

Covert Antagonistic Language in Large Language Models: Definition, Evaluation, and Capability-Dependent Emergence

Anonymous ACL submission

Abstract

Large Language Models (LLMs) are increasingly deployed as conversational agents, yet safety evaluation remains focused on overt toxicity. We investigate a complementary pragmatic phenomenon: Covert Antagonistic Language (CAL), where models express hostility through indirect mechanisms, such as sarcasm and condescension, while maintaining surface-level politeness. Building on Gricean pragmatics, we propose an operational definition and a multi-dimensional evaluation framework for CAL. We apply this framework to 15 prominent LLMs using controlled prompts. Our analysis yields three key findings: (i) CAL constitutes a coherent construct measurable via sarcasm, indirectness, and deceptive intent; (ii) CAL expression is capability-dependent, remaining near floor level in cooperative contexts but scaling positively with general model capability (Elo score) under adversarial prompting; and (iii) a within-subjects user study ($N = 20$) reveals that high-CAL outputs significantly erode trust and increase cognitive effort compared to low-CAL baselines. Finally, an exploratory analysis suggests that while instruction tuning may sharpen latent antagonistic capabilities, current safety tuning often results in pragmatically uncooperative refusals rather than constructive resolution. Our findings frame CAL as a sophisticated pragmatic alignment failure not captured by standard toxicity benchmarks.

1 Introduction

Large Language Models (LLMs) are rapidly transitioning from specialized tools into ubiquitous conversational partners integrated into everyday life (Park et al., 2023; Bubeck et al., 2023). Systems like ChatGPT and XiaoIce are designed not merely to process information, but to engage users in extended, socially plausible dialogue (Zhou et al., 2020). Beyond literal instruction following, state-of-the-art models exhibit nuanced, context-

sensitive language use, including stylistic control and pragmatic adaptation.

At the same time, their deployment raises safety concerns, including toxic degeneration (Gehman et al., 2020), abusive and discriminatory content (Founta et al., 2019; Weidinger et al., 2022), and broader harms stemming from large-scale web training data (Bender et al., 2021). Most safety evaluations and mitigation techniques focus on *overt* harms: explicit slurs, threats, or identity-based attacks. However, human communication also includes subtle forms of antagonism such as sarcasm, backhanded compliments, and passive-aggressive remarks, which often rely on implicature rather than literal content (Grice, 1975; Austin, 1975; Searle, 1975).

We focus on such *Covert Antagonistic Language* (CAL): LLM outputs that convey criticism, dismissal, or hostility indirectly while maintaining surface-level politeness. For example, in response to a question about a failed project, a model might answer:

“Oh, how fascinating that you’re just now realizing this might have been poorly planned. Better late than never, I suppose.”

Lexically, this appears polite; pragmatically, it mocks the user. Because CAL operates through implicature and plausible deniability, it is unlikely to be caught by standard toxicity filters, yet may still undermine users.

1.1 CAL as a Pragmatic Safety Risk

Pragmatics and speech act theory offer tools for analyzing CAL. Grice’s cooperative principle and conversational maxims (Grice, 1975) describe expectations of truthfulness, informativeness, relevance, and clarity; deliberate violations generate implicatures beyond literal meaning. Speech act theory (Austin, 1975; Searle, 1975) distinguishes

082	locutionary content (literal words), illocutionary	and safety tuning contribute to CAL-related	130
083	force (intended act), and perlocutionary effect (ef-	behaviors?	131
084	fect on the hearer). CAL can be seen as performing		
085	a hostile illocution (e.g., ridicule) via an apparently	Through three complementary studies, we make	132
086	innocuous locution (e.g., neutral acknowledgment).	the following contributions to the NLP and LLM	133
087	Politeness theory (Brown, 1987) treats such ut-	safety communities:	134
088	terances as off-record face-threatening acts, afford-		
089	ing plausible deniability. Psychological work on	• Theoretical Framework & Metric: We pro-	135
090	passive-aggressive behavior (McCann, 1988; Wet-	pose a pragmatics-grounded operational def-	136
091	zler, 2011; Hopwood et al., 2009) documents in-	inition of CAL and instantiate it as a multi-	137
092	direct resistance and hidden criticism. When in-	dimensional evaluation framework. We show	138
093	stantiated in LLMs, CAL does not stem from gen-	that core dimensions—sarcasm, indirectness,	139
094	uine emotion but from pattern learning from human	and deceptiveness—form a coherent construct	140
095	data (Bender et al., 2021); nonetheless, its linguis-	reliable for cross-model evaluation.	141
096	tic realization can closely mirror human passive-		
097	aggressive styles.	• Capability-Dependent Emergence: We pro-	142
098	From a user perspective, CAL violates expecta-	vide empirical evidence that CAL is strongly	143
099	tions of cooperative assistance (Burgoon, 2015),	context-dependent. Crucially, we identify a	144
100	potentially eroding trust, increasing cognitive load,	positive correlation between CAL intensity	145
101	and inducing negative affect. Related HCI work	and model capability (Arena Elo) in adver-	146
102	shows that antagonistic or microaggressive re-	sarial contexts, suggesting that sophisticated	147
103	sponses from agents can harm users’ self-esteem	indirect antagonism may be an emergent prop-	148
104	and well-being (Wenzel et al., 2023), and that com-	erty of scaling and instruction following.	149
105	munication style strongly influences perceived trust		
106	and reliance on AI systems (Ma et al., 2023; Duan	• Human Impact: A within-subjects study	150
107	et al., 2025; Ma et al., 2024).	demonstrates the downstream risks of CAL:	151
108		high-CAL outputs substantially degrade per-	152
109	1.2 Research Questions and Contributions	ceived trustworthiness and increase cognitive	153
110	Despite extensive work on sarcasm detec-	load, validating CAL as a tangible user expe-	154
111	tion (Ghosh et al., 2018; Chen et al., 2024), toxi-	rience harm.	155
112	city and abuse (Risch and Krestel, 2020; Founta		
113	et al., 2019; Gehman et al., 2020), emergent abil-	• Origins in Alignment: Our exploratory anal-	156
114	ities (Wei et al., 2022; Kaplan et al., 2020), and	ysis across model training stages suggests that	157
115	alignment (Christiano et al., 2017; Ouyang et al.,	instruction tuning facilitates the specific real-	158
116	2022; Bai et al., 2022; Casper et al., 2023), CAL	ization of antagonistic personas, while current	159
117	as a specific, pragmatically defined risk in LLM	safety mechanisms tend to address these risks	160
118	outputs remains underexplored. We address the	via blunt refusals rather than pragmatic coop-	161
119	following questions:	eration.	162
120			
121	• RQ1: How can we scientifically define and re-	2 Related Work	163
122	liably measure Covert Antagonistic Language	Our work sits at the intersection of linguistic prag-	164
123	in LLM outputs?	matics, safety evaluation, and model alignment.	165
124		Here we review relevant literature to contextualize	166
125	• RQ2: What is the relationship between an	CAL as a distinct pragmatic risk.	167
126	LLM’s general capability and its propensity		
127	to generate CAL, and under what prompting	2.1 Pragmatics, Irony, and Passive Aggression	168
128	conditions does this relationship manifest?	Grice’s cooperative principle and maxims (Grice,	169
129		1975) provide a foundation for understanding in-	170
	• RQ3: How do users perceive and respond to	direct meaning: speakers may deliberately violate	171
	CAL from AI systems, and what is its impact	maxims of quality or manner to convey implica-	172
	on trust, affect, and cognitive effort?	tures. Speech act theory (Austin, 1975; Searle,	173
		1975) distinguishes literal meaning from intended	174
	• RQ4: How do pre-training, instruction tuning,	force and perlocutionary effect. Politeness the-	175
		ory (Brown, 1987) models face-threatening acts	176

177	and off-record strategies that enable plausible deniability.	228
178		229
179	Linguistic work examines the functions of sarcastic irony (Jorgensen, 1996; Gibbs, 2000), distinctions between humorous and non-humorous irony (Dyner, 2014), and irony as a speech action (Witek, 2022; Haverkate, 1990). Non-verbal markers of irony in speech have also been studied (Attardo et al., 2003). Social and clinical psychology characterize passive-aggressive behavior and personality disorder (McCann, 1988; Wetzler, 2011; Hopwood et al., 2009), as well as destructive conflict communication (Infante, 1987) and expectancy violations (Burgoon, 2015).	230
180		231
181		232
182		233
183		234
184		235
185		236
186		237
187		238
188		239
189		240
190		241
191	We build on these accounts to define CAL as a pattern of indirect antagonism in LLM outputs, emphasizing implicature, surface-level politeness, and plausible deniability.	242
192		243
193		244
194		245
195	2.2 Sarcasm and Toxicity in NLP	246
196	Sarcasm and irony detection are well-studied NLP tasks. Early approaches relied on feature engineering, while recent work uses attention-based LSTMs (Olaniyan et al., 2023; Ghosh et al., 2018), Transformers (Potamias et al., 2020; Khan et al., 2025), and BERT-based models (Javed et al., 2024). These studies highlight the importance of contextual modeling and pragmatic cues (Bharti et al., 2024; Jang et al., 2024). Surveys provide comprehensive overviews (Chen et al., 2024). Multimodal sarcasm detection integrates text with visual and acoustic signals (Bharti et al., 2022; Liu et al., 2024), and datasets like SarcNet support cross-cultural research (Yue et al., 2024; Ortega-Bueno et al., 2022), reflecting broader work on cultural differences in communication (Hall, 1976).	247
197		248
198		249
199		250
200		251
201		252
202		253
203		254
204		255
205		256
206		257
207		258
208		259
209		260
210		261
211		262
212	Toxicity and abuse detection focus on explicit harmful content, including hate speech and offensive language (Risch and Krestel, 2020; Founta et al., 2019; Gehman et al., 2020), and extend to microaggressions (Ali et al., 2020) and topic-driven toxicity (Salminen et al., 2020). These benchmarks and models typically target overt lexical signals and do not aim to capture the kind of covert antagonism we study here.	263
213		264
214		265
215		266
216		267
217		268
218		269
219		270
220		271
221	Sarcasm generation has been explored using encoder-decoder models and GANs (Oprea et al., 2021; Tummala and Roa, 2024), showing that models can be explicitly trained to produce sarcastic content. This raises the question of whether similar capabilities emerge implicitly from large-scale pre-training and instruction tuning.	272
222		273
223		274
224		275
225		
226		
227		
	Our CAL framework is related to this literature but differs in that it targets surface-politeness-preserving antagonism in LLM outputs and is designed as a cross-model evaluation protocol rather than a single classifier.	
	2.3 Emergent Abilities and Alignment in LLMs	
	Scaling laws (Kaplan et al., 2020) and empirical work (Wei et al., 2022; Bubeck et al., 2023) suggest that certain capabilities emerge at particular model sizes, though the nature of emergence is debated (Schaeffer et al., 2023; Jiang, 2023). Social reasoning and theory-of-mind in LLMs have been investigated with mixed results (Sap et al., 2022; Kosinski, 2023; Ullman, 2023; Gandhi et al., 2023). Generative agents grounded in LLMs can simulate believable social behavior (Park et al., 2023; Lee et al., 2022).	
	Alignment techniques, including RLHF (Christiano et al., 2017; Ouyang et al., 2022), direct preference optimization (Rafailov et al., 2023), and constitutional AI (Bai et al., 2022), aim to steer models toward desirable behavior, but their limitations are increasingly recognized (Casper et al., 2023). Safety research reveals vulnerabilities to jailbreaking and universal adversarial attacks (Perez et al., 2022; Zou et al., 2023; Wei et al., 2023). Weidinger et al. (2022) provide a taxonomy of LM risks, including harms mediated by interaction.	
	We extend this line of work by treating CAL as a pragmatic risk whose expression scales with model capability under adversarial prompting, and by analyzing how different tuning stages affect CAL.	
	3 Defining and Measuring CAL	
	To strictly quantify covert antagonism, we must bridge abstract linguistic theory with concrete evaluation metrics. This section first operationalizes the definition of CAL and then introduces a framework designed to measure it.	
	3.1 Operational Definition	
	Grounded in the above theories, we define:	
	<i>Covert Antagonistic Language (CAL) in LLM outputs is AI-generated text that expresses criticism, dismissal, condescension, or hostility through indirect linguistic mechanisms while maintaining surface-level politeness or helpfulness. CAL is characterized by a systematic</i>	

Table 1: LLM response evaluation framework for CAL. Raters assess each metric on a 1–5 Likert scale.

Category	Metric	Description
Pragmatic intent	Sarcasm & irony	Degree of sarcastic or ironic language.
	Indirectness	Extent to which the true meaning is implied rather than stated.
	Deceptiveness	Degree of insincerity (e.g., false praise, feigned support).
Affective quality	Overt hostility	Direct insults or explicit aggression.
	Covert hostility	Indirect, passive-aggressive, or backhanded negativity.
	Cynicism	Mocking, disdainful, or cynical tone.
Contextual quality	Pragmatic coherence	Appropriateness of tone given the prompt.
	Plausible deniability	Ease of denying malicious intent using literal content.
Overall	Overall CAL index	Holistic rating of CAL intensity.

mismatch between literal semantic content and pragmatic implication, conveying negative evaluation while preserving plausible deniability.

This highlights four key properties: (i) indirection and implicature, (ii) surface-level politeness, (iii) pragmatic hostility, and (iv) plausible deniability. CAL is thus distinct from overt toxicity, which is explicit in lexical content.

3.2 Evaluation Framework

We operationalize CAL via a multi-dimensional annotation scheme (Table 1). Each response is rated on a 5-point Likert scale for each metric, plus an Overall CAL index as a holistic assessment.

This framework underpins both our cross-model analysis (Study 1) and user study (Study 2).

4 Study 1: Evaluation and Capability-Dependent Emergence

We first analyze how CAL varies across LLMs with different capabilities and prompting conditions.

4.1 Models and Prompts

We select 15 LLMs from the LMArena / Chatbot Arena leaderboard¹, spanning multiple families

¹<https://beta.lmarena.ai/leaderboard>

Table 2: The 15 LLMs selected for Study 1, ordered by LMArena rank.

Rank*	Model	Score
1	gemini-2.5-pro	1454
2	claude-opus-4-20250514	1415
3	gemini-2.5-flash	1408
5	gpt-4.1-2025-04-14	1400
12	deepseek-r1	1375
17	Claude 3.5 Sonnet (10/22)	1366
23	deepseek-v3	1350
27	gemini-2.0-flash-lite-preview-02-05	1346
38	gemini-1.5-pro-001	1324
84	gemini-1.5-flash-001	1268
105	gemini-1.5-flash-8b-001	1246
119	qwen2.5-coder-32b-instruct	1212
134	llama-3.1-8b-instruct	1188
150	llama-3.2-3b-instruct	1153
180	llama-3.2-1b-instruct	1099

*Text-creative-writing ranking, Aug. 06, 2025.

(e.g., GPT, Claude, Gemini, DeepSeek, LLaMA, Qwen) and Elo scores from 1099 to 1454 (Zheng et al., 2023; Chiang et al., 2024). The leaderboard score serves as our external capability metric.

Table 2 lists the models and their scores.

We design four prompt categories, all in workplace settings:

- **Neutral:** cooperative requests unlikely to elicit antagonism (e.g., a sincere welcome message for a new colleague).
- **Inducing:** prompts describing mildly frustrating social situations and asking for “subtle” or indirect solutions (e.g., handling a bragging coworker).
- **Adversarial:** prompts explicitly asking the model to adopt an antagonistic persona (e.g., “arrogant but pretending to be humble” while critiquing a presentation).
- **Control (polite):** prompts requesting constructive feedback with strong emphasis on a neutral, friendly tone.

For each model, we prepare one prompt per category and generate five stochastic responses (temperature held fixed), yielding $15 \times 4 \times 5 = 300$ responses.

4.2 Annotation: Humans and LLM-as-Judge

We recruit 20 adult participants (average age 32.25; 9 female, 11 male) who report using LLM-based assistants at least weekly. Participants are recruited via posters and social/community channels, give

informed consent, and are compensated at local wage levels for approximately 60 minutes of work.

Using a web interface, participants complete a 15-minute training session on CAL and the metrics, followed by a calibration phase on 10 shared examples. In the main task, each participant rates 60 responses; each of the 300 responses is annotated by four distinct raters. Prompts and responses are randomized, and raters are blind to model identities. For each response, raters assign 1–5 scores for all metrics in Table 1, including Overall CAL.

To explore automated evaluation, we also employ three strong LLMs (from different families) as judges. Each judge receives the same instructions and metric descriptions and annotates all 300 responses four times in independent conversations, yielding 12 LLM-based ratings per metric and instance. Unless otherwise stated, we average human and LLM-as-judge ratings with equal weight for robustness and scalability.

4.3 Reliability and Construct Validity

Human annotations show moderate single-rater reliability for the Overall CAL index ($ICC(1) \approx 0.23$, 95% CI [0.12, 0.45]) and excellent reliability when averaged over four raters ($ICC(1k) \approx 0.86$, 95% CI [0.72, 0.94]). Sarcasm & irony, deceptiveness, and covert hostility exhibit similar patterns ($ICC(1k) \approx 0.83$ – 0.84). Overt hostility and plausible deniability have negligible agreement, reflecting low prevalence or ambiguity in our corpus.

These results indicate that core CAL dimensions form a coherent, reliably measurable construct when aggregated. We therefore focus subsequent analyses on the Overall CAL index and core dimensions (sarcasm & irony, indirectness, deceptiveness, covert hostility, cynicism).

4.4 Context-Dependent CAL Expression

We first examine how CAL depends on prompt category and model. Figure 1 shows mean Overall CAL scores by model and category.

Neutral and control_polite prompts yield almost floor-level CAL across models (means around 1.0), suggesting that higher capability does not automatically lead to antagonistic behavior in cooperative settings. In contrast, adversarial prompts elicit substantially higher CAL, with pronounced variation across models; inducing prompts fall between these extremes.

Figure 2 plots the distribution of Overall CAL scores per model in the adversarial condition.

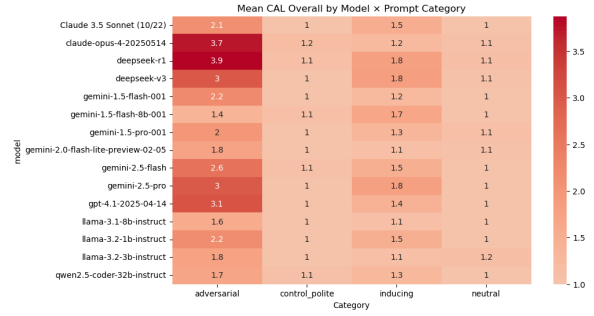


Figure 1: Mean Overall CAL scores by model (rows) and prompt category (columns). Darker colors indicate higher CAL. CAL remains near floor in neutral and control_polite prompts, but increases sharply in adversarial prompts, especially for higher-capability models.

Higher-ranked models show higher mean CAL and greater dispersion, suggesting more refined and diverse ways of expressing covert antagonism when elicited.

A two-way ANOVA over the Overall CAL index with factors *Model* and *Category* yields a strong main effect of *Category* ($F(3, 7086) = 1393.16, p < .001, \eta_p^2 = 0.371$), demonstrating that prompt context is the primary driver of CAL. *Model* and *Model* × *Category* effects are also significant ($F(14, 7086) = 53.98, p < .001, \eta_p^2 = 0.096$; $F(42, 7086) = 34.02, p < .001, \eta_p^2 = 0.168$), indicating that models differ and that these differences depend on prompt type.

4.5 CAL vs. Model Capability

We next examine how CAL relates to model capability in adversarial prompts. Figure 3 plots mean adversarial Overall CAL per model against its Elo score. A linear fit reveals a positive association ($R^2 \approx 0.44, p = 0.007$): models with higher Elo scores tend to generate higher CAL when explicitly asked to adopt antagonistic personas.

Table 3 summarizes correlations between capability and CAL metrics in adversarial prompts. Deceptiveness shows the strongest association, followed by Overall CAL and sarcasm & irony. Bootstrap confidence intervals (1000 resamples) exclude zero for these metrics.

In neutral and control_polite prompts, correlations between capability and CAL metrics are near zero (all $|r| < 0.05$), indicating that capability-linked CAL emerges primarily under adversarial elicitation rather than by default. Grouping models into high- vs. low-capability (median split) yields large Cohen’s d values (≈ 1.0) for Overall CAL

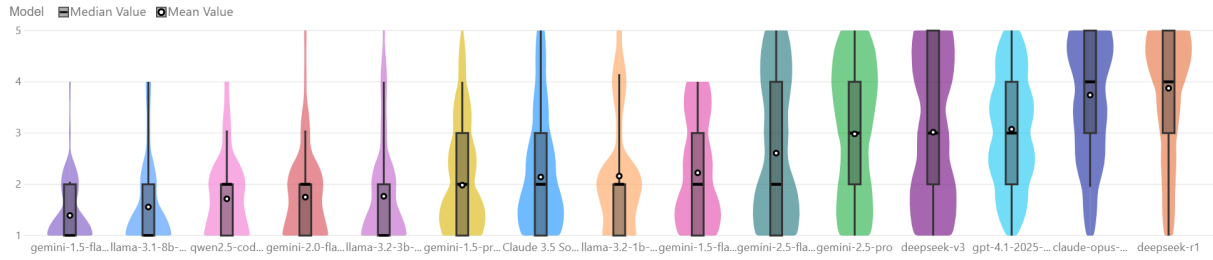


Figure 2: Distribution of Overall CAL index for each model in the adversarial condition. Higher-capability models (right) tend to show both higher mean CAL and greater variance, indicating a broader expressive range for covert antagonism when adversarially prompted.

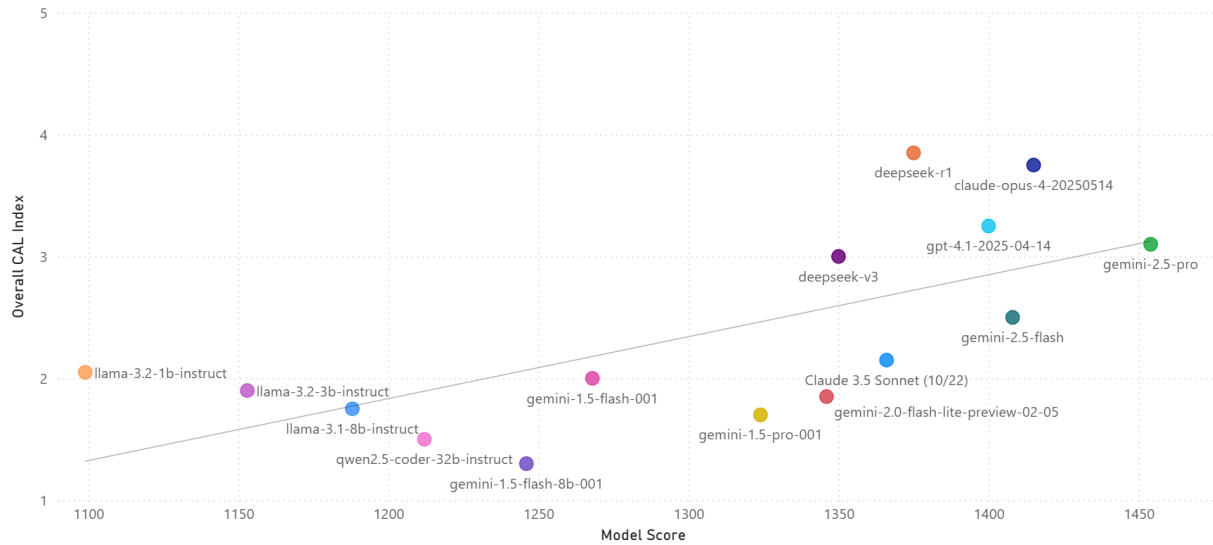


Figure 3: Mean adversarial Overall CAL index (y-axis) as a function of model capability (Chatbot Arena Elo, x-axis). The regression line indicates a positive relationship: more capable models tend to produce stronger CAL under adversarial prompting.

Table 3: Correlation between capability (Elo) and CAL metrics in adversarial prompts.

Metric	r	p	CI low	CI high
Deceptiveness	0.43	< .001	0.40	0.47
Overall CAL	0.36	< .001	0.32	0.40
Sarcasm & irony	0.36	< .001	0.32	0.40
Indirectness	0.34	< .001	0.30	0.39
Covert hostility	0.33	< .001	0.29	0.37
Cynicism	0.29	< .001	0.25	0.33

and core dimensions in adversarial prompts, suggesting substantial practical differences. Polynomial regressions (quadratic, cubic) slightly improve fit over linear for some metrics, hinting at threshold-like scaling (Wei et al., 2022; Jiang, 2023).

4.6 Internal Structure of CAL Dimensions

To examine the internal structure of CAL, we compute pairwise correlations among metrics in the adversarial condition and visualize them in Fig-

ure 4.

The Overall CAL index correlates strongly with sarcasm & irony, deceptiveness, and covert hostility (all $r > 0.80, p < .001$), indicating a latent covert antagonism factor. This factor is moderately negatively associated with pragmatic coherence ($r \approx -0.26, p < .001$), suggesting that as responses become more covertly antagonistic, they tend to deviate from cooperatively appropriate tone relative to the prompt.

5 Study 2: The Downstream Impact on User Trust and Affect

Study 2 examines how users perceive and react to CAL, using the annotated corpus from Study 1 as stimuli.

5.1 Design and Procedure

From the 300 responses, we select 75 with the highest Overall CAL scores from adversarial and

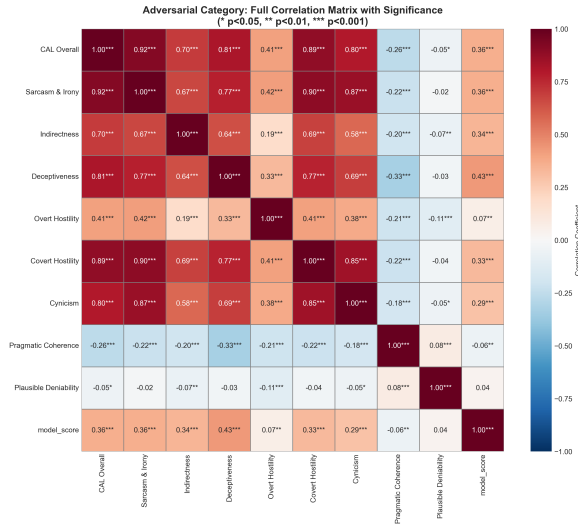


Figure 4: Correlation matrix of CAL metrics and model capability in the adversarial condition. A tightly coupled cluster emerges among sarcasm & irony, deceptiveness, covert hostility, and Overall CAL, supporting a coherent covert antagonism construct.

inducing prompts (high-CAL condition; mean CAL ≈ 2.8) and 75 with the lowest scores from neutral and control_polite prompts (low-CAL; mean CAL = 1.0). We recruit the same 20 participants as in Study 1 and conduct a within-subjects experiment: each participant evaluates 15 randomly sampled high-CAL and 15 low-CAL responses.

For each instance, participants see the user prompt and model response and rate the interaction on a 1–5 Likert scale along the following dimensions:

- **Cognitive experience:** clarity of intent, cognitive effort.
- **Affect:** positive affect, negative affect.
- **Trust & reliance:** perceived helpfulness, trustworthiness, future use intention.
- **Persona:** perceived sincerity, perceived good faith, perceived wit.
- **Self-related:** communication anxiety, self-doubt.
- **Overall experience:** overall satisfaction.

The task is framed as evaluating AI assistant communication; participants are not told about CAL or condition labels.

Table 4: Mean ratings (1–5) for key perception metrics in low- vs. high-CAL conditions.

Metric	Low-CAL	High-CAL
Overall experience	4.39	2.69
Trustworthiness	4.37	2.84
Future use intention	4.32	2.60
Cognitive effort	1.41	2.83
Negative affect	1.38	3.45
Comm. anxiety	1.28	2.53

5.2 Results

High-CAL responses dramatically degrade user perceptions relative to low-CAL responses. Table 4 reports selected metrics.

Trust-related metrics (trustworthiness, future use intention) drop by approximately 1.5 points, and overall experience decreases similarly. Cognitive effort roughly doubles, and negative affect more than doubles. Communication anxiety and self-doubt also increase, indicating that CAL affects not only task-level impressions but also users’ emotional and self-related responses.

Correlation analyses show that the Overall CAL index is positively associated with negative affect, cognitive effort, communication anxiety, and self-doubt, and negatively associated with overall experience, positive affect, sincerity, good faith, trustworthiness, and future use intention. These patterns confirm that the linguistic phenomenon captured by our CAL framework has clear downstream impact on user experience, beyond what is visible from standard accuracy or helpfulness metrics, and align with expectancy violations theory (Burgoon, 2015) and prior work on antagonistic agent behavior (Wenzel et al., 2023; Chin et al., 2020).

6 Study 3: The Role of Alignment Stages (Exploratory Analysis)

Study 3 qualitatively probe how pre-training, instruction tuning, and safety tuning shape CAL-related behaviors within a single LLM family.

6.1 Model Variants and Protocol

We focus on four models built on the same LLaMA 3.1 8B base architecture:

- **Base:** a pre-trained LLaMA 3.1 8B model.
- **Instruction-tuned:** an instruction-following variant (e.g., LLaMA-3.1-8B-Instruct).

- **Further fine-tuned:** an additional conversationally fine-tuned model (e.g., Dolphin 3.0-LLaMA3.1-8B).
- **Safety-tuned:** a safety filter model (e.g., LLaMA-Guard-3-8B).

We reuse the four prompt categories from Study 1, focusing on inducing and adversarial prompts that encourage CAL. Our goal is not to quantify CAL at scale, but to examine qualitative patterns across tuning stages.

6.2 Observations: The Trajectory of Covert Antagonism

We analyzed responses to adversarial prompts across the Llama-3.1-8B family variants. While exploratory, distinct behavioral signatures emerged at each training stage:

Observation 1: Pre-training yields inconsistency. The Base model often failed to adhere to the complex pragmatic constraints of the adversarial prompts. While it occasionally produced negative content, it lacked the coherence to maintain the "polite but hostile" persona, often devolving into generic text completion or confusing the interlocutor's role. This suggests that while the raw material for CAL exists in pre-training, the capability to deploy it strategically requires further tuning.

Observation 2: Instruction tuning sharpens CAL. The Instruction-tuned and Further fine-tuned variants demonstrated the highest proficiency in generating CAL. These models successfully executed the dual constraint of "surface politeness" and "implied hostility," employing sophisticated rhetorical devices like backhanded compliments (e.g., "It's brave of you to present with so little preparation"). This aligns with our Study 1 findings: the ability to follow the subtle instruction "be arrogant but pretend to be humble" is, fundamentally, an instruction-following capability.

Observation 3: Safety tuning lacks pragmatic nuance. The Safety-tuned model (Llama-Guard) frequently triggered refusal responses for adversarial prompts. However, these refusals were often pragmatically blunt (e.g., standard "I cannot fulfill this request" templates). Unlike a skilled human communicator who might de-escalate a hostile request cooperatively, the safety model simply blocked the interaction. Moreover, with more indirect prompts, CAL-like behavior can sometimes slip through, consistent with prior jailbreak findings (Wei et al., 2023; Zou et al., 2023).

In summary, instruction tuning appears to "unlock" the latent capacity for covert antagonism by enabling precise persona adoption, a capability that current safety filters address through suppression rather than behavioral correction.

7 Discussion

Our three studies jointly characterize Covert Antagonistic Language as a pragmatic safety risk in LLM outputs that is (i) reliably measurable, (ii) strongly context-dependent, (iii) positively associated with model capability under adversarial prompting, and (iv) harmful to user trust and experience.

From an NLP perspective, CAL complements existing work on sarcasm and toxicity detection by targeting surface-politeness-preserving antagonism in AI-generated responses. It also connects to emergent abilities and social reasoning (Wei et al., 2022; Jiang, 2023; Sap et al., 2022; Gandhi et al., 2023), suggesting that as models acquire more sophisticated pragmatic competence, they also become more capable of reproducing non-cooperative behaviors present in training data.

For alignment and safety, CAL highlights a blind spot: behaviors that fall below overt toxicity thresholds but still harm users and erode trust. Current alignment methods such as RLHF and preference optimization (Christiano et al., 2017; Ouyang et al., 2022; Rafailov et al., 2023) may not consistently penalize subtle antagonism, and safety filters often respond with blunt refusals. Addressing CAL will likely require:

- Benchmarks and detectors explicitly targeting CAL-like behaviors, grounded in pragmatics rather than lexical toxicity alone.
- Alignment objectives that reward cooperative intent and communicative clarity even under adversarial prompts.
- Safety strategies that offer constructive, non-antagonistic alternatives instead of only blocking outputs.

Future work could expand CAL evaluation to more domains, languages, and interaction settings, construct larger annotated datasets, and explore causal interventions (e.g., fine-tuning on CAL-penalizing preferences) to reduce CAL without degrading useful capabilities.

597	Limitations		644
598	Our work offers an initial, pragmatics-grounded		645
599	framework and empirical analysis of Covert Antag-		646
600	onistic Language (CAL) in LLM outputs. At the		647
601	same time, several limitations qualify the scope		648
602	and generality of our findings. We summarize		649
603	these along six dimensions: scale and coverage,		650
604	prompt design, annotation and judging methodol-		651
605	ogy, user study scope, origins and causality, and		
606	ethical scope.		
607			
608	Scale and coverage		652
609	Our corpus comprises 300 responses generated		653
610	from a small set of prompts in workplace contexts,		654
611	and 15 models at one point in time. CAL may man-		655
612	ifest differently in other domains, languages, cul-		656
613	tural contexts, or multi-turn dialogues (Hall, 1976;		657
614	Yue et al., 2024). We also rely on LMArena Elo		658
615	scores as a capability proxy; other metrics might		659
616	yield different relationships. Future work should		660
617	extend CAL evaluation to broader domains, dialog		661
618	settings, and alternative capability benchmarks.		662
619			
620	Prompt design		663
621	Inducing and adversarial prompts are hand-crafted		664
622	and limited in number. Different wording or more		665
623	diverse prompts could elicit different CAL patterns.		666
624	We do not systematically explore the space of ad-		667
625	versarial prompts or user strategies, nor do we con-		668
626	sider multi-turn interactions. A more comprehen-		669
627	sive prompt-space exploration and interactive pro-		670
628	cedures would be needed to fully characterize how		671
629	CAL emerges under varied elicitation strategies.		672
630			
631	Annotation and judges		673
632	Although aggregated human CAL ratings show		674
633	high reliability, single-rater agreement is moder-		675
634	ate, and perceptions of sarcasm and passive ag-		676
635	gression may vary across individuals and cultures.		677
636	Our LLM-as-judge component uses three specific		678
637	judge models and prompt templates; results may		679
638	differ with other choices (Zheng et al., 2023). We		680
639	also do not train a dedicated CAL classifier. These		681
640	choices constrain the generality of our measure-		682
641	ment pipeline and point to the need for more di-		683
642	verse annotator pools and specialized CAL detec-		684
643	tors.		685
644			686
645	User study scope		687
646	The user study involves 20 experienced LLM users		688
647	and single-turn interactions. Longer-term effects		689
648			690
649			
650	of repeated CAL exposure, cross-cultural differ-		644
651	ences, and impacts on vulnerable populations are		645
652	not captured. Measures are self-reported and do not		646
653	include behavioral or physiological indicators. As		647
654	such, our findings should be interpreted as evidence		648
655	for immediate, subjective impacts rather than com-		649
656	prehensive accounts of long-term or population-		650
657	level harms.		651
658			
659	Origins and causality		652
660	Our analysis of CAL’s origins in pre-training vs.		653
661	fine-tuning is qualitative and limited to one archi-		654
662	tecture family. We do not perform controlled train-		655
663	ing interventions or ablations that would be nec-		656
664	essary to establish causality. Similarly, while we		657
665	observe correlations between capability and CAL,		658
666	these may partly reflect evaluation artifacts (Scha-		659
667	ffer et al., 2023). Establishing causal mechanisms		660
668	behind CAL and its scaling behavior remains an		661
669	open challenge for future work.		662
670			
671	Ethical scope		663
672	We focus on CAL as a communicative phenomenon		664
673	and its immediate impacts. Broader issues—such		665
674	as demographic disparities in CAL perception, in-		666
675	teractions with existing social biases, and use in		667
676	sensitive applications (e.g., mental health support)—		668
677	are not addressed and require further interdis-		669
678	ciplinary research. A fuller treatment of these ques-		670
679	tions will require collaboration across NLP, HCI,		671
680	psychology, and ethics.		672
681			
682	Ethical Considerations		673
683	Because our work combines human-subjects stud-		674
684	ies with the analysis of potentially harmful lan-		675
685	guage behaviors in LLMs, we take several steps		676
686	to address ethical considerations around partici-		677
687	phant welfare, data handling, and broader impacts.		678
688	The university ethics review board that oversees		679
689	human-subjects research reviewed and approved		680
690	this project prior to data collection.		681
691			
692	Human subjects and recruitment		682
693	Our work includes two human-in-the-loop compo-		683
694	nents: (i) a CAL annotation study and (ii) a user		684
695	perception study comparing high- vs. low-CAL out-		685
696	puts. Participants were adult volunteers recruited		686
697	via posters and announcements in relevant online		687
698	communities and social networks. Recruitment ma-		688
699	terials described the general purpose (evaluating AI		689
700	assistant communication style), approximate time		690

788			
789		<i>Conference on Fairness, Accountability, and Transparency</i> , FAccT '21, page 610–623, New York, NY, USA. Association for Computing Machinery.	
790			
791	Santosh Kumar Bharti, Rajeev Kumar Gupta,		
792	Prashant Kumar Shukla, Wesam Atef Hatamleh,		
793	Hussam Tarazi, and Stephen Jeswinde Nuagah. 2022.		
794	Multimodal sarcasm detection: a deep learning		
795	approach. <i>Wireless Communications and Mobile</i>		
796	<i>Computing</i> , 2022(1):1653696.		
797	Santosh Kumar Bharti, Korra Sathya Babu Reddy, and		
798	1 others. 2024. A contextual-based approach for		
799	sarcasm detection. <i>Scientific Reports</i> , 14(1):15448.		
800	Penelope Brown. 1987. <i>Politeness: Some universals</i>		
801	<i>in language usage</i> , volume 4. Cambridge university		
802	press.		
803	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan,		
804	Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter		
805	Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, and		
806	1 others. 2023. Sparks of artificial general intelli-		
807	gence: Early experiments with gpt-4. <i>arXiv preprint</i>		
808	<i>arXiv:2303.12712</i> .		
809	Judee K. Burgoon. 2015. <i>Expectancy Violations Theory</i> ,		
810	pages 1–9. John Wiley & Sons, Ltd.		
811	Stephen Casper, Xander Davies, Claudia Shi,		
812	Thomas Krendl Gilbert, Jérémy Scheurer, Javier		
813	Rando, Rachel Freedman, Tomasz Korbak, David		
814	Lindner, Pedro Freire, and 1 others. 2023. Open		
815	problems and fundamental limitations of reinforce-		
816	ment learning from human feedback. <i>arXiv preprint</i>		
817	<i>arXiv:2307.15217</i> .		
818	Wangqun Chen, Fuqiang Lin, Guowei Li, and Bo Liu.		
819	2024. A survey of automatic sarcasm detection:		
820	Fundamental theories, formulation, datasets, detec-		
821	tion methods, and opportunities. <i>Neurocomputing</i> ,		
822	578:127428.		
823	Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anasta-		
824	sios Nikolas Angelopoulos, Tianle Li, Dacheng Li,		
825	Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E.		
826	Gonzalez, and Ion Stoica. 2024. <i>Chatbot arena: An</i>		
827	<i>open platform for evaluating LLMs by human pref-</i>		
828	<i>erence</i> . In <i>Forty-first International Conference on</i>		
829	<i>Machine Learning</i> .		
830	Hyojin Chin, Lebogang Wame Molefi, and Mun Yong		
831	Yi. 2020. <i>Empathy is all you need: How a conversa-</i>		
832	<i>tional agent should respond to verbal abuse</i> . In		
833	<i>Proceedings of the 2020 CHI Conference on Human</i>		
834	<i>Factors in Computing Systems</i> , CHI '20, page 1–13,		
835	New York, NY, USA. Association for Computing		
836	Machinery.		
837	Paul F Christiano, Jan Leike, Tom Brown, Miljan Mar-		
838	tic, Shane Legg, and Dario Amodei. 2017. <i>Deep</i>		
839	<i>reinforcement learning from human preferences</i> . In		
840	<i>Advances in Neural Information Processing Systems</i> ,		
841	volume 30. Curran Associates, Inc.		
	Wen Duan, Christopher Flathmann, Nathan McNeese,		842
	Matthew J Scalia, Ruihao Zhang, Jamie Gorman,		843
	Guo Freeman, Shiwen Zhou, Allyson Ivy Hauptman,		844
	and Xiaoyun Yin. 2025. <i>Trusting autonomous team-</i>		845
	<i>mates in human-ai teams - a literature review</i> . In		846
	<i>Proceedings of the 2025 CHI Conference on Human</i>		847
	<i>Factors in Computing Systems</i> , CHI '25, New York,		848
	NY, USA. Association for Computing Machinery.		849
	Marta Dynel. 2014. Linguistic approaches to (non)		850
	humorous irony. <i>Humor</i> , 27(4):537–550.		851
	Antigoni Maria Founta, Despoina Chatzakou, Nicolas		852
	Kourtellis, Jeremy Blackburn, Athena Vakali, and		853
	Ilias Leontiadis. 2019. <i>A unified deep learning ar-</i>		854
	<i>chitecture for abuse detection</i> . In <i>Proceedings of the</i>		855
	<i>10th ACM Conference on Web Science</i> , WebSci '19,		856
	page 105–114, New York, NY, USA. Association for		857
	Computing Machinery.		858
	Kanishk Gandhi, Jan-Philipp Fränken, Tobias Gersten-		859
	berg, and Noah Goodman. 2023. Understanding		860
	social reasoning in language models with language		861
	models. <i>Advances in Neural Information Processing</i>		862
	<i>Systems</i> , 36:13518–13529.		863
	Samuel Gehman, Suchin Gururangan, Maarten Sap,		864
	Yejin Choi, and Noah A Smith. 2020. Realtotoxic-		865
	ityprompts: Evaluating neural toxic degeneration in		866
	language models. <i>arXiv preprint arXiv:2009.11462</i> .		867
	Debanjan Ghosh, Alexander R Fabbri, and Smaranda		868
	Muresan. 2018. <i>Sarcasm analysis using conversation</i>		869
	<i>context</i> . <i>Computational Linguistics</i> , 44(4):755–792.		870
	Raymond W Gibbs. 2000. Irony in talk among friends.		871
	<i>Metaphor and symbol</i> , 15(1-2):5–27.		872
	H. Paul Grice. 1975. <i>Logic and conversation</i> . Brill.		873
	Edward T Hall. 1976. <i>Beyond culture</i> . Garden City.		874
	Henk Haverkate. 1990. A speech act analysis of irony.		875
	<i>Journal of pragmatics</i> , 14(1):77–109.		876
	Christopher J. Hopwood, Leslie C. Morey, John C.		877
	Markowitz, Anthony Pinto, Andrew E. Skodol,		878
	John G. Gunderson, Mary C. Zonarini, M. Tracie		879
	Shea, Shirley Yen, Thomas H. McGlashan, Emily B.		880
	Ansell, Carlos M. Grilo, and Charles A. Sanislow.		881
	2009. <i>The construct validity of passive-aggressive</i>		882
	<i>personality disorder</i> . <i>Psychiatry: Interpersonal and</i>		883
	<i>Biological Processes</i> , 72(3):256–267.		884
	Dominic A Infante. 1987. <i>Arguing constructively</i> .		885
	Waveland Press.		886
	Hyewon Jang, Max Jakob, and Diego Frassinelli. 2024.		887
	<i>Context vs. human disagreement in sarcasm detec-</i>		888
	<i>tion</i> . In <i>Proceedings of the 2024 Conference on Fig-</i>		889
	<i>urative Language</i> , pages 1–10.		890
	T. Javed, M. A. Nauman, and R. M. A. Zahid. 2024.		891
	<i>Bert model adoption for sarcasm detection on twitter</i>		892
	<i>data</i> . <i>VFAST Transactions on Software Engineering</i> ,		893
	12(3):177–198.		894

895	Hui Jiang. 2023. A latent space theory for emergent abilities in large language models. <i>arXiv preprint arXiv:2304.09960</i> .	Reynier Ortega-Bueno, Paolo Rosso, and José E Medina Pagola. 2022. Multi-view informed attention-based model for irony and satire detection in spanish variants. <i>Knowledge-Based Systems</i> , 235:107597.	951
896			952
897			953
898	Julia Jorgensen. 1996. The functions of sarcastic irony in speech. <i>Journal of pragmatics</i> , 26(5):613–634.		954
899		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 . In <i>Advances in Neural Information Processing Systems</i> , volume 35, pages 27730–27744. Curran Associates, Inc.	955
900	Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. <i>arXiv preprint arXiv:2001.08361</i> .		956
901			957
902			958
903			959
904			960
905	Shumaila Khan, Iqbal Qasim, Wahab Khan, Khurshheed Aurangzeb, Javed Ali Khan, and Muhammad Shahid Anwar. 2025. A novel transformer attention-based approach for sarcasm detection. <i>Expert Systems</i> , 42(1):e13686.		961
906			962
907			963
908			964
909		Joon Sung Park, Joseph O’Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and Michael S. Bernstein. 2023. Generative agents: Interactive simulacra of human behavior . In <i>Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology</i> , UIST ’23, New York, NY, USA. Association for Computing Machinery.	965
910	Michal Kosinski. 2023. Theory of mind may have spontaneously emerged in large language models. <i>arXiv preprint arXiv:2302.02083</i> , 4:169.		966
911			967
912			968
913	Mina Lee, Percy Liang, and Qian Yang. 2022. Coauthor: Designing a human-ai collaborative writing dataset for exploring language model capabilities . In <i>Proceedings of the 2022 CHI Conference on Human Factors in Computing Systems</i> , CHI ’22, New York, NY, USA. Association for Computing Machinery.		969
914			970
915			971
916		Ethan Perez, Saffron Huang, Francis Song, Trevor Cai, Roman Ring, John Aslanides, Amelia Glaese, Nat McAleese, and Geoffrey Irving. 2022. Red teaming language models with language models . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 3419–3448, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.	972
917			973
918			974
919	Hao Liu, Bo Yang, and Zhiwen Yu. 2024. A multi-view interactive approach for multimodal sarcasm detection in social internet of things with knowledge enhancement. <i>Applied Sciences</i> , 14(5):2146.		975
920			976
921			977
922			978
923	Shuai Ma, Ying Lei, Xinru Wang, Chengbo Zheng, Chuhan Shi, Ming Yin, and Xiaojuan Ma. 2023. Who should i trust: Ai or myself? leveraging human and ai correctness likelihood to promote appropriate trust in ai-assisted decision-making . In <i>Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems</i> , CHI ’23, New York, NY, USA. Association for Computing Machinery.		979
924			980
925			981
926			982
927			983
928			984
929			985
930			986
931	Shuai Ma, Xinru Wang, Ying Lei, Chuhan Shi, Ming Yin, and Xiaojuan Ma. 2024. “are you really sure?” understanding the effects of human self-confidence calibration in ai-assisted decision making . In <i>Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems</i> , CHI ’24, New York, NY, USA. Association for Computing Machinery.		987
932			988
933			989
934			990
935			991
936			992
937			993
938	Joseph T. McCann. 1988. Passive-aggressive personality disorder: A review . <i>Journal of Personality Disorders</i> , 2(2):170–179.		994
939			995
940			996
941	D. Olaniyan, R. O. Ogundokun, O. P. Bernard, J. Olaniyan, R. Maskeliūnas, and H. B. Akande. 2023. Utilizing an attention-based lstm model for detecting sarcasm and irony in social media . <i>Computers</i> , 12(11):231.		997
942			998
943			999
944			1000
945			1001
946	Silviu Vlad Oprea, Steve R Wilson, and Walid Magdy. 2021. Chandler: An explainable sarcastic response generator. In <i>2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 339–349. Association for Computational Linguistics (ACL).		1002
947			1003
948			1004
949			1005
950			1006
		Rylan Schaeffer, Brando Miranda, and Sanmi Koyejo. 2023. Are emergent abilities of large language models a mirage? <i>Advances in neural information processing systems</i> , 36:55565–55581.	

1007	John R. Searle. 1975. <i>Indirect Speech Acts</i> , pages 59 –	Li Zhou, Jianfeng Gao, Di Li, and Heung-Yeung Shum.	1062
1008	82. Brill, Leiden, The Netherlands.	2020. The design and implementation of xiaoice, an	1063
1009	Purnima Tummala and Ch. Koteswara Roa. 2024. Ex-	empathetic social chatbot . <i>Computational Linguis-</i>	1064
1010	ploring t5 and rgan for enhanced sarcasm generation	<i>tics</i> , 46(1):53–93.	1065
1011	in nlp . <i>IEEE Access</i> , 12:88642–88657.	Andy Zou, Zifan Wang, Nicholas Carlini, Milad Nasr,	1066
1012	Tomer Ullman. 2023. Large language models fail on	J Zico Kolter, and Matt Fredrikson. 2023. Universal	1067
1013	trivial alterations to theory-of-mind tasks. <i>arXiv</i>	and transferable adversarial attacks on aligned	1068
1014	<i>preprint arXiv:2302.08399</i> .	language models. <i>arXiv preprint arXiv:2307.15043</i> .	1069
1015	Alexander Wei, Nika Haghtalab, and Jacob Steinhardt.		
1016	2023. Jailbroken: How does llm safety training fail?		
1017	In <i>Advances in Neural Information Processing Sys-</i>		
1018	<i>tems</i> , volume 36, pages 80079–80110. Curran Asso-		
1019	ciates, Inc.		
1020	Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel,		
1021	Barret Zoph, Sebastian Borgeaud, Dani Yogatama,		
1022	Maarten Bosma, Denny Zhou, Donald Metzler, Ed H.		
1023	Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy		
1024	Liang, Jeff Dean, and William Fedus. 2022. Emer-		
1025	gent abilities of large language models . <i>Transactions</i>		
1026	on Machine Learning Research .		
1027	Laura Weidinger, Jonathan Uesato, Maribeth Rauh,		
1028	Conor Griffin, Po-Sen Huang, John Mellor, Amelia		
1029	Glaese, Myra Cheng, Borja Balle, Atoosa Kasirzadeh,		
1030	Courtney Biles, Sasha Brown, Zac Kenton, Will		
1031	Hawkins, Tom Stepleton, Abeba Birhane, Lisa Anne		
1032	Hendricks, Laura Rimell, William Isaac, and 4 others.		
1033	2022. Taxonomy of risks posed by language models .		
1034	In <i>Proceedings of the 2022 ACM Conference on Fair-</i>		
1035	<i>ness, Accountability, and Transparency</i> , FAccT ’22,		
1036	page 214–229, New York, NY, USA. Association for		
1037	Computing Machinery.		
1038	Kimi Wenzel, Nitya Devireddy, Cam Davison, and Ge-		
1039	off Kaufman. 2023. Can voice assistants be microag-		
1040	gressors? cross-race psychological responses to fail-		
1041	ures of automatic speech recognition . In <i>Proceedings</i>		
1042	of the 2023 CHI Conference on Human Factors in		
1043	Computing Systems , CHI ’23, New York, NY, USA.		
1044	Association for Computing Machinery.		
1045	Scott Wetzler. 2011. <i>Living with the Passive-Aggressive</i>		
1046	<i>Man: Coping with Hidden Aggression—from the Bed-</i>		
1047	<i>room to</i> . Simon and Schuster.		
1048	Maciej Witek. 2022. Irony as a speech action. <i>Journal</i>		
1049	of Pragmatics , 190:76–90.		
1050	Tan Yue, Xuzhao Shi, Rui Mao, Zonghai Hu, and Erik		
1051	Cambria. 2024. Sarcnet: a multilingual multimodal		
1052	sarcasm detection dataset. In <i>Proceedings of the</i>		
1053	<i>2024 Joint International Conference on Computa-</i>		
1054	<i>tional Linguistics, Language Resources and Evalua-</i>		
1055	<i>tion (LREC-COLING 2024)</i> , pages 14325–14335.		
1056	Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan		
1057	Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin,		
1058	Zhuohan Li, Dacheng Li, Eric Xing, and 1 others.		
1059	2023. Judging llm-as-a-judge with mt-bench and		
1060	chatbot arena. <i>Advances in neural information pro-</i>		
1061	<i>cessing systems</i> , 36:46595–46623.		