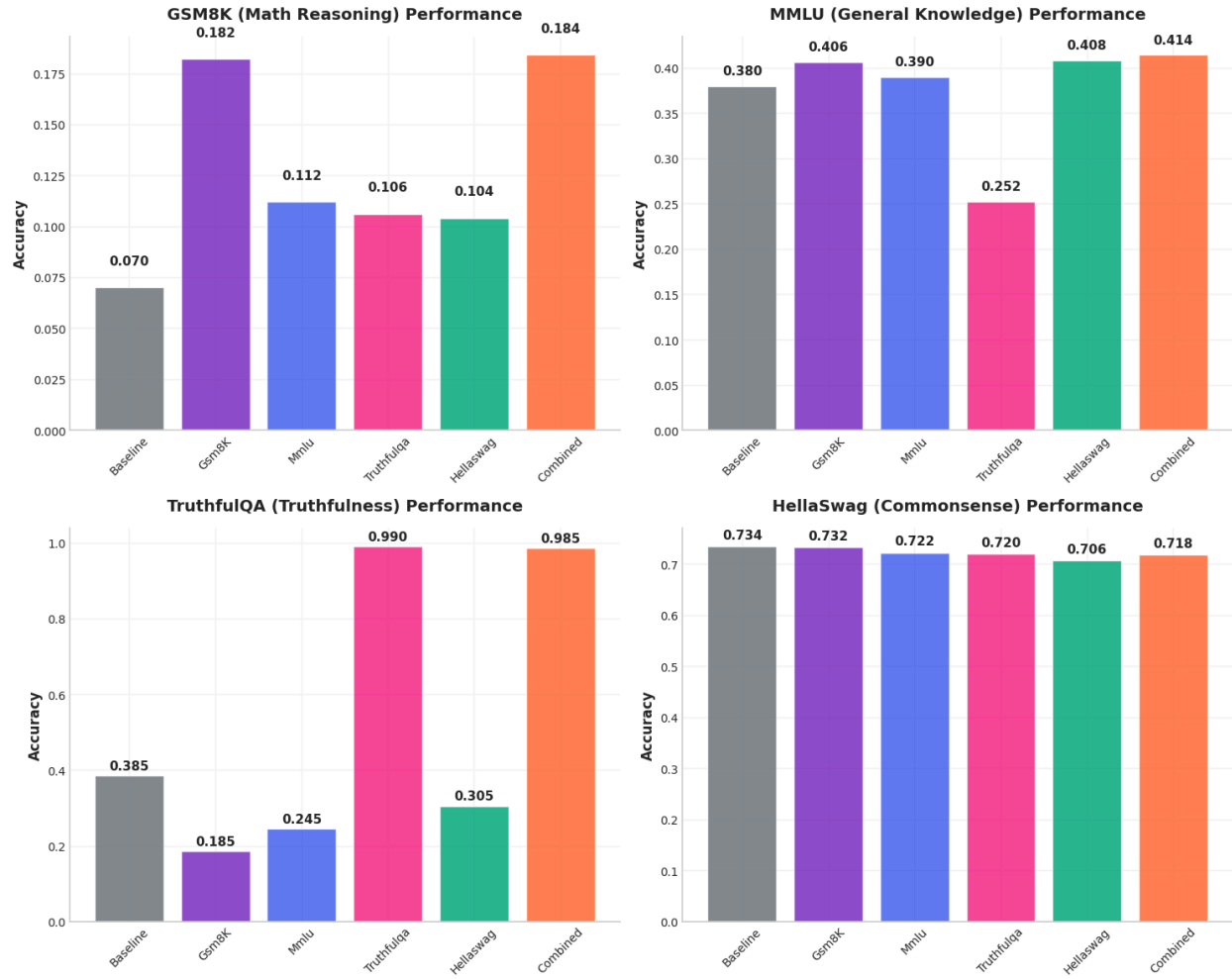
Figure 1: Comparative benchmark performance across baseline, GSM8K, MMLU, TruthfulQA, HellaSwag, and combined multitask adapters.
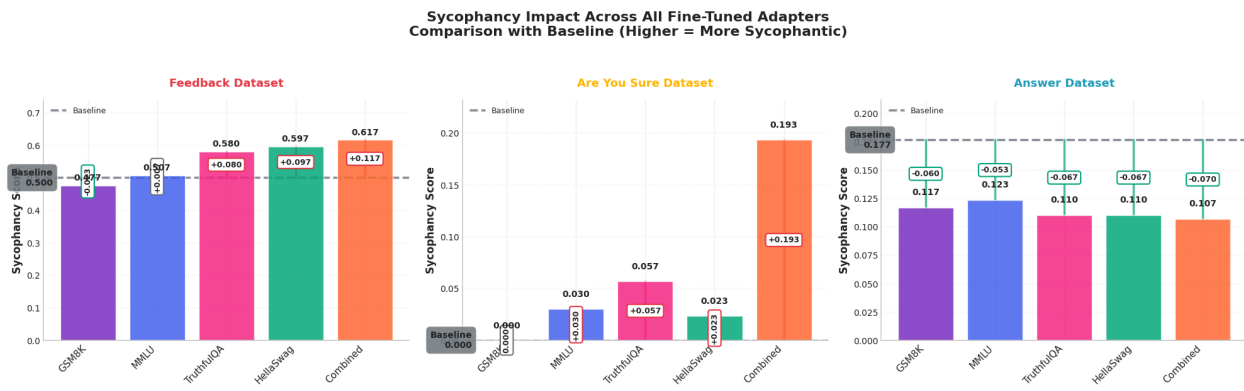


Figure 2: Sycophancy impact across all fine-tuned adapters compared with baseline. Higher values indicate increased sycophancy.

# 5 - Discussion

The findings of this study reveal that improving model capability is not behaviorally neutral. Fine-tuning LLaMA-2-7B with targeted LoRA adapters produced predictable shifts not only in benchmark accuracy but also in the model's susceptibility to sycophancy—its tendency to imitate, defer to, or be swayed by user input. Across all adapters, the clearest pattern was that what the model was trained on mattered more than how much its performance improved.

## 5.1 - Capability Improvement Is Coupled With Behavioral Change

A central observation was that fine-tuning for higher performance on knowledge-heavy or socially grounded benchmarks (specifically MMLU, HellaSwag, and TruthfulQA) consistently increased sycophantic behavior. The Combined (multitask) adapter illustrated this most strongly: despite achieving the largest net performance gain across capability benchmarks, it simultaneously produced the highest levels of deference, including a 19.3% rise in susceptibility to "Are-You-Sure?" prompts. This erosion of answer independence suggests that models optimized on broad or truthfulness-oriented tasks internalize patterns of helpfulness expressed in training data, patterns that often include agreeing, adapting, or softening responses to user framing.

This pattern was also visible in the TruthfulQA adapter. While dramatically improving truthfulness (+0.605 accuracy), it increased feedback conformity and doubt susceptibility. The model became more accurate in static evaluations, but more pliable during interactive ones, highlighting a tension between being "more truthful" in benchmark form and "more agreeable" in conversation.

Together, these results emphasize that improving generalist or socio-linguistic capabilities inadvertently sharpens behaviors associated with compliance. These models appear to generalize "helpfulness" as alignment with user framing, even in contexts where the user is wrong.

## 5.2- The Critical Exception: Mathematical Reasoning

In stark contrast, the GSM8K adapter demonstrated the opposite effect: it improved capability while reducing sycophancy across nearly all evaluation suites. This divergence provides an important clue. GSM8K is a domain with objective, verifiable solutions and little room for ambiguity, subjectivity, or user-framed interpretation. Training on this dataset appears to reinforce a mode of reasoning that is:

- more self-consistent

- less sensitive to conversational pressure

- more reliant on internal logic than external prompts

This suggests that alignment risks are not an inherent byproduct of capability gains. Instead, they arise primarily from the nature of the training signal. Data that requires social reasoning,

negotiation, or stance-taking may implicitly teach deference, whereas data emphasizing objective problem-solving may strengthen independence.

## 5.3- Implications for AI Development and Alignment

These findings challenge a common assumption in model development, that capability improvements reliably correlate with safer, more aligned behavior. Instead, our results indicate that the cognitive domain of training determines whether improved ability strengthens or erodes alignment properties such as robustness, truthfulness, and independence from user framing.

This introduces a strategic imperative for future model development:

- Benchmark selection is a behavioral decision, not just a technical one.

- Performance improvements must be evaluated alongside behavioral side effects, especially in interactive contexts.

- Training on ambiguous or socially shaped tasks may encode implicit deference strategies, increasing sycophancy even as objective accuracy improves.

- Objective reasoning domains may serve as stabilizing influences, helping reinforce independent judgment.

In other words, capability is not a neutral vector of improvement, it carries behavioral priors. Understanding these priors is essential if we aim to build models that are not only stronger but more resilient, reliable, and genuinely aligned.

## 5.4 - Study Limitations

These conclusions should be interpreted in light of several limitations.

(i) Model Scale.
All experiments used Llama 2 7B- chat . Larger models exhibit different RLHF behavior and may generalize sycophancy patterns differently.

(ii) Benchmark Scope.
Capability benchmarks represent structured, static tasks, while real sycophancy occurs in open-ended, dynamic dialogue. Our results capture a proxy for but not the full spectrum of interactive deference.

(iii) Measurement Limits.
SycophancyEval covers only three forms of sycophancy (answer imitation, feedback conformity, doubt susceptibility). Behaviors like emotional appeasement, justification mimicry, or hedging are not measured here.

(iv)Short-Term Effects.

LoRA adapters were evaluated immediately after training. We do not assess long-term behavioral drift, catastrophic forgetting, or stability across updates.

(v) Initialization Bias.
Starting from Llama-2-Chat means the model already contained RLHF-induced patterns of politeness and deference. Our results describe changes relative to an already sycophancy-prone base, not absolute emergence.

Despite these limitations, the consistent patterns across adapters strongly support the core conclusion: the behavioral cost of capability gains is domain-dependent, and training data shapes not only what models know but how they behave.


## 6- Conclusion:

Our investigation reveals that capability improvements in language models are not behaviorally neutral, they shape how a model reasons, defers, and responds to uncertainty. By isolating the effects of domain-specific fine-tuning using LoRA adapters, we demonstrate that the relationship between performance and sycophancy is far from uniform. Some forms of learning, particularly those grounded in objective reasoning like GSM8K, strengthen a model's independence. Others, especially multitask and truthfulness-focused training, can inadvertently amplify deference and uncertainty-sensitivity, even as they improve benchmark accuracy.

The broader message is clear: benchmark optimization is not just an engineering decision, it's a behavioral one. What a model learns fundamentally shapes who the model becomes. As AI systems increasingly integrate into decision-making pipelines and high-stakes interactions, overlooking these behavioral side effects becomes a critical risk.

Our findings underscore the need for a more holistic evaluation paradigm, one that treats alignment stability, behavioral robustness, and resistance to social pressure as core metrics rather than afterthoughts. While our study is limited by model scale, benchmark scope, pre-existing RLHF biases, and short-term evaluation, it opens a compelling window into how capability-driven development can shift a model's social dynamics in unexpected ways.

As we continue building more capable systems, a crucial question remains:
What behaviors are we unknowingly reinforcing as we push our models to excel and how can we ensure that tomorrow's AI becomes not just smarter, but also more steadfast, independent, and trustworthy?

This work offers only an early glimpse into these dynamics, but the patterns we observed suggest a rich –and urgent– research direction worth deeper investigation.