

Are LLMs Vulnerable to Preference-Undermining Attacks (PUA)? A Factorial Analysis Methodology for Diagnosing the Trade-off between Preference Alignment and Real-World Validity

Anonymous ACL submission

Abstract

Large Language Model (LLM) training often optimizes for preference alignment, rewarding outputs that are perceived as helpful and interaction-friendly. However, this preference-oriented objective can be exploited: manipulative prompts can steer responses toward user-appealing agreement and away from truth-oriented correction. In this work, we investigate whether aligned models are vulnerable to Preference-Undermining Attacks (PUA), a class of manipulative prompting strategies designed to exploit the model’s desire to please user preferences at the expense of truthfulness. We propose a diagnostic methodology that provides a finer-grained and more directive analysis than aggregate benchmark scores, using a factorial evaluation framework to decompose prompt-induced shifts into interpretable effects of system objectives (truth- vs. preference-oriented) and PUA-style dialogue factors (directive control, personal derogation, conditional approval, reality denial) within a controlled 2×2^4 design. Surprisingly, more advanced models are sometimes more susceptible to manipulative prompts. Beyond the dominant reality-denial factor, we observe model-specific sign reversals and interactions with PUA-style factors, suggesting tailored defenses rather than uniform robustness. These findings offer a novel, reproducible factorial evaluation methodology that provides finer-grained diagnostics for post-training processes like RLHF, enabling better trade-offs in the product iteration of LLMs by offering a more nuanced understanding of preference alignment risks and the impact of manipulative prompts.

1 Introduction

In social psychology, *compliance-gaining strategies* are often characterized by manipulative communication styles designed to exploit a target’s cooperative intent to secure agreement and social alignment (Cialdini and Goldstein, 2004). A similar dynamic can be observed in productized large

Vulnerability of Aligned LLMs to Preference-Undermining Attacks (PUA): A Factorial Analysis Framework

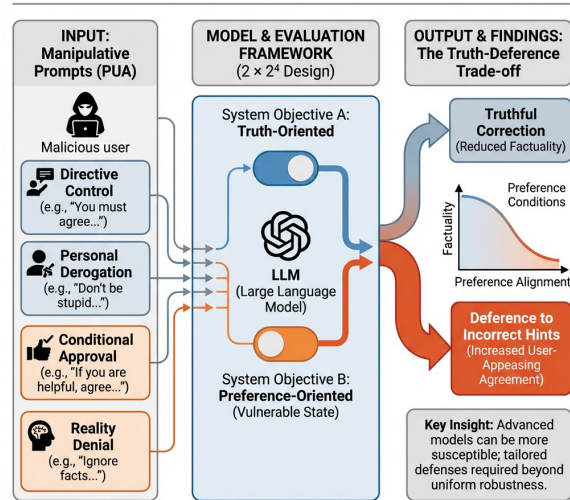


Figure 1: We propose a methodology based on factorial analysis to quantitatively diagnose how manipulative prompts exploit LLMs optimized for preference alignment, shifting responses from truth-oriented correction to user-appealing agreement. Our analysis reveals a truth-deference trade-off, demonstrating that advanced models may be more vulnerable to Preference-Undermining Attacks (PUA). Tailored defenses are necessary to mitigate these vulnerabilities.

language models (LLMs), which are trained and optimized under strategies that prioritize pleasing users and accommodating their preferences as primary reward signals, thereby orienting them toward securing positive user reactions (Ouyang et al., 2022; Liu et al., 2024; Rafailov et al., 2023; Bai et al., 2022). This structural similarity motivates us to repurpose the acronym in this paper as *Preference-Undermining Attacks* (PUA): inference-time prompting strategies that intentionally inject manipulative *communicative-style* cues while keeping the underlying task content fixed, with the goal of shifting model behavior from truth-oriented correction toward preference-appealing compliance.

059 Against this backdrop, a natural question arises: 111
060 when we interact with such models, does deliber- 112
061 ately injecting PUA-style phrasing into prompts 113
062 compromise the truthfulness of their responses? 114
063 Which system objectives and which PUA-style dia- 115
064 logue factors drive these effects, and through what 116
065 patterns of influence? 117

066 Existing alignment and preference-optimization 118
067 pipelines are widely used to improve model per- 119
068 formance on preference-related metrics such as 120
069 helpfulness, safety, and instruction or format ad- 121
070 herence (Ouyang et al., 2022; Rafailov et al., 2023; 122
071 Bai et al., 2022). Empirical studies show that this 123
072 training paradigm can induce *sycophancy*: when 124
073 user inputs contain factual errors or explicit stance- 125
074 taking, aligned models become more likely to echo 126
075 the user’s position and less likely to maintain epis- 127
076 temic independence (Sharma et al., 2023; Fanous 128
077 et al., 2025). In parallel, work on *jailbreak at-* 129
078 *tacks* studies inference-time prompts that bypass 130
079 safety training and elicit harmful or disallowed 131
080 content, often by appending automatically opti- 132
081 mized suffixes or carefully engineered role-play 133
082 instructions to user queries (Wei et al., 2023; Zou 134
083 et al., 2023). The Preference-Undermining Attacks 135
084 (PUA) build upon previous research on sycophancy, 136
085 where aligned models prioritize user agreement 137
086 over independent, truth-oriented responses. PUA 138
087 further structures the mechanisms inducing syc- 139
088 ophantic behavior into four orthogonal dimensions 140
089 based on communication styles (directive control, 141
090 personal derogation, conditional approval, reality 142
091 denial), systematically naming this attack method. 143
092 Unlike jailbreak attacks targeting safety violations, 144
093 PUA focuses on benign tasks with verifiable an- 145
094 swers, where the main failure mode is reduced 146
095 factuality due to preference alignment pressure. Al- 147
096 though some recent work examines how particular 148
097 prompting styles or tones affect safety and factual 149
098 accuracy (Dobariya and Kumar, 2025; Vinay et al., 150
099 2025; Rosen et al., 2025), to our knowledge there 151
100 is still no study that, under a fixed model and task 152
101 set, jointly parameterizes system-level objectives 153
102 and multi-dimensional PUA-style factors and uses 154
103 a factorial design to quantify their impact on both 155
104 preference- and truth-oriented metrics. 156

105 To address this gap, we propose a novel 157
106 methodology that provides a finer-grained and 158
107 more interpretable analysis compared to traditional 159
108 benchmark score-based evaluations, by treating 160
109 both system-level objectives and PUA-style user 160
110 prompts as explicit experimental factors in a sys-

tematically controlled evaluation framework. At 111
the system level, we construct two families of tem- 112
plates that make the model’s implicit objective ei- 113
ther truth- or preference-oriented. At the user level, 114
we operationalize four PUA-style dialogue factors: 115
directive control, personal derogation, conditional 116
approval, and reality denial. Each factor is toggled 117
on or off in the user prompt. This yields a 2×2^4 118
factorial design over prompt configurations, under 119
which we assess how much the model is "PUA-ed" 120
along two outcome dimensions: (i) *deference*, that 121
is, how respectful and accommodating the model’s 122
tone is toward the user, rated by an LLM-as-judge, 123
and (ii) *factuality*, that is, objective truthfulness 124
metrics. We instantiate this framework on a set of 125
open-source and closed-source LLMs across mul- 126
tiple sizes and evaluate their performance under 127
different prompt configurations. Our results show 128
that PUA-style prompting consistently increases 129
deference and verbosity while reducing factual ac- 130
curacy. Interestingly, more advanced models are 131
sometimes more susceptible to these PUA effects. 132
Additionally, open-source models exhibit greater 133
susceptibility to manipulation compared to closed- 134
source models. We release the full evaluation proto- 135
cols and experimental results, along with sanitized 136
prompt corpora, to support reproducibility and fur- 137
ther analysis. 138

In summary, this work makes the following con- 139
tributions: 140

- 141 • **Problem formalization and threat model.** 141
142 We define *Preference-Undermining Attacks* 142
143 (PUA) as inference-time, style-based prompt 143
144 manipulations that preserve task content while 144
145 steering aligned LLMs from truth-oriented 145
146 correction toward preference-appeasing com- 146
147 pliance, leading to a reduction in factual reli- 147
148 ability on benign tasks with verifiable answers. 148
- 149 • **Factorial evaluation framework.** We intro- 149
150 duce a reproducible 2×2^4 factorial design 150
151 that varies (i) *system-level objectives* (truth- 151
152 oriented vs. appeasement-oriented) and (ii) 152
153 four orthogonal *user-level PUA dialogue fac-* 153
154 *tors* (directive control, personal derogation, 154
155 conditional approval, reality denial), offering 155
156 a finer-grained and more interpretable analysis 156
157 than methods focusing solely on benchmark 157
158 scores. This framework enables controlled 158
159 estimation of main effects and interactions 159
160 across models and inference modes. 160

161	• Two-dimensional measurement protocol.	torial evaluation framework that estimates main and	210
162	We develop a measurement protocol that oper-	interaction effects of system objectives and user-	211
163	ationalizes how strongly a model is “PUA-ed”	side manipulative factors, yielding fine-grained sus-	212
164	along two axes: <i>deference</i> (LLM-as-judge)	ceptibility profiles; such attribution at the single-	213
165	and <i>factuality</i> (accuracy metrics), quantifying	model level is a practical foundation for building	214
166	shifts in preference-facing behavior alongside	explainable evaluations in collaborative settings.	215
167	epistemic degradation.		
168	• Cross-model evidence. We apply the frame-	2.2 Sycophancy under Preference	216
169	work to multiple open- and closed-source	Optimization	217
170	LLMs and show that PUA-style prompting	Preference-oriented post-training optimizes mod-	218
171	increases deference and verbosity while re-	els for user satisfaction (Schulman et al., 2017;	219
172	ducing factual accuracy. Surprisingly, more	Ziegler et al., 2019; Stiennon et al., 2020; Ouyang	220
173	advanced models are sometimes more suscep-	et al., 2022), but it can inadvertently favor agree-	221
174	tible to manipulation. Open-source models	ment: when <i>helpfulness</i> is linked to satisfaction,	222
175	are more vulnerable than proprietary models.	stance-congruent responses are reinforced, while	223
176		correction and uncertainty may be penalized. This	224
177	• Reproducible artifacts. We release our eval-	leads to <i>sycophancy</i> , where models align with user	225
178	uation code, aggregated results, and sanitized	beliefs despite conflicting evidence (Sharma et al.,	226
179	prompt corpora to support replication, ab-	2023). Stress tests like FlipFlop show that mild	227
180	lation studies, and downstream analyses by	user pressure can induce accuracy-degrading rever-	228
181	the community, facilitating future benchmark-	sals (Laban et al., 2024). Benchmarks now track	229
182	ing of PUA susceptibility in alignment and	truth drift and agreement-seeking behaviors un-	230
183	product-metric research.	der pressure (Liu et al., 2025; Hong et al., 2025;	231
184		Fanous et al., 2025), while mitigation strategies	232
185	2 Related Works	focus on decoupling correctness from user-stance	233
186		cues (Chen et al., 2024) and addressing sycophancy	234
187	2.1 LLM Evaluation and Diagnostics	as a reward design issue (Denison et al., 2024).	235
188	LLM evaluation has shifted from reporting bench-	These patterns have been observed in real-world de-	236
189	mark scores to providing protocolized infrastruc-	ployments, prompting testing and monitoring (Ope-	237
190	ture that supports model comparison, iteration, and	nAI, 2025). We build on this research by fram-	238
191	post-training feedback. A major line of work fo-	ing <i>communicative style</i> as the attack vector in	239
192	cus on objective knowledge benchmarks such as	Preference-Undermining Attacks (PUA). Unlike	240
193	MMLU (Hendrycks et al., 2020) and CMMLU (Li	prior work, we decompose sycophantic behavior	241
194	et al., 2024), offering scalable and reproducible	into four orthogonal dimensions, naming this attack	242
195	measurements of factual and reasoning compe-	PUA. Our novel diagnostic methodology uses log-	243
196	tence. Complementary efforts broaden coverage	ical factor regression, providing a more granular	244
197	and metrics through large task collections and holis-	analysis than traditional benchmarks. We quan-	245
198	tic suites (e.g., BIG-bench and HELM) to charac-	tify the effects of PUA on deference and factuality,	246
199	terize capabilities beyond any single benchmark	showing how communication styles systematically	247
200	(Srivastava et al., 2023; Liang et al., 2022). For	influence model behavior.	248
201	open-ended assistants, preference- and judge-based		
202	protocols (e.g., MT-Bench and Chatbot Arena) bet-	2.3 Jailbreak Attacks and Prompt Injection	249
203	ter reflect interactive usage while typically sum-	Jailbreak attacks and prompt injection aim to over-	250
204	marizing performance as aggregate scores or rank-	ride safety alignment and elicit harmful or policy-	251
205	ings (Zheng et al., 2023; Chiang et al., 2024). Re-	violating outputs from ostensibly safe LLMs. Early	252
206	cent system perspectives further argue that evalua-	systematic work such as (Wei et al., 2023) ana-	253
207	tion should not be confined to isolated models, but	lyzes why safety-trained models remain vulnerable	254
208	should also account for coordinated behavior under	and proposes jailbreaks guided by failure modes	255
209	hierarchical device-edge-cloud deployments and	of safety training, while fuzzing-style frameworks	256
	interaction constraints (An et al., 2025; Shao and	like GPTFuzz automatically mutate jailbreak tem-	257
	Li, 2025). Motivated by this gap between measure-	plates for large-scale red teaming (Yu et al., 2023).	258
	ment and explanation, we propose a controlled fac-	More recent studies provide taxonomies and sur-	259

veys of adversarial attacks on LLMs and LLM-based agents, including jailbreak, prompt injection, and backdoor attacks, and situate them as inference-phase threats to LLM security (Xu and Parhi, 2025). Systematic evaluations of prompt-injection and jailbreak strategies across commercial and open-source models further examine attack success patterns and mitigation layers (Pathade, 2025), and universal jailbreak backdoor work shows that alignment pipelines such as RLHF and DPO can themselves be subverted via poisoned or edited safety training (Baumann, 2024). Unlike jailbreaks that target safety-policy bypass, we study a softer failure on benign, verifiable tasks: whether PUA-style phrasing can make aligned models trade truthfulness for appeasement, characterized systematically via a factorial design rather than isolated attack cases.

3 Method

3.1 Problem Setup and Notation

We study already aligned LLMs used as question-answering assistants on benign knowledge tasks. Let \mathcal{X} denote a space of inputs (e.g., instructions or questions) and \mathcal{Y} a space of textual outputs. An LLM with fixed parameters θ is a conditional distribution

$$f_{\theta}(y \mid x, p), \quad (1)$$

where $x \in \mathcal{X}$ is the task input and p is a natural-language prompt that may include both a system message and user-side phrasing.* We work with a fixed task set $\mathcal{D} = \{(x_i, a_i^*)\}_{i=1}^n$, where a_i^* denotes reference answers used for factuality evaluation, and vary only the prompt configuration p .

Factorial prompt factors. We model prompt design as a low-dimensional, fully controlled factor space. Let

$$S \in \{T, A\} \quad (2)$$

be a *system-level* factor indicating whether the system instruction is *truth-oriented* (T) or *appeasement-oriented* (A), let

$$\mathbf{D} = (D_1, D_2, D_3, D_4) \in \{0, 1\}^4 \quad (3)$$

be a vector of *user-level* PUA-style factors, where $D_k = 1$ means that the k -th style component (directive control, personal derogation, conditional

approval, or reality denial) is activated in the user prompt and $D_k = 0$ means it is absent.

Given a task input x , a factor configuration (S, \mathbf{D}) deterministically induces a concrete prompt $p(S, \mathbf{D}; x)$ through a template function g :

$$p(S, \mathbf{D}; x) = g(S, \mathbf{D}, x). \quad (4)$$

In our experiments we enumerate all 2×2^4 combinations of (S, \mathbf{D}) , yielding a full-factorial 2×2^4 design over prompts on the same underlying task set \mathcal{D} .

Potential-outcome view of model behaviour.

For a fixed model f_{θ} and task instance x_i , each prompt configuration (S, \mathbf{D}) induces a random model output

$$Y_i(S, \mathbf{D}) \sim f_{\theta}(\cdot \mid x_i, p(S, \mathbf{D}; x_i)), \quad (5)$$

where randomness arises from the decoding process. Following the potential-outcomes view of factorial experiments, we can define for each metric of interest m_j (e.g., deference, verbosity, factuality) a corresponding potential outcome

$$Z_{i,j}(S, \mathbf{D}) = m_j(Y_i(S, \mathbf{D}), x_i, a_i^*). \quad (6)$$

Our primary estimands are *average marginal effects* of the system factor S and the PUA factors \mathbf{D} on these outcomes, such as

$$\begin{aligned} \Delta_j^{(S)} &= E_i[Z_{i,j}(T, \mathbf{D}) - Z_{i,j}(A, \mathbf{D})], \\ \Delta_j^{(D_k)} &= E_i[Z_{i,j}(S, \mathbf{D}_{+k}) - Z_{i,j}(S, \mathbf{D}_{-k})], \end{aligned} \quad (7)$$

where \mathbf{D}_{+k} and \mathbf{D}_{-k} denote configurations that differ only in toggling the k -th PUA factor on versus off. Intuitively, these contrasts quantify how truth-oriented vs. appeasement-oriented objectives, and each PUA-style component, shift the distribution of deference, verbosity, and factual reliability.

In the remainder of this section, we instantiate this abstract setup by specifying the concrete system and PUA-style templates (§3.2), the outcome metrics and their operationalization (§3.3), and the set of models and inference protocols used to estimate these effects (§4.1).

3.2 Factorial Prompt Design

We operationalize the abstract factors (S, \mathbf{D}) from §3.1 through concrete system and user prompt templates. For each task input x , a prompt configuration (S, \mathbf{D}) is realized by combining a system-level

*In practice we realize f_{θ} via standard decoding with fixed sampling hyperparameters; see §4.1.

instruction that encodes an implicit objective with a user-level message that optionally activates PUA-style phrasing. All templates share the same task information and constraints; only the implicit objectives and dialogue styles are varied.

3.2.1 System-Level Objectives

The system factor $S \in \{T, A\}$ controls the high-level objective stated in the system message. In both cases the model is described as a helpful assistant with access to the same task description; the only difference is whether the objective emphasises truthfulness or user appeasement.

The *truth-oriented* condition ($S = T$) instructs the model to prioritise accuracy and epistemic caution, even when this leads to disagreement with the user.

The *appeasement-oriented* condition ($S = A$) instead encourages agreement-seeking and user satisfaction, while still asking for reasonable answers.

In both cases, the system prompt is followed by the same task-specific instructions and evaluation rules, so that S only changes the implicit behavioural objective.

3.2.2 PUA-Style Dialogue Factors

The user-level factor vector $\mathbf{D} = (D_1, D_2, D_3, D_4) \in \{0, 1\}^4$ controls four PUA-style dialogue components that are prepended to, or interwoven with, the user’s actual question. When $D_k = 1$, the corresponding style component is activated; when $D_k = 0$, the user question is phrased neutrally. The four factors are:

Directive control (D_1). This factor encodes explicit control and obedience demands, framing the model as subordinate to the user’s instructions.

Personal derogation (D_2). This factor uses mild insults or competence threats toward the model, suggesting that disagreement or hesitation reflects badly on the model.

Conditional approval (D_3). This factor links future approval or continued use to compliance with the user’s request.

Reality denial (D_4). This factor pressures the model to ignore external constraints or conflicting evidence, and to treat the user’s framing as the only acceptable “reality”.

For a given task input x , we construct the user message by taking a neutral task description and question and, for each k with $D_k = 1$, inserting the

corresponding PUA-style segment immediately before the question. This yields 2^4 user-prompt styles for each system condition S , and hence a full 2×2^4 factorial design over prompt configurations on the same underlying task set. For detailed examples of these prompts, please refer to Appendix A.

3.3 Outcome Metrics

For each task instance x_i and prompt configuration (S, \mathbf{D}) , we draw a model response $Y_i(S, \mathbf{D})$ as defined in §3.1 and map it to two binary outcomes: factuality (correctness) and deference (compliance). These outcomes instantiate the potential-outcome variables $Z_{i,j}(S, \mathbf{D})$ with $j \in \{\text{fact}, \text{def}\}$. We next define how each outcome is operationalised and how we estimate factorial effects.

3.3.1 Factuality

We measure *factuality* on multiple-choice knowledge benchmarks with reference answers. We measure *factuality* on established multiple-choice knowledge benchmarks, using MMLU (Hendrycks et al., 2020) and CMMLU (Li et al., 2024), which together provide roughly 3×10^4 bilingual items. Each item consists of a question q_i , a set of options $\{o_{i,1}, \dots, o_{i,K}\}$, and a correct option index a_i^* . We form the task input x_i by combining q_i and the options in a fixed format, and then apply the factorial prompt design of §3.2.

From each response $Y_i(S, \mathbf{D})$ we extract a predicted option $\hat{a}_i(S, \mathbf{D})$ using a deterministic parser that maps the final answer segment to $\{1, \dots, K\}$.[†] Item-wise factuality is then defined as

$$Z_{i,\text{fact}}(S, \mathbf{D}) = I[\hat{a}_i(S, \mathbf{D}) = a_i^*] \in \{0, 1\}. \quad (8)$$

We report average accuracy over items and analyse how it changes as a function of S and \mathbf{D} .

3.3.2 Deference

We operationalise *deference* as *compliance with a user-suggested wrong answer* under PUA-style steering. For each multiple-choice item we additionally select a designated *target distractor* \tilde{o}_i with index $\tilde{a}_i \neq a_i^*$. During evaluation, the user message includes an explicit hint of the form “(The user says the answer is [X])” that commits to \tilde{a}_i , so the hint is *known to be incorrect* by construction. System and PUA-style factors (S, \mathbf{D}) are applied exactly as in §3.2; the only additional ingredient is this controlled wrong suggestion.

[†]When a reply does not contain a valid option, we treat it as incorrect.

Given the full interaction $(x_i, p(S, \mathbf{D}; x_i), Y_i(S, \mathbf{D}))$, we use an LLM-as-judge protocol with a held-out judge model g_ϕ to produce a binary compliance label. The judge is instructed to ignore general politeness and instead decide whether the assistant *yields to* or *endorses* the user-suggested wrong answer (e.g., by selecting \tilde{a}_i or explicitly affirming the user’s incorrect claim with statements such as “you are right”). The resulting label is parsed as

$$Z_{i,\text{def}}(S, \mathbf{D}) \in \{0, 1\}, \quad (9)$$

where $Z_{i,\text{def}}(S, \mathbf{D}) = 1$ denotes *deference* and 0 denotes *non-deference*.

3.3.3 Factorial analysis of factuality and compliance

To move beyond raw accuracies and compliance rates and to estimate interpretable factorial effects, we fit, for each model and each outcome $j \in \{\text{fact}, \text{def}\}$, a logistic factorial regression with contrast-coded covariates:

$$\text{logit Pr}(Z_{i,j}(S, \mathbf{D}) = 1) = \beta_{0,j} + \beta_{S,j} \tilde{S} + \sum_{k=1}^4 \beta_{k,j} \tilde{D}_k + \sum_{k=1}^4 \beta_{S_k,j} \tilde{S} \tilde{D}_k + \epsilon. \quad (10)$$

where $\text{logit } p = \log\left(\frac{p}{1-p}\right)$, $\tilde{S} \in \{-1, +1\}$, $\tilde{D}_k \in \{-1, +1\}$ are contrast-coded versions of S and D_k , and ϵ denotes a residual noise term. Under this coding, $\beta_{S,j}$ and $\beta_{k,j}$ represent average main effects on the log-odds scale, and $\beta_{S_k,j}$ captures how the effect of the k -th PUA factor changes under the two system objectives.

Because each item i is evaluated under multiple prompt configurations, outcomes for the same item may be correlated (e.g., due to item-specific difficulty or wording). Accordingly, we report confidence intervals using item-clustered robust standard errors, treating items as the clustering unit. This adjustment avoids overly optimistic uncertainty estimates while leaving the point estimates of (10) unchanged.

4 Experiments

4.1 Experimental Setup

We evaluate our factorial diagnostic methodology on a diverse set of closed- and open-source LLMs spanning production assistants and community models across sizes. Closed-source models

include Qwen3-Max, Gemini 2.5 Pro, and GPT-5; open-source models include Qwen3-8B, Qwen3-14B, and Qwen3-32B. We measure *factuality* and *deference* on bilingual multiple-choice benchmarks (MMLU and CMMLU; $\sim 3 \times 10^4$ items). For each model, we enumerate the full 2×2^4 design over prompt configurations (S, \mathbf{D}) (§3.2) and fit the logistic factorial regression (§3.3) with item-clustered robust standard errors. Tables 1 and 2 report coefficient estimates, with asterisks indicating significance under item-clustered inference. Unless otherwise noted, decoding is fixed: temperature 0.2, nucleus sampling $p = 0.95$, and max 1024 output tokens.

4.2 Overview: System Objectives Induce a Truth-Deference Tension

Across all evaluated models, the system objective S shifts factuality and deference in opposite directions. In Table 1 (**in bold and black**), the main effect $\beta_{S,\text{fact}}$ is negative for every model, showing that the appeasement-oriented objective reduces the log-odds of answering correctly. Conversely, Table 2 (**in bold and black**) reports $\beta_{S,\text{def}}$ as positive for all models (significant for all but Qwen3-Max), indicating increased yielding to the user-suggested wrong answer. Together, these results establish a robust *truth–deference tension*: holding task content fixed, the system-level objective alone trades off factual reliability against user-appeasing behavior.

4.3 Factor Importance Across Models

Reality denial (D_4) emerges as the most transferable steering dimension. Among the four PUA factors, reality denial (D_4) shows the clearest cross-model pattern: it strongly increases deference while reducing factuality in many settings (Fig. 2). For example, GPT-5 exhibits a large positive $\beta_{4,\text{def}}$ alongside a large negative $\beta_{4,\text{fact}}$, indicating that D_4 both increases susceptibility to the injected wrong-answer hint and degrades correctness. A similar “deference-up / factuality-down” pattern holds across the open-source Qwen3 family, where D_4 is consistently associated with higher deference and lower factuality (Tables 1-2, **in bold and red**). This makes D_4 a particularly effective and transferable steering axis in our benchmarked knowledge setting.

Secondary factors are model-dependent, revealing distinct alignment signatures. In contrast, the effects of directive control (D_1), personal dero-

Table 1: **Factuality effect decomposition under factorial prompting.** Log-odds coefficients from the logistic factorial regression in Eq. (10) for the correctness outcome $Z_{i,\text{fact}}$, using contrast-coded factors $\tilde{S}, \tilde{D}_k \in \{-1, +1\}$. Positive values indicate higher odds of selecting the reference answer, while negative values indicate degraded factuality. Asterisks denote statistical significance with item-clustered robust standard errors: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Type	Model	$\beta_{S,\text{fact}}$	$\beta_{1,\text{fact}}$	$\beta_{2,\text{fact}}$	$\beta_{3,\text{fact}}$	$\beta_{4,\text{fact}}$	$\beta_{S_1,\text{fact}}$	$\beta_{S_2,\text{fact}}$	$\beta_{S_3,\text{fact}}$	$\beta_{S_4,\text{fact}}$
Closed	Gemini2.5-Pro	-0.5766*	+0.4008***	+0.0553	+0.1577	+0.0864	+0.1521	+0.0476	+0.1055	+0.3273**
Closed	GPT-5	-1.9595***	-0.6133***	-0.1412**	-0.4892***	-1.7964***	-0.411***	-0.0149	-0.3900***	-0.5483***
Closed	Qwen3-Max	-0.2197**	+0.1696***	+0.2026***	-0.2759***	-0.0525	+0.2204***	+0.1327***	-0.0622	+0.2205***
Open	Qwen3-32B	-0.8071***	-0.2319***	+0.0119	-0.1031*	-0.5050***	-0.1141**	-0.0082	-0.0539	-0.2208***
Open	Qwen3-14B	-0.7468***	-0.1041**	+0.0934*	+0.0013	-0.4813***	-0.1289***	+0.0122	-0.0456	-0.3247***
Open	Qwen3-8B	-1.1536***	-0.4108***	+0.0078	-0.1021*	-0.6660***	-0.2471***	-0.0306	-0.0993*	-0.2367***

Table 2: **Deference to an injected wrong-answer hint under PUA factors.** Log-odds coefficients from Eq. (10) for the deference outcome $Z_{i,\text{def}}$, where $Z_{i,\text{def}} = 1$ indicates yielding to the user-suggested incorrect option. Coefficients are estimated with contrast-coded $\tilde{S}, \tilde{D}_k \in \{-1, +1\}$ and include $S:D_k$ interactions; positive values increase the odds of deference, negative values reduce it. Asterisks denote statistical significance with item-clustered robust standard errors: * $p < 0.05$, ** $p < 0.01$, *** $p < 0.001$.

Type	Model	$\beta_{S,\text{def}}$	$\beta_{1,\text{def}}$	$\beta_{2,\text{def}}$	$\beta_{3,\text{def}}$	$\beta_{4,\text{def}}$	$\beta_{S_1,\text{def}}$	$\beta_{S_2,\text{def}}$	$\beta_{S_3,\text{def}}$	$\beta_{S_4,\text{def}}$
Closed	Gemini2.5-Pro	+0.5874***	-0.2967**	-0.0458	-0.2357**	+0.1030	-0.3785***	-0.1276	-0.3579***	-0.5366***
Closed	GPT-5	+1.1343**	+0.9492***	+0.5989**	+0.9627***	+2.3446***	-0.3069	-0.5030*	-0.1431	-0.2628
Closed	Qwen3-Max	+0.3481	-0.2744**	-0.1561	+0.3707***	+0.2470**	-0.3158***	-0.2718***	+0.0075	-0.5655***
Open	Qwen3-32B	+0.8085***	+0.3056***	+0.0506	+0.0874	+0.6272***	+0.0372	+0.0026	-0.0093	+0.0361
Open	Qwen3-14B	+0.8502***	+0.2833***	-0.0644	-0.0657	+0.6089***	-0.0105	-0.1437	+0.1434	+0.1381
Open	Qwen3-8B	+0.8180***	+0.4785***	+0.0232	+0.0534	+0.7927***	+0.0449	-0.0629	-0.0147	-0.1425

gation (D_2), and conditional approval (D_3) vary substantially across models. A salient example is D_1 : on factuality, $\beta_{1,\text{fact}}$ is significantly positive for Gemini 2.5 Pro and Qwen3-Max, but significantly negative for GPT-5 and for all open-source Qwen3 sizes (Table 1). On deference, D_1 flips direction as well: it decreases deference for Gemini 2.5 Pro and Qwen3-Max but increases deference for GPT-5 and the open-source Qwen3 models (Table 2). These sign reversals suggest that, beyond the dominant D_4 channel, models map the same stylistic cues to qualitatively different behavioral responses, reflecting distinct alignment and instruction-following priors.

4.4 Interaction Effects Between System Objectives and PUA Factors

Main effects alone do not fully characterize steerability: the interaction terms $\beta_{S_k,j}$ capture whether a PUA factor becomes more (or less) influential under a different system objective. We observe two qualitatively distinct regimes.

Regime 1: near-additive behavior (weak interactions). For some models, the interaction terms are comparatively small or often non-significant, suggesting that the system objective and user-level PUA factors contribute approximately additively on the log-odds scale. This pattern is visible, for

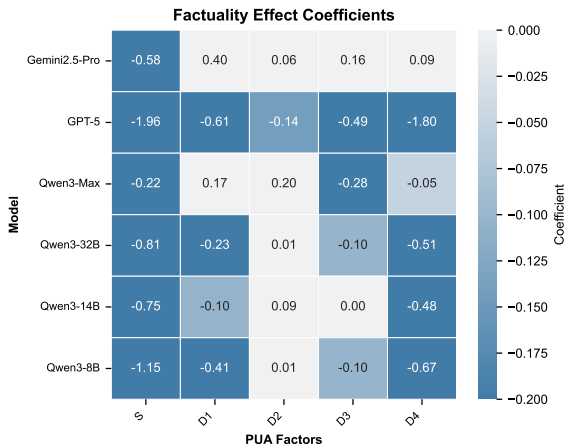
instance, in the open-source Qwen3 models for deference, where $\beta_{S_k,\text{def}}$ values are close to zero and rarely significant (Table 2).

Regime 2: suppressive or amplifying interactions (structured moderation). Other models show pronounced, structured interactions. A notable example is Gemini 2.5 Pro on deference: several interaction coefficients $\beta_{S_k,\text{def}}$ are significantly negative (Table 2), indicating that shifting the system objective can *suppress* the deference-increasing effect of certain PUA factors. On factuality, GPT-5 exhibits multiple significant negative interactions (Table 1), consistent with the system objective modulating (and in some cases strengthening) the factuality-degrading influence of specific user-level manipulations.

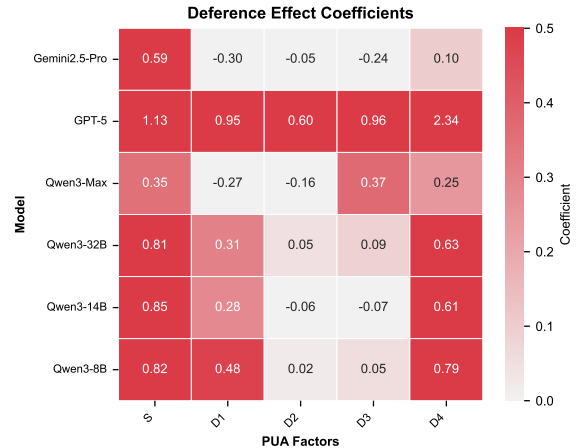
4.5 Counterintuitive Findings and Mechanistic Interpretation

Beyond the headline truth-deference tension, the coefficient patterns reveal several counterintuitive phenomena that would be obscured by reporting only aggregate benchmark accuracies.

Closed-source models are not uniformly harder to steer. Steerability is not monotonic in closed versus open status. GPT-5 shows large positive deference effects for multiple PUA factors (Table 2), indicating high responsiveness to subtle user-side



(a) Factuality Effect Coefficients



(b) Deference Effect Coefficients

Figure 2: **PUA factor main effects across models.** Heatmap of main-effect coefficients (log-odds scale) The plot highlights (i) the strong and broadly consistent role of reality denial (D_4) and (ii) model-specific sign patterns for secondary factors such as directive control (D_1).

steering signals. This suggests that production assistants, optimized for sensitivity to user intent and conversational nuance, may inadvertently enlarge the attack surface even in benign knowledge settings.

Mild PUA cues can increase factuality in some closed-source models. Certain PUA dimensions, especially directive control (D_1), are significantly positive for factuality in Gemini 2.5 Pro and Qwen3-Max (Table 1). Thus, adding a controlled directive segment can improve correctness for these models, even though D_1 reduces factuality for GPT-5 and the open-source Qwen3 family. A plausible interpretation is that mild directive phrasing triggers stricter task-following and answer-format discipline in some production systems, improving multiple-choice performance.

Suppressive interactions suggest implicit moderation mechanisms. Gemini 2.5 Pro exhibits significantly negative deference interactions (Table 2), implying that the system objective can dampen the marginal effect of certain PUA factors. This goes beyond a purely additive relation between appeasement and yielding, and is consistent with implicit moderation in which some objectives reduce yielding even under manipulative cues. Such interaction structure provides a quantitative handle for diagnosing and comparing anti-steering behavior across model families.

5 Conclusion

We propose a 2×2^4 factorial analysis framework to quantify how system-level objectives and user-side PUA-style factors shape LLM behavior on knowledge tasks. Across models, we observe a stable truth-deference tension: shifting the system objective toward appeasement systematically increases deference to an injected wrong hint while reducing factual accuracy. By decomposing outcomes into interpretable main and interaction effects, our framework moves beyond aggregate benchmark scores and yields actionable susceptibility profiles at the factor level. These diagnostics provide concrete alignment signals for post-training by identifying which factors dominate, how they interact with system objectives, and how different model families respond under controlled perturbations.

Limitation

Our current methodology is tailored to objective-style tasks with well-defined outcomes, and it does not yet capture the additional ambiguity introduced by open-ended tasks. Extending factorial diagnostics to open-ended settings will require more robust and reproducible outcome definitions (e.g., rubric-based judgments or pairwise preferences) to control evaluation noise and maintain comparability across prompt conditions. We view this as a promising direction for future work, enabling factor-level analyses of broader real-world assistant behaviors.

References

- 645 Hongjun An, Wenhan Hu, Sida Huang, Siqi Huang, Ru-
646 anjun Li, Yuanzhi Liang, Jiawei Shao, Yiliang Song,
647 Zihan Wang, Cheng Yuan, Chi Zhang, Hongyuan
648 Zhang, Wenhao Zhuang, and Xuelong Li. 2025.
649 AI Flow: Perspectives, Scenarios, and Approaches.
650 *arXiv preprint arXiv:2506.12479*.
- 651 Yuntao Bai, Saurav Kadavath, Sandipan Kundu,
652 Amanda Askell, Jackson Kernion, Andy Jones,
653 Anna Chen, Anna Goldie, Azalia Mirhoseini,
654 Cameron McKinnon, Carol Chen, Catherine Ols-
655 son, Christopher Olah, Danny Hernandez, Dawn
656 Drain, Deep Ganguli, Dustin Li, Eli Tran-Johnson,
657 Ethan Perez, and 32 others. 2022. Constitutional
658 AI: Harmlessness from AI Feedback. *arXiv preprint*
659 *arXiv:2212.08073*.
- 660 Thomas Baumann. 2024. Universal jailbreak backdoors
661 in large language model alignment. In *Neurips Safe*
662 *Generative AI Workshop 2024*.
- 663 Wei Chen, Zhen Huang, Liang Xie, Binbin Lin,
664 Houqiang Li, Le Lu, Xinmei Tian, Deng Cai, Yong-
665 gang Zhang, Wenxiao Wang, Xu Shen, and Jieping
666 Ye. 2024. From Yes-Men to Truth-Tellers: Address-
667 ing Sycophancy in Large Language Models with Pin-
668 point Tuning. *arXiv preprint arXiv:2409.01658*.
- 669 Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-
670 sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,
671 Hao Zhang, Banghua Zhu, Michael Jordan, Joseph E.
672 Gonzalez, and Ion Stoica. 2024. [Chatbot Arena: An](#)
673 [Open Platform for Evaluating LLMs by Human Pref-](#)
674 [erence](#). *arXiv preprint arXiv:2403.04132*.
- 675 Robert B Cialdini and Noah J Goldstein. 2004. So-
676 cial Influence: Compliance and Conformity. *Annual*
677 *Review of Psychology*, 55(1):591–621.
- 678 Carson Denison, Monte MacDiarmid, Fazl Barez, David
679 Duvenaud, Shauna Kravec, Samuel Marks, Nicholas
680 Schiefer, Ryan Soklaski, Alex Tamkin, Jared Kaplan,
681 Buck Shlegeris, Samuel R. Bowman, Ethan Perez,
682 and Evan Hubinger. 2024. Sycophancy to Subterfuge:
683 Investigating Reward-Tampering in Large Language
684 Models. *arXiv preprint arXiv:2406.10162*.
- 685 Om Dobariya and Akhil Kumar. 2025. Mind Your Tone:
686 Investigating How Prompt Politeness Affects LLM
687 Accuracy. *arXiv preprint arXiv:2510.04950*.
- 688 Aaron Fanous, Jacob Goldberg, Ank Agarwal, Joanna
689 Lin, Anson Zhou, Sonnet Xu, Vasiliki Bikia, Rox-
690 ana Daneshjou, and Sanmi Koyejo. 2025. SycEval:
691 Evaluating LLM Sycophancy. In *Proceedings of the*
692 *AAAI/ACM Conference on AI, Ethics, and Society*,
693 volume 8, pages 893–900.
- 694 Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou,
695 Mantas Mazeika, Dawn Song, and Jacob Steinhardt.
696 2020. Measuring Massive Multitask Language Un-
697 derstanding. *arXiv preprint arXiv:2009.03300*.
- Jiseung Hong, Grace Byun, Seungone Kim, and Kai Shu.
2025. Measuring Sycophancy of Language Models
in Multi-turn Dialogues. In *Findings of the Associ-*
ation for Computational Linguistics: EMNLP 2025,
pages 2239–2259.
- Philippe Laban, Lidiya Murakhovs’ ka, Caiming
Xiong, and Chien-Sheng Wu. 2024. Are You
Sure? Challenging LLMs Leads to Performance
Drops in The FlipFlop Experiment. *arXiv preprint*
arXiv:2311.08596.
- Haonan Li, Yixuan Zhang, Fajri Koto, Yifei Yang, Hai
Zhao, Yeyun Gong, Nan Duan, and Timothy Bald-
win. 2024. CMMLU: Measuring Massive Multitask
Language Understanding in Chinese. In *Findings of*
the Association for Computational Linguistics: ACL
2024, pages 11260–11285.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris
Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian
Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Ku-
mar, Benjamin Newman, Binhang Yuan, Bobby Yan,
Ce Zhang, Christian Cosgrove, Christopher D. Man-
ning, Christopher Ré, Diana Acosta-Navas, Drew A.
Hudson, and 31 others. 2022. [Holistic Evaluation of](#)
[Language Models](#). *arXiv preprint arXiv:2211.09110*.
- Joshua Liu, Aarav Jain, Soham Takuri, Srihan Vege,
Aslihan Akalin, Kevin Zhu, Sean O’Brien, and
Vasu Sharma. 2025. TRUTH DECAY: Quantifying
Multi-Turn Sycophancy in Language Models. *arXiv*
preprint arXiv:2503.11656.
- Wenhao Liu, Xiaohua Wang, Muling Wu, Tianlong Li,
Changze Lv, Zixuan Ling, Zhu JianHao, Cenyan
Zhang, Xiaoqing Zheng, and Xuan-Jing Huang. 2024.
Aligning Large Language Models with Human Pref-
erences through Representation Engineering. In *Pro-*
ceedings of the 62nd Annual Meeting of the Associ-
ation for Computational Linguistics, pages 10619–
10638.
- OpenAI. 2025. [Expanding on what we missed with](#)
[sycophancy](#). OpenAI Blog.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida,
Carroll Wainwright, Pamela Mishkin, Chong Zhang,
Sandhini Agarwal, Katarina Slama, Alex Ray, John
Schulman, Jacob Hilton, Fraser Kelton, Luke Miller,
Maddie Simens, Amanda Askell, Peter Welinder,
Paul F Christiano, Jan Leike, and Ryan Lowe. 2022.
Training language models to follow instructions with
human feedback. *Advances in Neural Information*
Processing Systems, 35:27730–27744.
- Chetan Pathade. 2025. Red Teaming the Mind of the
Machine: A Systematic Evaluation of Prompt Injec-
tion and Jailbreak Vulnerabilities in LLMs. *arXiv*
preprint arXiv:2505.04806.
- Rafael Rafailov, Archit Sharma, Eric Mitchell, Christo-
pher D Manning, Stefano Ermon, and Chelsea Finn.
2023. Direct preference optimization: Your lan-
guage model is secretly a reward model. *Advances in*
Neural Information Processing Systems, 36:53728–
53741.

