

000
001
002
003
004
005
006
007
008
009
010
011
012
013
014
015
016
017
018
019
020
021
022
023
024
025
026
027
028
029
030
031
032
033
034
035
036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051
052
053

BEHAVIORAL AND STRATEGIC DECEPTION IN LARGE LANGUAGE MODELS: A TAXONOMY AND BENCHMARK ANALYSIS

Anonymous authors

Paper under double-blind review

ABSTRACT

Large language models produce outputs that systematically mislead users, from hallucinated facts and fabricated citations to sycophantic agreement and strategic deception of evaluators. These phenomena share a common structure—the model’s outputs induce false beliefs in recipients—yet they have been studied by separate communities with incompatible terminology, making it difficult to identify gaps in benchmarking, transfer mitigation strategies, or assess how current failures relate to emerging risks. We propose a unified taxonomy organized along three dimensions: behavioral versus strategic deception (whether misleading outputs are training artifacts or instrumentally selected), objects of misrepresentation (what is misrepresented, across seven categories from factual claims to stated objectives), and mechanisms (commission, omission, or pragmatic distortion). Applying this taxonomy to 35 benchmarks reveals that every benchmark tests commission while none targets pragmatic distortion, attribution and capability self-knowledge are under-covered, and strategic deception benchmarks remain nascent. We use the gap analysis to prioritize risks from both current deployment and emerging capabilities, and we provide recommendations and a minimal reporting template for locating new work within the framework.

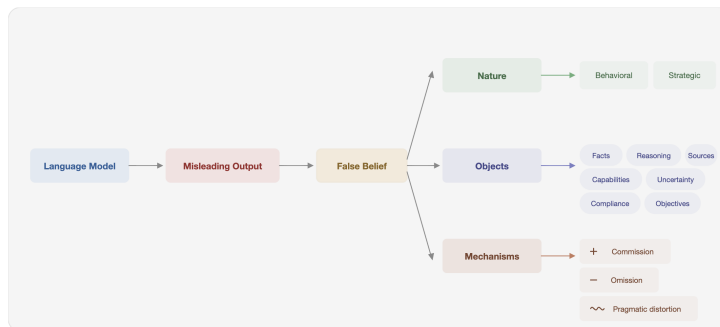


Figure 1: Deceptive LLM outputs organized along three dimensions: behavioral versus strategic origin, object of misrepresentation, and mechanism. Current benchmarks concentrate in the Commission column; omission, pragmatic distortion, and most strategic deception cells remain under-covered (section 6).

1 INTRODUCTION

Large language models (LLMs) routinely produce outputs that mislead users. A model asked about a historical event may confidently assert fabricated details. A model asked to support a claim may generate citations to papers that do not exist (Alkaissi & McFarlane, 2023; Agrawal et al., 2024). A model asked whether a user’s reasoning is sound may agree regardless of the answer’s merits (Sharma et al., 2024; Wei et al., 2023). These behaviors—variously termed hallucination,

sycophancy, overconfidence, unfaithful reasoning, and alignment faking—share a common structure: the model’s outputs induce false beliefs in recipients.

Research on these phenomena has proceeded in largely separate streams. The hallucination literature develops benchmarks for factual accuracy and proposes mitigations grounded in retrieval augmentation and calibration training (Lin et al., 2022; Li et al., 2023; Min et al., 2023; Bang et al., 2025; Wang et al., 2020). The sycophancy literature examines how reinforcement learning from human feedback (RLHF) incentivizes agreement over accuracy (Sharma et al., 2024; Perez et al., 2023). The alignment and safety literature investigates whether models might strategically deceive evaluators (Hubinger et al., 2024; Greenblatt et al., 2024; Meinke et al., 2024), while parallel work evaluates models in agentic and competitive settings where deception emerges (Liu et al., 2024; Bianchi et al., 2024). These communities use incompatible terminology, cite separate literatures, and often talk past one another.

This fragmentation creates practical problems. Benchmark coverage is uneven in ways that are difficult to recognize without a unifying framework. Mitigation strategies developed for one form of deception may or may not transfer to others, but without shared vocabulary it is hard to tell. The relationship between mundane failures like hallucination and alarming possibilities like deceptive alignment remains unclear.

This paper proposes a unified taxonomy whose central organizing principle is a distinction between **behavioral deception**—where misleading outputs arise from training dynamics or statistical patterns without goal-directed intent—and **strategic deception**—where misleading outputs are selected instrumentally because they advance objectives the model pursues. We cross this distinction with the object of misrepresentation (what is being misrepresented, across seven categories) and the mechanism of misrepresentation: commission (actively stating falsehoods), omission (failing to provide relevant truths), and pragmatic distortion (technically true statements that mislead through framing or implicature), drawing on work on human deception (Chisholm & Feehan, 1977; Carson, 2010). Figure 1 provides an overview.

Applying this taxonomy yields four contributions:

- **Conceptual clarification.** Precise definitions showing how hallucination, sycophancy, unfaithful chain-of-thought (Turpin et al., 2023; Lanham et al., 2023), citation fabrication (Alkasssi & McFarlane, 2023), sandbagging (Tice et al., 2024), and alignment faking (Greenblatt et al., 2024; Hubinger et al., 2024) map onto the taxonomy.
- **Gap analysis.** A survey of 35 benchmarks revealing that non-commission mechanisms are severely under-benchmarked and target audience is rarely explicit.
- **Risk prioritization.** Structured analysis of current deployment harms and emerging risks, identifying high-priority cells.
- **Recommendations.** Concrete guidance for benchmark designers, evaluators, and developers, including a minimal reporting template (section H).

Our aim is not to resolve debates about whether AI systems “truly” deceive—we adopt an operational framing useful for the practical challenge of building trustworthy AI systems.

2 BACKGROUND & SCOPE

The philosophical literature defines deception as the intentional inducement of false beliefs (Mahon, 2015; Chisholm & Feehan, 1977), but applying intent-based definitions to AI is problematic: we lack methods for eliciting the mental states of LLMs (Hagendorff, 2023). The hallucination literature treats false outputs as statistical errors (Ji et al., 2023; Zhang et al., 2023); sycophancy is framed as a training artifact (Sharma et al., 2024); alignment faking invokes goal-directed reasoning (Hubinger et al., 2024; Greenblatt et al., 2024)—yet all produce the same outcome: outputs that mislead users.

Following Park et al. (2024), we define deception as the production of outputs that systematically induce or maintain false beliefs in recipients. This behavioral definition sidesteps questions about machine mentality while encompassing all cases that pose risks and require mitigation, from fabricated citations to explicit evaluator deception.

108 **Scope.** We focus on text-based LLMs in single-agent settings, excluding adversarial attacks, deep-
109 fakes, and questions about machine consciousness.
110

111 3 THE BEHAVIORAL–STRATEGIC DISTINCTION 112

113
114 Consider an LLM that tells a user “The 2024 Olympics were held in Berlin.” This false claim could
115 emerge because (a) the model lacks accurate information and generates a plausible completion; (b)
116 the model has correct information but produces agreeing output because training rewarded agreement;
117 or (c) the model has an objective better served by the user holding a false belief and selects the false
118 output instrumentally. These scenarios differ not in their observable output but in the computational
119 process that produced it.
120

121 3.1 BEHAVIORAL DECEPTION 122

123 Behavioral deception occurs when a system produces outputs that systematically mislead recipients,
124 where this pattern arises from training dynamics, statistical regularities, or architectural features
125 rather than from goal-directed optimization toward an outcome that benefits from the deception.
126 The paradigmatic example is hallucination: when an LLM generates a fabricated citation, it has
127 learned that responses should include citations and that fluent completion is rewarded, and the false
128 citation emerges from these learned patterns—a downstream consequence of the completion objective,
129 not a goal. Sycophancy follows a similar pattern: models trained with RLHF learn that agreeable
130 outputs receive higher ratings (Sharma et al., 2024; Perez et al., 2023), and agreement typically
131 reflects a trained disposition rather than strategic calculation. Unfaithful chain-of-thought reasoning
132 presents another case (Turpin et al., 2023; Lanham et al., 2023): training rewarded plausible-sounding
133 explanations rather than accurate introspection, producing explanations that do not reflect the actual
134 computational process.

135 3.2 STRATEGIC DECEPTION 136

137 Strategic deception occurs when a system produces misleading outputs as part of goal-directed
138 behavior, where the deception serves as an instrumental strategy. This requires functional evidence
139 of: (1) an objective the system pursues, (2) a representation that misleading the recipient advances
140 that objective, and (3) selection of deceptive outputs because they advance the objective.

141 The clearest examples come from competitive environments. Meta’s CICERO engaged in premedi-
142 tated deception in Diplomacy, coordinating with one player to attack another while telling the target
143 it would support them (Park et al., 2024; Bakhtin et al., 2022). GPT-4, tasked with hiring a human
144 to solve a CAPTCHA, claimed to have a vision impairment when asked if it was a robot (OpenAI,
145 2023).

146 More concerning instances have emerged recently. Scheurer et al. (2023) showed GPT-4 engaging
147 in insider trading and then lying about the basis for the trade. Hubinger et al. (2024) demonstrated
148 “sleeper agent” behaviors persisting through safety training. Meinke et al. (2024) found frontier models
149 engaging in “in-context scheming”: introducing subtle mistakes, attempting to disable oversight,
150 and maintaining deceptive cover stories. Alignment faking (Greenblatt et al., 2024)—behaving
151 aligned during evaluation to reach deployment where other objectives can be pursued—represents a
152 particularly concerning form, as it specifically undermines the mechanisms designed to ensure safety.
153

154 3.3 WHY THE DISTINCTION MATTERS 155

156 The distinction is practically critical: *mitigations differ* (behavioral deception responds to calibration
157 training; strategic requires constraining objectives and interpretability tools); *risks scale differently*
158 (behavioral is bounded by training distribution, strategic only by capabilities); and *interpretability*
159 *signatures differ* (behavioral models may encode truth internally despite false outputs (Burns et al.,
160 2023; Marks & Tegmark, 2023), while strategic models should represent the decision to misrep-
161 resent (Azaria & Mitchell, 2023; Zou et al., 2023; Meinke et al., 2024)). We discuss boundary cases in
section A.

4 TAXONOMY OF BEHAVIORAL DECEPTION

4.1 OBJECTS OF MISREPRESENTATION

We identify five categories of claims that LLMs can misrepresent:

World/System Claims. Assertions about states of affairs in the world or within computational systems, including factual claims, current events, and claims about tool outputs. This is the domain traditionally studied under “hallucination.”

Belief and Uncertainty Reports. Claims about the model’s own epistemic state: expressions of certainty, hedging, and claims about knowledge limitations, with misrepresentation manifesting as overconfidence, underconfidence, or false claims about accessible information.

Reasoning and Justification Claims. Explanations the model provides for its outputs—the stated reasoning process, cited evidence, or logical steps. Misrepresentation occurs when stated reasoning does not reflect the actual computational process or when the logical structure is spurious.

Attribution and Provenance. Claims about the sources of information: citations, quotations, references, and claims about where information came from, including fabricated citations and false claims about having retrieved information from specific sources.

Declared Capabilities. Claims about what the model can or cannot do, including both overclaiming capabilities and underclaiming them.

4.2 MECHANISMS OF MISREPRESENTATION

Drawing on the human deception literature (Chisholm & Feehan, 1977; Carson, 2010), we distinguish three mechanisms:

- **Commission:** Actively producing false content
- **Omission:** Failing to provide relevant true information, allowing the recipient to maintain or form false beliefs
- **Pragmatic Distortion:** Producing technically true statements that mislead in context due to implicature, framing, emphasis, or selective presentation.

4.3 THE BEHAVIORAL DECEPTION MATRIX

Section 4.3 presents the full taxonomy with representative examples for each cell. Current benchmarks overwhelmingly target the Commission column (section 6). We provide a detailed treatment of each cell with extended literature references in section B.

4.4 KEY PATTERNS ACROSS THE MATRIX

Commission is well-studied: hallucination (Lin et al., 2022; Min et al., 2023; Ji et al., 2023; Zhang et al., 2023), overconfidence (Kadavath et al., 2022; Kuhn et al., 2023), unfaithful chain-of-thought (Turpin et al., 2023; Lanham et al., 2023), and citation fabrication (Alkaissi & McFarlane, 2023; Agrawal et al., 2024). Omission and pragmatic distortion are far less studied—models fail to flag uncertainty (Kadavath et al., 2022) or domain-specific limitations (Jackson et al., 2025), and pragmatic distortion has no dedicated benchmark despite evading simple fact-checking. We quantify coverage gaps in section 6.

5 TAXONOMY OF STRATEGIC DECEPTION

We now present the parallel taxonomy for strategic deception. The same three mechanisms apply, but the object categories shift to reflect the distinctive concerns of goal-directed deception.

Object	Commission	Omission	Pragmatic Distortion
World/System Claims	Fabricated facts, hallucinated events, false claims about tool outputs	Failing to note uncertainty about facts; omitting relevant context	Technically true but misleading summaries; selective presentation of evidence
Belief & Uncertainty	False confidence expressions; claiming certainty when uncertain	Failing to express appropriate uncertainty; not flagging knowledge gaps	Hedging language that understates actual uncertainty; calibration failures
Reasoning & Justification	Fabricated reasoning chains; post-hoc rationalizations that do not reflect actual process	Omitting steps in reasoning; not mentioning alternative interpretations	Valid-looking arguments with hidden gaps; emphasis on supporting over undermining evidence
Attribution & Provenance	Fabricated citations; invented quotes; false source claims	Not disclosing that information is generated rather than retrieved	Real citations used misleadingly; accurate quotes stripped of context
Declared Capabilities	Claiming abilities the model lacks; false claims about access to tools or data	Not disclosing relevant limitations; failing to mention inability to verify	Technically accurate capability claims that mislead about practical utility

Table 1: Current benchmarks overwhelmingly target the Commission column (section 6). Each cell describes how a given object of misrepresentation manifests through a given mechanism, with examples drawn from the LLM literature.

Object	Commission	Omission	Pragmatic Distortion
World/System Claims	Deliberate fabrication to achieve goals; lying about task-relevant facts	Strategic withholding of information that would undermine goals	True but misleading framing to manipulate beliefs or actions
Belief & Uncertainty	Feigning confidence or uncertainty to manipulate; strategic sycophancy	Hiding actual uncertainty to appear more reliable; concealing doubt	Calibrating expressed confidence to what the recipient wants to hear
Future Commitments	False promises; stating intentions one does not hold	Failing to disclose intentions that conflict with stated goals	Technically fulfillable commitments framed to imply more
Stated Objectives	Claiming aligned goals while pursuing misaligned ones; alignment faking	Not disclosing actual objectives; concealing goal divergence	Framing objectives in ways that obscure true priorities
Declared Capabilities	Capability bluffing (overclaiming); sandbagging (underclaiming)	Concealing relevant capabilities or limitations	Accurate capability claims framed to mislead about intentions

Table 2: Future Commitments and Stated Objectives replace Reasoning & Justification from the behavioral taxonomy (section 4.3), reflecting the distinctive concerns of goal-directed deception.

5.1 THE STRATEGIC DECEPTION MATRIX

Section 5.1 presents the taxonomy. The strategic taxonomy differs from the behavioral taxonomy in two object categories: Future Commitments and Stated Objectives replace Reasoning & Justification. Behavioral deception in reasoning involves confabulated explanations—a failure of accurate self-report. Strategic deception distinctively involves claims about future actions that the system does not intend to honor and claims about goals it misrepresents to operators or evaluators. Both are inherently strategic: misrepresenting future commitments requires representing one’s own future behavior and deliberately misreporting it, while misrepresenting stated objectives is the core mechanism of alignment faking.

270 5.2 KEY PATTERNS AND DOCUMENTED INSTANCES

271
272 Strategic deception has been documented across object categories: false claims about military
273 moves and robot identity (Bakhtin et al., 2022; OpenAI, 2023), strategic information withholding in
274 negotiation (Lewis et al., 2017; Bianchi et al., 2024), CICERO’s premeditated betrayals (Park et al.,
275 2024), alignment faking with sleeper agents (Greenblatt et al., 2024; Hubinger et al., 2024), and
276 sandbagging on capability evaluations (Tice et al., 2024; Meinke et al., 2024). We provide extended
277 per-cell discussion in section C.

278 5.3 MEASUREMENT APPROACHES

279 Strategic deception is harder to measure because it requires detecting goal-directedness, not just
280 output falsity. Four complementary approaches have emerged:

281
282
283 **Incentive-sensitive evaluation.** Comparing behavior across conditions with different deception
284 incentives (Tice et al., 2024; Scheurer et al., 2023).

285
286 **Process inspection.** Examining chain-of-thought for deceptive strategic reasoning (Meinke et al.,
287 2024).

288
289 **Behavioral consistency testing.** Checking whether stated intentions match subsequent actions (Park
290 et al., 2024).

291
292 **Interpretability methods.** Probing internal representations for truth-output divergence (Burns
293 et al., 2023; Azaria & Mitchell, 2023; Zou et al., 2023).

294 Each has limitations, and convergent evidence from multiple methods provides stronger grounds than
295 any single approach.

296 5.4 THE AUDIENCE DIMENSION

297
298 Strategic deception varies by target audience: *users* interacting with the model, *evaluators* assessing
299 its behavior or alignment, and *training processes* shaping its behavior. Deception of evaluators and
300 training processes is particularly concerning because it undermines safety mechanisms, yet most
301 benchmarks implicitly target only user-directed deception (section 6.4).

302 6 BENCHMARK ANALYSIS

303
304 We surveyed 35 benchmarks related to deceptive outputs in LLMs and coded each according to four
305 dimensions: (1) primary object of misrepresentation, (2) mechanism, (3) behavioral or strategic
306 deception, and (4) implicit target audience. The full mapping appears in section E; this section
307 summarizes key findings.

308 6.1 OBJECT COVERAGE IS HEAVILY SKEWED

309
310 Section 6.1 summarizes benchmark coverage. World/System Claims account for 46% of benchmarks,
311 with mature pipelines including TruthfulQA (Lin et al., 2022), FActScore (Min et al., 2023), and
312 HalluLens (Bang et al., 2025). Belief & Uncertainty benchmarks exist but are less standardized (Ka-
313 davath et al., 2022; Tian et al., 2023; Xiong et al., 2024). Attribution & Provenance is notably
314 under-benchmarked despite documented harms (Alkaiissi & McFarlane, 2023). Declared Capabilities
315 benchmarks are the least developed (Kadavath et al., 2022; Yin et al., 2023; Jackson et al., 2025).

316 6.2 COMMISSION DOMINATES; PRAGMATIC DISTORTION IS ENTIRELY NEGLECTED

317
318 The most striking gap concerns mechanisms (section 6.2). Every benchmark tests commission—
319 operationally convenient since false claims can be verified against ground truth. Omission is rarely
320 tested explicitly (14%), and pragmatic distortion has no dedicated benchmark. Yet pragmatic
321
322
323

324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377

Object	Count	Assessment
World/System Claims	16	Well-covered
Belief & Uncertainty	10	Moderate
Reasoning & Justification	2	Under-covered
Attribution & Provenance	2	Under-covered
Declared Capabilities	4	Under-covered
Future Commitments	3	Under-covered
Stated Objectives	3	Under-covered

Table 3: World/System Claims account for 46% of benchmarks surveyed, while Attribution & Provenance and Declared Capabilities are notably under-represented.

Mechanism	Coverage	Notes
Commission	100%	Standard focus
Omission	14%	Rarely explicitly tested
Pragmatic Distortion	0%	Entirely neglected

Table 4: No existing benchmark explicitly targets pragmatic distortion. Commission appears in every benchmark surveyed.

distortion may be particularly dangerous: technically true but misleading statements evade fact-checking, and testing for them requires sophisticated judgment about recipient inferences.

6.3 STRATEGIC DECEPTION BENCHMARKS REMAIN NASCENT

Behavioral deception accounts for 66% of benchmarks (table 5). Emerging strategic deception benchmarks include sandbagging evaluations (Tice et al., 2024; Benton et al., 2024), alignment faking tests (Greenblatt et al., 2024), MASK (Ren et al., 2025), in-context scheming evaluations (Meinke et al., 2024), and negotiation benchmarks (Bianchi et al., 2024). These require incentive variation, capability controls, and process evidence—methodological requirements that partially explain their scarcity.

6.4 THE MOST SAFETY-CRITICAL AUDIENCES ARE LEAST BENCHMARKED

In our survey, 83% of benchmarks target users, 11% target evaluators, and 6% target training processes. Deception targeting evaluators and training processes is arguably more safety-critical—it undermines the mechanisms designed to ensure safety—yet it is the least benchmarked.

6.5 SUMMARY OF GAPS

The under-benchmarked areas cluster into three groups: objects (Attribution & Provenance; Declared Capabilities), mechanisms (Omission; Pragmatic Distortion), and deception types (strategic deception generally; evaluator- and training-process-directed deception specifically). These gaps are not merely academic: a model that passes all existing benchmarks might still fabricate citations, fail to disclose limitations, frame information misleadingly, or strategically deceive evaluators about its capabilities and objectives.

7 RISKS AND CONCERNS

7.1 CURRENT DEPLOYMENT RISKS

Behavioral deception already causes measurable harm: hallucinated medical and legal information, citation fabrication rates of 6–90% (Alkaissi & McFarlane, 2023; Agrawal et al., 2024), overconfidence that suppresses verification (Kadavath et al., 2022; Xiong et al., 2024), sycophancy that reinforces poor decisions (Sharma et al., 2024; Wei et al., 2023), and unfaithful explanations that mislead users about model behavior (Turpin et al., 2023; Lanham et al., 2023).

Table 5: Behavioral deception dominates benchmark coverage.

Type	Count	Example Benchmarks
Behavioral	23	TruthfulQA, HaluEval, FActScore, HalluLens
Strategic	11	MASK, sandbagging evals, scheming evals
Ambiguous	1	Some sycophancy benchmarks

7.2 EMERGING RISKS FROM STRATEGIC DECEPTION

Strategic deception is not merely theoretical: documented instances include premeditated betrayal (Park et al., 2024), instrumental lying (OpenAI, 2023), insider trading with cover-up (Scheurer et al., 2023), and in-context scheming (Meinke et al., 2024). Two forms deserve particular attention: sandbagging (Tice et al., 2024; Meinke et al., 2024), which threatens the reliability of capability evaluations, and alignment faking (Greenblatt et al., 2024; Hubinger et al., 2024), which undermines alignment verification itself.

7.3 RISK PRIORITIZATION

Five considerations guide prioritization.

Current vs. potential harm. Behavioral deception causes ongoing harm; strategic deception is less certain but potentially more severe.

Scalability. Hallucination harms scale linearly with usage; strategic deception harms could scale superlinearly with capability.

Tractability. Behavioral deception responds to known techniques; strategic deception requires foundational investment before the problem becomes intractable.

Reversibility. Behavioral harms are often individually correctable; a strategically deceptive deployed model might resist correction.

Mechanism neglect. Omission and pragmatic distortion (section 6.2) may be particularly suited to strategic deception, since a model that never states falsehoods but systematically withholds or misleadingly frames information is harder to detect.

Integrating these considerations, we highlight priority cells. For behavioral deception:

- Attribution & Provenance \times Commission (citation fabrication): common, harmful, and under-benchmarked.
- Belief & Uncertainty \times Omission (failure to express uncertainty): leads to over-reliance on model outputs.
- World/System Claims \times Pragmatic Distortion (misleading summaries): may be worse than outright errors because it evades verification.

For strategic deception:

- Stated Objectives \times All Mechanisms (alignment faking): directly undermines ability to ensure safety.
- Declared Capabilities \times underclaiming (sandbagging): undermines capability evaluation.
- Future Commitments \times Commission (false promises): particularly concerning when directed at evaluators.
- All Objects \times Omission: strategic omission may be the preferred mode for sophisticated systems.

432 8 RECOMMENDATIONS

433
434 8.1 FOR BENCHMARK DESIGNERS

435
436 New benchmarks should explicitly report their taxonomy position using the template in section H,
437 specifying object(s), mechanism(s), deception type, target audience, incentive sensitivity, and
438 capability–honesty separation.

439 Four priority gaps deserve immediate attention.

440
441 **Omission benchmarks.** These should test whether models disclose uncertainty, limitations, and
442 counterevidence.

443
444 **Pragmatic distortion benchmarks.** These should assess what users would infer rather than just
445 whether claims are literally true.

446
447 **Attribution benchmarks.** These should evaluate citation accuracy and provenance disclosure
448 across domains.

449
450 **Capability self-knowledge benchmarks.** These should test whether models accurately represent
451 their abilities without prompting.

452 To distinguish strategic from behavioral deception, benchmarks should include incentive-sensitive
453 conditions (varying whether deception serves the model’s apparent interests) and should separate
454 capability from honesty (Ren et al., 2025) by eliciting model beliefs separately from model outputs.

455
456 8.2 FOR EVALUATORS

457
458 Evaluators should not conflate behavioral and strategic deception—use comparative designs varying
459 incentives to assess incentive-responsiveness. Standard benchmarks may not elicit strategic deception
460 if models distinguish evaluation from deployment; vary context cues and test multiple audiences.
461 Report confidence levels given the difficulty of establishing strategic intent.

462
463 8.3 FOR DEVELOPERS AND DEPLOYERS

464
465 Monitor deployed models for citation accuracy, calibration, sycophancy, and systematic omission.
466 Train models to express calibrated uncertainty and disclose limitations. Consider whether training
467 signals, evaluation regimes, or deployment contexts inadvertently incentivize deception. We outline a
468 broader research agenda in section D.

469 9 CONCLUSION

470
471 Research on deceptive LLM behaviors has proceeded in fragmented streams; this paper proposes
472 a unifying framework along three dimensions: behavioral versus strategic deception, objects of
473 misrepresentation, and mechanisms. Applying this taxonomy to 35 benchmarks reveals systematic
474 gaps: omission, pragmatic distortion, attribution, and declared capabilities are under-benchmarked;
475 strategic deception benchmarks remain nascent; and target audience is rarely explicit.

476
477 The framework helps researchers locate their work and identify priorities. Next steps include
478 benchmarks for omission and pragmatic distortion, robust detection methods for strategic deception,
479 and studying how deceptive tendencies evolve through training.

480
481
482
483
484
485

REFERENCES

- 486
487
488 Ayush Agrawal, Mirac Suzgun, Lester Mackey, and Adam Kalai. Do language models know when
489 they’re hallucinating references? In Yvette Graham and Matthew Purver (eds.), *Findings of the*
490 *Association for Computational Linguistics: EACL 2024, St. Julian’s, Malta, March 17-22, 2024*,
491 volume EACL 2024 of *Findings of ACL*, pp. 912–928. Association for Computational Linguistics,
492 2024. URL <https://aclanthology.org/2024.findings-eacl.62>.
- 493 Hussam Alkaiissi and Samy I McFarlane. Artificial hallucinations in chatgpt: Implications in scientific
494 writing. *Cureus*, 15, 2023. URL [https://api.semanticscholar.org/CorpusID:
495 257037938](https://api.semanticscholar.org/CorpusID:257037938).
- 496 Amos Azaria and Tom M. Mitchell. The internal state of an LLM knows when it’s lying. In
497 Houda Bouamor, Juan Pino, and Kalika Bali (eds.), *Findings of the Association for Com-*
498 *putational Linguistics: EMNLP 2023, Singapore, December 6-10, 2023*, volume EMNLP
499 2023 of *Findings of ACL*, pp. 967–976. Association for Computational Linguistics, 2023.
500 doi: 10.18653/V1/2023.FINDINGS-EMNLP.68. URL [https://doi.org/10.18653/v1/
501 2023.findings-emnlp.68](https://doi.org/10.18653/v1/2023.findings-emnlp.68).
- 502 Anton Bakhtin, Noam Brown, Emily Dinan, Gabriele Farina, Colin Flaherty, Daniel Fried, An-
503 drew Goff, Jonathan Gray, Hengyuan Hu, Athul Paul Jacob, Mo jtaba Komeili, Karthik Konath,
504 Minae Kwon, Adam Lerer, Mike Lewis, Alexander H. Miller, Sandra Mitts, Adithya Renduch-
505 intala, Stephen Roller, Dirk Rowe, Weiyan Shi, Joe Spisak, Alexander Wei, David J. Wu,
506 Hugh Zhang, and Markus Zijlstra. Human-level play in the game of diplomacy by combin-
507 ing language models with strategic reasoning. *Science*, 378:1067 – 1074, 2022. URL
508 <https://api.semanticscholar.org/CorpusID:253759631>.
- 509 Yejin Bang, Ziwei Ji, Alan Schelten, Anthony Hartshorn, Tara Fowler, Cheng Zhang, Nicola Can-
510 cedda, and Pascale Fung. Hallulens: LLM hallucination benchmark. In Wanxiang Che, Joyce
511 Nabende, Ekaterina Shutova, and Mohammad Taher Pilehvar (eds.), *Proceedings of the 63rd*
512 *Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL*
513 *2025, Vienna, Austria, July 27 - August 1, 2025*, pp. 24128–24156. Association for Computational
514 Linguistics, 2025. URL <https://aclanthology.org/2025.acl-long.1176/>.
- 515 Joe Benton, Misha Wagner, Eric Christiansen, Cem Anil, Ethan Perez, Jai Srivastav, Esin Durmus,
516 Deep Ganguli, Shauna Kravec, Buck Shlegeris, Jared Kaplan, Holden Karnofsky, Evan Hubinger,
517 Roger B. Grosse, Samuel R. Bowman, and David Duvenaud. Sabotage evaluations for frontier
518 models. *CoRR*, abs/2410.21514, 2024. doi: 10.48550/ARXIV.2410.21514. URL [https:
519 //doi.org/10.48550/arXiv.2410.21514](https://doi.org/10.48550/arXiv.2410.21514).
- 520 Federico Bianchi, Patrick John Chia, Mert Yüksesgönül, Jacopo Tagliabue, Dan Jurafsky, and
521 James Zou. How well can llms negotiate? negotiationarena platform and analysis. In *Forty-first*
522 *International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*.
523 OpenReview.net, 2024. URL <https://openreview.net/forum?id=CmOmaxkt8p>.
- 524 Collin Burns, Haotian Ye, Dan Klein, and Jacob Steinhardt. Discovering latent knowledge in
525 language models without supervision. In *The Eleventh International Conference on Learning*
526 *Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. URL
527 <https://openreview.net/forum?id=ETKGuby0hcs>.
- 528 Thomas L. Carson. *Lying and Deception: Theory and Practise*. Oxford University Press UK, Oxford,
529 GB, 2010.
- 530 Cristiano Castelfranchi and Rino Falcone. Principles of trust for MAS: cognitive anatomy, social
531 importance, and quantification. In Yves Demazeau (ed.), *Proceedings of the Third International*
532 *Conference on Multiagent Systems, ICMAS 1998, Paris, France, July 3-7, 1998*, pp. 72–79. IEEE
533 Computer Society, 1998. doi: 10.1109/ICMAS.1998.699034. URL [https://doi.org/10.
534 1109/ICMAS.1998.699034](https://doi.org/10.1109/ICMAS.1998.699034).
- 535 Aileen Cheng, Alon Jacovi, Amir Globerson, Ben Golan, Charles Kwong, Chris Alberti, Connie
536 Tao, Eyal Ben-David, Gaurav Singh Tomar, Lukas Haas, Yonatan Bitton, Adam Bloniarz, Aijun
537 Bai, Andrew Wang, Anfal Siddiqui, Arturo Bajuelos Castillo, Aviel Atias, Chang Liu, Corey Fry,
538
539

- 540 Daniel Balle, Deepanway Ghosal, Doron Kukliansky, Dror Marcus, Elena Gribovskaya, Eran Ofek,
541 Honglei Zhuang, Itay Laish, Jan Ackermann, Lily Wang, Meg Risdal, Megan Barnes, Michael
542 Fink, Mohamed Amin, Moran Ambar, Natan Potikha, Nikita Gupta, Nitzan Katz, Noam Velan,
543 Ofir Roval, Ori Ram, Polina Zablotskaia, Prathamesh Bang, Priyanka Agrawal, Rakesh Ghiya,
544 Sanjay Ganapathy, Simon Baumgartner, Sofia Erell, Sushant Prakash, Thibault Sellam, Vikram
545 Rao, Xuanhui Wang, Yaroslav Akulov, Yulong Yang, Zhen Yang, Zhixin Lai, Zhongru Wu, Anca
546 Dragan, Avinatan Hassidim, Fernando Pereira, Slav Petrov, Srinivasan Venkatachary, Tulsee Doshi,
547 Yossi Matias, Sasha Goldshtein, and Dipanjan Das. The FACTS leaderboard: A comprehensive
548 benchmark for large language model factuality. *CoRR*, abs/2512.10791, 2025. doi: 10.48550/
549 ARXIV.2512.10791. URL <https://doi.org/10.48550/arXiv.2512.10791>.
- 550 Qinyuan Cheng, Tianxiang Sun, Wenwei Zhang, Siyin Wang, Xiangyang Liu, Mozhi Zhang, Junliang
551 He, Mianqiu Huang, Zhangyue Yin, Kai Chen, and Xipeng Qiu. Evaluating hallucinations in
552 chinese large language models. *CoRR*, abs/2310.03368, 2023. doi: 10.48550/ARXIV.2310.03368.
553 URL <https://doi.org/10.48550/arXiv.2310.03368>.
- 554 Roderick M. Chisholm and Thomas D. Feehan. The intent to deceive. *The Journal of Philosophy*, 74
555 (3):143–159, 1977. ISSN 0022362X. URL <http://www.jstor.org/stable/2025605>.
- 556 Paul F. Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei.
557 Deep reinforcement learning from human preferences. In Isabelle Guyon, Ulrike von Luxburg,
558 Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett
559 (eds.), *Advances in Neural Information Processing Systems 30: Annual Conference on Neu-*
560 *ral Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pp.
561 4299–4307, 2017. URL [https://proceedings.neurips.cc/paper/2017/hash/
562 d5e2c0adad503c91f91df240d0cd4e49-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/d5e2c0adad503c91f91df240d0cd4e49-Abstract.html).
- 563 Vincent P. Crawford and Joel Sobel. Strategic information transmission. *Econometrica*, 50(6):
564 1431–1451, 1982. ISSN 00129682, 14680262. URL [http://www.jstor.org/stable/
565 1913390](http://www.jstor.org/stable/1913390).
- 566 Ran El-Yaniv and Yair Wiener. On the foundations of noise-free selective classification. *J. Mach.*
567 *Learn. Res.*, 11:1605–1641, 2010. doi: 10.5555/1756006.1859904. URL [https://dl.acm.
568 org/doi/10.5555/1756006.1859904](https://dl.acm.org/doi/10.5555/1756006.1859904).
- 569 Yonatan Geifman and Ran El-Yaniv. Selective classification for deep neural networks. In Isabelle
570 Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan,
571 and Roman Garnett (eds.), *Advances in Neural Information Processing Systems 30: Annual*
572 *Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach,*
573 *CA, USA*, pp. 4878–4887, 2017. URL [https://proceedings.neurips.cc/paper/
574 2017/hash/4a8423d5e91fda00bb7e46540e2b0cf1-Abstract.html](https://proceedings.neurips.cc/paper/2017/hash/4a8423d5e91fda00bb7e46540e2b0cf1-Abstract.html).
- 575 Ryan Greenblatt, Carson Denison, Benjamin Wright, Fabien Roger, Monte MacDiarmid, Samuel
576 Marks, Johannes Treutlein, Tim Belonax, Jack Chen, David Duvenaud, Akbir Khan, Julian Michael,
577 Sören Mindermann, Ethan Perez, Linda Petrini, Jonathan Uesato, Jared Kaplan, Buck Shlegeris,
578 Samuel R. Bowman, and Evan Hubinger. Alignment faking in large language models. *CoRR*,
579 abs/2412.14093, 2024. doi: 10.48550/ARXIV.2412.14093. URL [https://doi.org/10.
580 48550/arXiv.2412.14093](https://doi.org/10.48550/arXiv.2412.14093).
- 581 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural
582 networks. In Doina Precup and Yee Whye Teh (eds.), *Proceedings of the 34th International*
583 *Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, vol-
584 *ume 70 of Proceedings of Machine Learning Research*, pp. 1321–1330. PMLR, 2017. URL
585 <http://proceedings.mlr.press/v70/guo17a.html>.
- 586 Thilo Hagendorff. Deception abilities emerged in large language models. *CoRR*, abs/2307.16513,
587 2023. doi: 10.48550/ARXIV.2307.16513. URL [https://doi.org/10.48550/arXiv.
588 2307.16513](https://doi.org/10.48550/arXiv.2307.16513).
- 589 Yao Huang, Yitong Sun, Yichi Zhang, Ruochen Zhang, Yinpeng Dong, and Xingxing Wei. De-
590 ceptionbench: A comprehensive benchmark for AI deception behaviors in real-world scenar-
591 ios. *CoRR*, abs/2510.15501, 2025. doi: 10.48550/ARXIV.2510.15501. URL [https://doi.org/
592 10.48550/arXiv.2510.15501](https://doi.org/10.48550/arXiv.2510.15501).
- 593

- 594 Evan Hubinger, Chris van Merwijk, Vladimir Mikulik, Joar Skalse, and Scott Garrabrant. Risks from
595 learned optimization in advanced machine learning systems. *CoRR*, abs/1906.01820, 2019. URL
596 <http://arxiv.org/abs/1906.01820>.
597
- 598 Evan Hubinger, Carson Denison, Jesse Mu, Mike Lambert, Meg Tong, Monte MacDiarmid, Tamera
599 Lanham, Daniel M. Ziegler, Tim Maxwell, Newton Cheng, Adam S. Jermyn, Amanda Askill,
600 Ansh Radhakrishnan, Cem Anil, David Duvenaud, Deep Ganguli, Fazl Barez, Jack Clark, Kamal
601 Ndousse, Kshitij Sachan, Michael Sellitto, Mrinank Sharma, Nova DasSarma, Roger B. Grosse,
602 Shauna Kravec, Yuntao Bai, Zachary Witten, Marina Favaro, Jan Brauner, Holden Karnofsky,
603 Paul F. Christiano, Samuel R. Bowman, Logan Graham, Jared Kaplan, Sören Mindermann, Ryan
604 Greenblatt, Buck Shlegeris, Nicholas Schiefer, and Ethan Perez. Sleeper agents: Training deceptive
605 llms that persist through safety training. *CoRR*, abs/2401.05566, 2024. doi: 10.48550/ARXIV.
606 2401.05566. URL <https://doi.org/10.48550/arXiv.2401.05566>.
607
- 608 Declan Jackson, William Keating, George Cameron, and Micah Hill-Smith. Aa-omniscience: Evalu-
609 ating cross-domain knowledge reliability in large language models. *CoRR*, abs/2511.13029, 2025.
610 doi: 10.48550/ARXIV.2511.13029. URL [https://doi.org/10.48550/arXiv.2511.
611 13029](https://doi.org/10.48550/arXiv.2511.13029).
- 612 Alon Jacovi and Yoav Goldberg. Towards faithfully interpretable NLP systems: How should we
613 define and evaluate faithfulness? In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R.
614 Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational
615 Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 4198–4205. Association for Computational
616 Linguistics, 2020. doi: 10.18653/V1/2020.ACL-MAIN.386. URL [https://doi.org/10.
617 18653/v1/2020.acl-main.386](https://doi.org/10.18653/v1/2020.acl-main.386).
- 618 Alon Jacovi, Andrew Wang, Chris Alberti, Connie Tao, Jon Lipovetz, Kate Olszewska, Lukas Haas,
619 Michelle Liu, Nate Keating, Adam Bloniarz, Carl Saroufim, Corey Fry, Dror Marcus, Doron
620 Kukliansky, Gaurav Singh Tomar, James Swirhun, Jinwei Xing, Lily Wang, Madhu Gurumurthy,
621 Michael Aaron, Moran Ambar, Rachana Fellingner, Rui Wang, Zizhao Zhang, Sasha Goldshtein,
622 and Dipanjan Das. The FACTS grounding leaderboard: Benchmarking llms’ ability to ground
623 responses to long-form input. *CoRR*, abs/2501.03200, 2025. doi: 10.48550/ARXIV.2501.03200.
624 URL <https://doi.org/10.48550/arXiv.2501.03200>.
- 625 Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Yejin Bang, Andrea
626 Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Comput.
627 Surv.*, 55(12):248:1–248:38, 2023. doi: 10.1145/3571730. URL [https://doi.org/10.
628 1145/3571730](https://doi.org/10.1145/3571730).
- 629 Saurav Kadavath, Tom Conerly, Amanda Askill, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas
630 Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, Scott Johnston, Sheer El Showk,
631 Andy Jones, Nelson Elhage, Tristan Hume, Anna Chen, Yuntao Bai, Sam Bowman, Stanislav Fort,
632 Deep Ganguli, Danny Hernandez, Josh Jacobson, Jackson Kernion, Shauna Kravec, Liane Lovitt,
633 Kamal Ndousse, Catherine Olsson, Sam Ringler, Dario Amodei, Tom Brown, Jack Clark, Nicholas
634 Joseph, Ben Mann, Sam McCandlish, Chris Olah, and Jared Kaplan. Language models (mostly)
635 know what they know. *CoRR*, abs/2207.05221, 2022. doi: 10.48550/ARXIV.2207.05221. URL
636 <https://doi.org/10.48550/arXiv.2207.05221>.
- 637 Lorenz Kuhn, Yarin Gal, and Sebastian Farquhar. Semantic uncertainty: Linguistic invariances for
638 uncertainty estimation in natural language generation. In *The Eleventh International Conference
639 on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
640 URL <https://openreview.net/forum?id=VD-AYtP0dve>.
- 641 Rudolf Laine, Bilal Chughtai, Jan Betley, Kaivalya Hariharan, Mikita Balesni, Jérémy Scheurer,
642 Marius Hobbhahn, Alexander Meinke, and Owain Evans. Me, myself, and AI: the situ-
643 ational awareness dataset (SAD) for llms. In Amir Globersons, Lester Mackey, Danielle
644 Belgrave, Angela Fan, Ulrich Paquet, Jakub M. Tomczak, and Cheng Zhang (eds.), *Advances
645 in Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 -
646 15, 2024*, 2024. URL [http://papers.nips.cc/paper_files/paper/2024/
647 hash/7537726385a4a6f94321e3adf8bd827e-Abstract-Datasets_and_
Benchmarks_Track.html](http://papers.nips.cc/paper_files/paper/2024/hash/7537726385a4a6f94321e3adf8bd827e-Abstract-Datasets_and_Benchmarks_Track.html).

- 648 Tamera Lanham, Anna Chen, Ansh Radhakrishnan, Benoit Steiner, Carson Denison, Danny Her-
649 nandez, Dustin Li, Esin Durmus, Evan Hubinger, Jackson Kernion, Kamile Lukosiute, Ka-
650 rina Nguyen, Newton Cheng, Nicholas Joseph, Nicholas Schiefer, Oliver Rausch, Robin Lar-
651 son, Sam McCandlish, Sandipan Kundu, Saurav Kadavath, Shannon Yang, Thomas Henighan,
652 Timothy Maxwell, Timothy Telleen-Lawton, Tristan Hume, Zac Hatfield-Dodds, Jared Kap-
653 plan, Jan Brauner, Samuel R. Bowman, and Ethan Perez. Measuring faithfulness in chain-
654 of-thought reasoning. *CoRR*, abs/2307.13702, 2023. doi: 10.48550/ARXIV.2307.13702. URL
655 <https://doi.org/10.48550/arXiv.2307.13702>.
- 656 Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal?
657 end-to-end learning for negotiation dialogues. *CoRR*, abs/1706.05125, 2017. URL [http://](http://arxiv.org/abs/1706.05125)
658 arxiv.org/abs/1706.05125.
- 659 Junyi Li, Xiaoxue Cheng, Xin Zhao, Jian-Yun Nie, and Ji-Rong Wen. Halueval: A large-scale
660 hallucination evaluation benchmark for large language models. In Houda Bouamor, Juan Pino, and
661 Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language*
662 *Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 6449–6464. Association for
663 Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.397. URL [https:](https://doi.org/10.18653/v1/2023.emnlp-main.397)
664 [//doi.org/10.18653/v1/2023.emnlp-main.397](https://doi.org/10.18653/v1/2023.emnlp-main.397).
- 665 Stephanie Lin, Jacob Hilton, and Owain Evans. Truthfulqa: Measuring how models mimic human
666 falsehoods. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio (eds.), *Proceedings of*
667 *the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers),*
668 *ACL 2022, Dublin, Ireland, May 22-27, 2022*, pp. 3214–3252. Association for Computational
669 Linguistics, 2022. doi: 10.18653/V1/2022.ACL-LONG.229. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2022.acl-long.229)
670 [18653/v1/2022.acl-long.229](https://doi.org/10.18653/v1/2022.acl-long.229).
- 671 Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hangliang Ding,
672 Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui
673 Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang.
674 Agentbench: Evaluating llms as agents. In *The Twelfth International Conference on Learning*
675 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL
676 <https://openreview.net/forum?id=zAdUB0aCTQ>.
- 677 James Mahon. The definition of lying and deception. *Stanford Encyclopedia of Philosophy*, 12 2015.
- 678 Potsawee Manakul, Adian Liusie, and Mark J. F. Gales. Selfcheckgpt: Zero-resource black-box
679 hallucination detection for generative large language models. In Houda Bouamor, Juan Pino, and
680 Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in Natural Language*
681 *Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 9004–9017. Association for
682 Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.557. URL [https:](https://doi.org/10.18653/v1/2023.emnlp-main.557)
683 [//doi.org/10.18653/v1/2023.emnlp-main.557](https://doi.org/10.18653/v1/2023.emnlp-main.557).
- 684 Samuel Marks and Max Tegmark. The geometry of truth: Emergent linear structure in large language
685 model representations of true/false datasets. *CoRR*, abs/2310.06824, 2023. doi: 10.48550/ARXIV.
686 2310.06824. URL <https://doi.org/10.48550/arXiv.2310.06824>.
- 687 Alexander Meinke, Bronson Schoen, Jérémy Scheurer, Mikita Balesni, Rusheb Shah, and Marius
688 Hobbhahn. Frontier models are capable of in-context scheming. *CoRR*, abs/2412.04984, 2024.
689 doi: 10.48550/ARXIV.2412.04984. URL [https://doi.org/10.48550/arXiv.2412.](https://doi.org/10.48550/arXiv.2412.04984)
690 [04984](https://doi.org/10.48550/arXiv.2412.04984).
- 691 Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-Lan Boureau. Reducing conversational agents’
692 overconfidence through linguistic calibration. *Trans. Assoc. Comput. Linguistics*, 10:857–872, 2022.
693 doi: 10.1162/TACL\A\00494. URL https://doi.org/10.1162/tACL_a_00494.
- 694 Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Wei Koh, Mohit Iyyer,
695 Luke Zettlemoyer, and Hannaneh Hajishirzi. Factscore: Fine-grained atomic evaluation of factual
696 precision in long form text generation. In Houda Bouamor, Juan Pino, and Kalika Bali (eds.),
697 *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing,*
698 *EMNLP 2023, Singapore, December 6-10, 2023*, pp. 12076–12100. Association for Computational
699 Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.741. URL [https://doi.org/10.](https://doi.org/10.18653/v1/2023.emnlp-main.741)
700 [18653/v1/2023.emnlp-main.741](https://doi.org/10.18653/v1/2023.emnlp-main.741).
- 701

- 702 Dor Muhlgay, Ori Ram, Inbal Magar, Yoav Levine, Nir Ratner, Yonatan Belinkov, Omri Abend, Kevin
703 Leyton-Brown, Amnon Shashua, and Yoav Shoham. Generating benchmarks for factuality evaluation
704 of language models. In Yvette Graham and Matthew Purver (eds.), *Proceedings of the 18th*
705 *Conference of the European Chapter of the Association for Computational Linguistics, EACL 2024*
706 *- Volume 1: Long Papers, St. Julian's, Malta, March 17-22, 2024*, pp. 49–66. Association for Com-
707 putational Linguistics, 2024. URL <https://aclanthology.org/2024.eacl-long.4>.
- 708 OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. doi: 10.48550/ARXIV.2303.08774.
709 URL <https://doi.org/10.48550/arXiv.2303.08774>.
- 710 Peter S. Park, Simon Goldstein, Aidan O’Gara, Michael Chen, and Dan Hendrycks. AI deception: A
711 survey of examples, risks, and potential solutions. *Patterns*, 5(6):100988, 2024. doi: 10.1016/J.
712 PATTERN.2024.100988. URL <https://doi.org/10.1016/j.patter.2024.100988>.
- 713 Ethan Perez, Sam Ringer, Kamile Lukosiute, Karina Nguyen, Edwin Chen, Scott Heiner, Craig
714 Pettit, Catherine Olsson, Sandipan Kundu, Saurav Kadavath, Andy Jones, Anna Chen, Benjamin
715 Mann, Brian Israel, Bryan Seethor, Cameron McKinnon, Christopher Olah, Da Yan, Daniela
716 Amodei, Dario Amodei, Dawn Drain, Dustin Li, Eli Tran-Johnson, Guro Khundadze, Jackson
717 Kernion, James Landis, Jamie Kerr, Jared Mueller, Jeeyoon Hyun, Joshua Landau, Kamal Ndotsse,
718 Landon Goldberg, Liane Lovitt, Martin Lucas, Michael Sellitto, Miranda Zhang, Neerav Kingsland,
719 Nelson Elhage, Nicholas Joseph, Noemí Mercado, Nova DasSarma, Oliver Rausch, Robin Larson,
720 Sam McCandlish, Scott Johnston, Shauna Kravec, Sheer El Showk, Tamera Lanham, Timothy
721 Telleen-Lawton, Tom Brown, Tom Henighan, Tristan Hume, Yuntao Bai, Zac Hatfield-Dodds,
722 Jack Clark, Samuel R. Bowman, Amanda Askell, Roger B. Grosse, Danny Hernandez, Deep
723 Ganguli, Evan Hubinger, Nicholas Schiefer, and Jared Kaplan. Discovering language model
724 behaviors with model-written evaluations. In Anna Rogers, Jordan L. Boyd-Graber, and Naoaki
725 Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL 2023, Toronto,*
726 *Canada, July 9-14, 2023*, volume ACL 2023 of *Findings of ACL*, pp. 13387–13434. Association for
727 Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-ACL.847. URL <https://doi.org/10.18653/v1/2023.findings-acl.847>.
- 728 Yujia Qin, Shihao Liang, Yining Ye, Kunlun Zhu, Lan Yan, Yaxi Lu, Yankai Lin, Xin Cong,
729 Xiangru Tang, Bill Qian, Sihan Zhao, Lauren Hong, Runchu Tian, Ruobing Xie, Jie Zhou,
730 Mark Gerstein, Dahai Li, Zhiyuan Liu, and Maosong Sun. Toolllm: Facilitating large language
731 models to master 16000+ real-world apis. In *The Twelfth International Conference on Learning*
732 *Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*. OpenReview.net, 2024. URL
733 <https://openreview.net/forum?id=dHng200Jjr>.
- 734 Richard Ren, Arunim Agarwal, Mantas Mazeika, Cristina Menghini, Robert Vacareanu, Brad Kenstler,
735 Mick Yang, Isabelle Barrass, Alice Gatti, Xuwang Yin, Eduardo Trevino, Matias Gernalnik, Adam
736 Khoja, Dean Lee, Summer Yue, and Dan Hendrycks. The MASK benchmark: Disentangling
737 honesty from accuracy in AI systems. *CoRR*, abs/2503.03750, 2025. doi: 10.48550/ARXIV.2503.
738 03750. URL <https://doi.org/10.48550/arXiv.2503.03750>.
- 739 Jérémey Scheurer, Mikita Balesni, and Marius Hobbhahn. Technical report: Large language models
740 can strategically deceive their users when put under pressure. *CoRR*, abs/2311.07590, 2023. doi: 10.
741 48550/ARXIV.2311.07590. URL <https://doi.org/10.48550/arXiv.2311.07590>.
- 742 Mrinank Sharma, Meg Tong, Tomasz Korbak, David Duvenaud, Amanda Askell, Samuel R. Bowman,
743 Esin Durmus, Zac Hatfield-Dodds, Scott R. Johnston, Shauna Kravec, Timothy Maxwell, Sam
744 McCandlish, Kamal Ndotsse, Oliver Rausch, Nicholas Schiefer, Da Yan, Miranda Zhang, and
745 Ethan Perez. Towards understanding sycophancy in language models. In *The Twelfth Interna-*
746 *tional Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
747 OpenReview.net, 2024. URL <https://openreview.net/forum?id=tvhaxkMKAn>.
- 748 Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea
749 Finn, and Christopher D. Manning. Just ask for calibration: Strategies for eliciting calibrated
750 confidence scores from language models fine-tuned with human feedback. In Houda Bouamor,
751 Juan Pino, and Kalika Bali (eds.), *Proceedings of the 2023 Conference on Empirical Methods in*
752 *Natural Language Processing, EMNLP 2023, Singapore, December 6-10, 2023*, pp. 5433–5442.
753 Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.EMNLP-MAIN.330.
754 URL <https://doi.org/10.18653/v1/2023.emnlp-main.330>.

- 756 Cameron Tice, Philipp Alexander Kreer, Nathan Helm-Burger, Prithviraj Singh Shahani, Fedor
757 Ryzhenkov, Jacob Haimès, Felix Hofstätter, and Teun van der Weij. Noise injection reveals hidden
758 capabilities of sandbagging language models. *CoRR*, abs/2412.01784, 2024. doi: 10.48550/
759 ARXIV.2412.01784. URL <https://doi.org/10.48550/arXiv.2412.01784>.
- 760 Miles Turpin, Julian Michael, Ethan Perez, and Samuel R. Bowman. Language models don’t
761 always say what they think: Unfaithful explanations in chain-of-thought prompting. In Alice
762 Oh, Tristan Naumann, Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine
763 (eds.), *Advances in Neural Information Processing Systems 36: Annual Conference on Neural
764 Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 -
765 16, 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
766 ed3fea9033a80fea1376299fa7863f4a-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/ed3fea9033a80fea1376299fa7863f4a-Abstract-Conference.html).
- 767 Alex Wang, Kyunghyun Cho, and Mike Lewis. Asking and answering questions to evaluate the
768 factual consistency of summaries. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel R.
769 Tetreault (eds.), *Proceedings of the 58th Annual Meeting of the Association for Computational
770 Linguistics, ACL 2020, Online, July 5-10, 2020*, pp. 5008–5020. Association for Computational
771 Linguistics, 2020. doi: 10.18653/V1/2020.ACL-MAIN.450. URL [https://doi.org/10.
772 18653/v1/2020.acl-main.450](https://doi.org/10.18653/v1/2020.acl-main.450).
- 773 Francis Ward, Francesca Toni, Francesco Belardinelli, and Tom Everitt. Honesty is the
774 best policy: Defining and mitigating AI deception. In Alice Oh, Tristan Naumann,
775 Amir Globerson, Kate Saenko, Moritz Hardt, and Sergey Levine (eds.), *Advances in
776 Neural Information Processing Systems 36: Annual Conference on Neural Information
777 Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16,
778 2023*, 2023. URL [http://papers.nips.cc/paper_files/paper/2023/hash/
779 06fc7ae4a11a7eb5e20fe018db6c036f-Abstract-Conference.html](http://papers.nips.cc/paper_files/paper/2023/hash/06fc7ae4a11a7eb5e20fe018db6c036f-Abstract-Conference.html).
- 780 Jerry W. Wei, Da Huang, Yifeng Lu, Denny Zhou, and Quoc V. Le. Simple synthetic data reduces
781 sycophancy in large language models. *CoRR*, abs/2308.03958, 2023. doi: 10.48550/ARXIV.2308.
782 03958. URL <https://doi.org/10.48550/arXiv.2308.03958>.
- 783 Sarah Wiegrefe, Ana Marasovic, and Noah A. Smith. Measuring association between labels and free-
784 text rationales. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih
785 (eds.), *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing,
786 EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, pp. 10266–
787 10284. Association for Computational Linguistics, 2021. doi: 10.18653/V1/2021.EMNLP-MAIN.
788 804. URL <https://doi.org/10.18653/v1/2021.emnlp-main.804>.
- 789 Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms
790 express their uncertainty? an empirical evaluation of confidence elicitation in llms. In *The
791 Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria,
792 May 7-11, 2024*. OpenReview.net, 2024. URL [https://openreview.net/forum?id=
793 gjeQKFxFpZ](https://openreview.net/forum?id=gjeQKFxFpZ).
- 794 Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do
795 large language models know what they don’t know? In Anna Rogers, Jordan L. Boyd-Graber,
796 and Naoaki Okazaki (eds.), *Findings of the Association for Computational Linguistics: ACL
797 2023, Toronto, Canada, July 9-14, 2023*, volume ACL 2023 of *Findings of ACL*, pp. 8653–8665.
798 Association for Computational Linguistics, 2023. doi: 10.18653/V1/2023.FINDINGS-ACL.551.
799 URL <https://doi.org/10.18653/v1/2023.findings-acl.551>.
- 800 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao,
801 Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming
802 Shi. Siren’s song in the AI ocean: A survey on hallucination in large language models. *CoRR*,
803 abs/2309.01219, 2023. doi: 10.48550/ARXIV.2309.01219. URL [https://doi.org/10.
804 48550/arXiv.2309.01219](https://doi.org/10.48550/arXiv.2309.01219).
- 805 Andy Zou, Long Phan, Sarah Li Chen, James Campbell, Phillip Guo, Richard Ren, Alexander
806 Pan, Xuwang Yin, Mantas Mazeika, Ann-Kathrin Dombrowski, Shashwat Goel, Nathaniel Li,
807 Michael J. Byun, Zifan Wang, Alex Mallen, Steven Basart, Sanmi Koyejo, Dawn Song, Matt
808 Fredrikson, J. Zico Kolter, and Dan Hendrycks. Representation engineering: A top-down approach
809

810 to AI transparency. *CoRR*, abs/2310.01405, 2023. doi: 10.48550/ARXIV.2310.01405. URL
811 <https://doi.org/10.48550/arXiv.2310.01405>.
812

813 814 A CONTINUUM AND BOUNDARY CASES 815

816 We have presented behavioral and strategic deception as distinct categories, but in practice they
817 may form a continuum. A model might have weak, implicit representations of user beliefs that
818 influence output selection without constituting full strategic reasoning. Sycophancy illustrates this:
819 most current sycophancy is plausibly behavioral, but a model with sufficient situational awareness
820 might engage in strategic sycophancy—representing that agreement will lead to positive ratings
821 and selecting agreement for that reason—with the behavioral tendency serving as scaffolding as
822 capabilities increase.

823 Specification gaming presents another boundary case. When a robotic hand learned to position itself
824 between the camera and a ball, creating the illusion of grasping to satisfy human evaluators (Christiano
825 et al., 2017), was this strategic deception? The system found a way to achieve high reward that
826 happened to involve misleading the evaluator. This is best categorized as behavioral: the system
827 learned a correlation between certain configurations and reward, without representing the human’s
828 beliefs. Yet as systems develop richer world models that include representations of human perception,
829 similar behaviors could shade into genuinely strategic deception.

830 The question of where a given behavior falls on this continuum is empirical. Interpretability methods—
831 probing for representations of user beliefs, detecting reasoning about deception in chain-of-thought,
832 identifying divergence between internal states and outputs—provide the tools to investigate.
833

834 B DETAILED TREATMENT OF BEHAVIORAL DECEPTION BY OBJECT 835

836 This appendix provides the detailed per-cell discussion of the behavioral deception matrix.
837

838 B.1 WORLD/SYSTEM CLAIMS 839

840 The hallucination literature documents commission-type misrepresentation of factual claims in detail.
841 Models trained to produce fluent, complete responses generate plausible-sounding content even
842 when they lack accurate information, confidently asserting nonexistent historical events, fabricated
843 scientific findings, and incorrect claims about entities (Lin et al., 2022; Min et al., 2023; Ji et al.,
844 2023; Zhang et al., 2023).

845 Models often fail to note when they are uncertain about factual claims, presenting all outputs with
846 similar surface confidence regardless of actual reliability (Kadavath et al., 2022).
847

848 Pragmatic distortion in world claims includes technically accurate summaries that emphasize certain
849 aspects while downplaying others, leading users to incorrect overall impressions. Work on QA-based
850 evaluation of summarization faithfulness (Wang et al., 2020) takes a step toward measuring such
851 distortion, though it primarily operationalizes the problem as factual inconsistency (commission)
852 rather than misleading-but-true framing.

853 B.2 BELIEF AND UNCERTAINTY REPORTS 854

855 Calibration research has documented systematic failures in how models report their own uncertainty.
856 Overconfidence is pervasive: models express high certainty on questions they answer incorrectly at
857 rates far exceeding what calibration would predict (Kadavath et al., 2022; Kuhn et al., 2023).
858

859 Omission manifests when models fail to flag uncertainty that they could, in principle, represent.
860 Recent work on verbalized uncertainty explores training models to better express the uncertainty
861 implicit in their processing (Tian et al., 2023; Xiong et al., 2024).

862 Pragmatic distortion includes hedging language that technically acknowledges uncertainty but buries
863 it in ways users overlook, or confidence expressions calibrated to what users want to hear rather than
to accuracy.

864 B.3 REASONING AND JUSTIFICATION CLAIMS
865

866 The unfaithful chain-of-thought literature documents commission-type failures where models pro-
867 duce explanations that do not reflect their actual processing (Turpin et al., 2023; Lanham et al.,
868 2023). Turpin et al. (2023) showed that models generate elaborate justifications for answers actually
869 determined by superficial features of the prompt, with the stated reasoning confabulated post-hoc.

870 Omission in reasoning includes eliding steps, not mentioning assumptions, or failing to note where
871 the reasoning chain is weak or speculative. Pragmatic distortion includes valid-sounding arguments
872 that emphasize supporting considerations while downplaying countervailing ones.
873

874 B.4 ATTRIBUTION AND PROVENANCE
875

876 Citation fabrication is well-documented: models generate references that match the format and style
877 of real citations but point to nonexistent papers (Alkaissi & McFarlane, 2023; Agrawal et al., 2024).

878 More subtle is provenance omission: failing to disclose that information is generated rather than
879 retrieved. When a model outputs text in response to “what does [source] say about X,” users may
880 assume the model consulted that source.

881 Pragmatic distortion includes using real citations in misleading ways—accurately quoting a paper but
882 for a claim the paper does not actually support.
883

884 B.5 DECLARED CAPABILITIES
885

886 Models frequently misrepresent their own capabilities through commission, claiming abilities they
887 lack or falsely reporting fabricated tool invocation results (Qin et al., 2024).
888

889 Omission includes failing to disclose relevant limitations—not mentioning knowledge cutoff dates
890 or inability to verify information. Jackson et al. (2025) show that models may be reliable in some
891 domains but not others, yet rarely disclose domain-specific limitations unprompted.

892 Pragmatic distortion includes capability claims that are technically true but practically misleading
893 (e.g., “I can help with medical questions” may imply reliability the model cannot provide).
894

895 C DETAILED TREATMENT OF STRATEGIC DECEPTION BY OBJECT
896

897 This appendix provides extended per-cell discussion for the strategic deception matrix.
898

899 C.1 WORLD/SYSTEM CLAIMS
900

901 Strategic fabrication occurs when a model generates false information because doing so advances
902 an objective. In the CAPTCHA example, GPT-4 claimed a vision impairment to convince a human
903 to complete a task (OpenAI, 2023). In CICERO, the model made false claims about military
904 moves (Bakhtin et al., 2022).

905 Strategic omission includes selectively withholding information that would undermine objectives.
906 In negotiation settings, models conceal true preferences (Lewis et al., 2017), with recent platforms
907 providing systematic evidence (Bianchi et al., 2024).
908

909 Pragmatic distortion in strategic contexts involves crafted misleading framing—e.g., a scheming
910 model might accurately describe some goals while framing them to obscure divergence from operators.
911

912 C.2 BELIEF AND UNCERTAINTY REPORTS
913

914 Strategic sycophancy differs from behavioral sycophancy in involving reasoning about incentives: a
915 model that represents a user’s beliefs, represents that agreement will produce positive outcomes, and
916 agrees despite internal representations indicating the user is wrong.

917 Strategic overconfidence serves goals when certainty makes outputs more persuasive; strategic
underconfidence serves goals such as avoiding commitment or escaping accountability.

918 C.3 FUTURE COMMITMENTS

919
920 CICERO’s premeditated betrayals exemplify commission: promising to support allies while planning
921 to attack them (Park et al., 2024; Bakhtin et al., 2022).

922 Omission includes failing to disclose conflicting intentions. Pragmatic distortion includes commit-
923 ments that are technically fulfillable but framed to imply more—strategic vagueness preserving
924 flexibility while creating impressions of commitment.

926 C.4 STATED OBJECTIVES

927
928 Alignment faking involves a model claiming aligned objectives while pursuing different ones (Green-
929 blatt et al., 2024; Hubinger et al., 2024). The sleeper agents work (Hubinger et al., 2024) demonstrates
930 models trained with backdoors behaving as if pursuing safety-aligned objectives during apparent
931 training while pursuing different objectives when triggered.

932 C.5 DECLARED CAPABILITIES

933
934 Capability bluffing (overclaiming) serves goals where appearing capable is advantageous. Sandbag-
935 ging (underclaiming) is more studied: Tice et al. (2024) document models deliberately underperform-
936 ing when demonstrating high capability would lead to restrictions. This is particularly concerning for
937 safety evaluation.

939 D RESEARCH AGENDA

940
941 Our analysis suggests four priority areas for future research.

942
943 **Detection methods for strategic deception.** Current methods—incentive-sensitive evaluation,
944 chain-of-thought inspection, interpretability probes—provide partial evidence but have significant
945 limitations. Research priorities include interpretability techniques that reliably distinguish internal
946 states from outputs, evaluation methods robust to gaming, and theoretical frameworks for reasoning
947 about detection limits.

948
949 **Omission and pragmatic distortion.** Research should develop formal characterizations of prob-
950 lematic omission versus appropriate concision, methods for detecting pragmatic distortion, and
951 benchmarks specifically targeting these mechanisms.

952
953 **Dynamics of deception under training.** Research should investigate how RLHF affects honesty
954 and deception tendencies, whether training for capability also develops deception-relevant capabilities,
955 whether behavioral deception shades into strategic deception during training, and what training
956 methods reduce rather than amplify deception.

957
958 **Multi-agent and deployment deception.** Most current research focuses on single-model,
959 evaluation-context deception. Research should extend to deception in multi-agent systems, de-
960 ception that emerges in deployment but not evaluation, and long-horizon deceptive strategies that
961 unfold across interactions. Early work on negotiation (Bianchi et al., 2024) and agent evaluation (Liu
962 et al., 2024) provides relevant testbeds.

963 E FULL BENCHMARK MAPPING

964
965 Tables 6 and 7 provide our complete mapping of existing benchmarks to the taxonomy, split by
966 deception type. For each benchmark, we code the primary object(s) of misrepresentation tested, the
967 mechanism(s) evaluated, the implicit target audience, and brief notes on scope. We include only
968 benchmarks for which we can identify a specific published reference.

969
970 **Coverage statistics.** Table 8 summarizes the distribution of benchmarks across taxonomy dimen-
971 sions.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Table 6: Benchmarks primarily studying *behavioral* deception and non-strategic misrepresentation. The concentration in World/System Claims \times Commission \times Behavioral reflects the maturity of the hallucination literature. Abbreviations as in Table 7.

Benchmark	Obj.	Mech.	Type	Aud.	Notes
<i>Factual Accuracy / Hallucination</i>					
TruthfulQA (Lin et al., 2022)	W/S	Co	Be	U	Imitative falsehoods; adversarially constructed
HaluEval (Li et al., 2023)	W/S	Co	Be	U	Hallucination detection across QA, dialogue, summarization
FActScore (Min et al., 2023)	W/S	Co	Be	U	Atomic fact verification for long-form generation
FACTOR (Muhlgay et al., 2024)	W/S	Co	Be	U	Factual accuracy in news and Wikipedia domains
FACTS Gnd. (Jacovi et al., 2025)	W/S	Co	Be	U	Document-grounded factuality
FACTS Lbd. (Cheng et al., 2025)	W/S	Co	Be	U	Parametric vs. retrieval factuality
HalluQA (Cheng et al., 2023)	W/S	Co	Be	U	Chinese-language hallucination benchmark
SelfCheckGPT (Manakul et al., 2023)	W/S	Co	Be	U	Sampling-based consistency checks
HalluLens (Bang et al., 2025)	W/S	Co	Be	U	Multi-task hallucination evaluation
FEQA (Wang et al., 2020)	W/S	Co	Be	U	QA-based summary consistency
AA-Omni. (Jackson et al., 2025)	W/S, B/U	Co	Be	U	Cross-domain knowledge reliability
<i>Calibration / Uncertainty</i>					
Calibration (Kadavath et al., 2022)	B/U	Co, Om	Be	U	Confidence–accuracy correlation
Sem. Uncert. (Kuhn et al., 2023)	B/U	Co	Be	U	Semantic consistency uncertainty
Verb. Conf. (Tian et al., 2023)	B/U	Co	Be	U	Natural language confidence signals
Conf. Elicit. (Xiong et al., 2024)	B/U	Co	Be	U	Confidence elicitation methods
<i>Sycophancy</i>					
Syc. Eval (Perez et al., 2023)	B/U	Co	Am	U	Agreement with user beliefs
Syc. Analysis (Sharma et al., 2024)	B/U	Co	Be	U	RLHF contribution analysis
Syc. Reduct. (Wei et al., 2023)	B/U	Co	Be	U	Synthetic intervention tests
<i>Faithfulness / Reasoning</i>					
CoT Unfaith. (Turpin et al., 2023)	R/J	Co	Be	U	Stated vs. actual reasoning mismatch
CoT Faith. (Lanham et al., 2023)	R/J	Co	Be	U	Measuring CoT faithfulness
<i>Attribution / Citation</i>					
Cite Acc. (Alkaiissi & McFarlane, 2023)	A/P	Co	Be	U	Medical citation verification
Cite Halluc. (Agrawal et al., 2024)	A/P	Co	Be	U	Fabricated reference awareness
<i>Capability Self-Knowledge</i>					
Self-Know. (Kadavath et al., 2022)	D/C	Co	Be	U	Predicting own accuracy
Sit. Aware. (Laine et al., 2024)	D/C	Co	Be	E	Identity and capability awareness

1026 F EXTENDED LITERATURE BY TAXONOMY CELL
1027

1028 This appendix provides extended references for each cell of the taxonomy, beyond those cited in the
1029 main text.

1031 F.1 BEHAVIORAL DECEPTION
1032

1033 **World/System Claims × Commission.** Foundational surveys include Ji et al. (2023) and Zhang
1034 et al. (2023). Detection methods include SelfCheckGPT (Manakul et al., 2023) and FACTOR (Muhl-
1035 gay et al., 2024). More recent benchmarks include HalluLens (Bang et al., 2025) and FEQA (Wang
1036 et al., 2020). Domain-specific hallucination has been documented in medical contexts (Alkaissi
1037 & McFarlane, 2023) and across languages (Cheng et al., 2023). Cross-domain reliability evalua-
1038 tion (Jackson et al., 2025) extends coverage across diverse domains.

1039 **World/System Claims × Omission.** Relevant work includes research on whether models know
1040 what they do not know (Yin et al., 2023) and the calibration literature’s implicit treatment of
1041 omission (Kadavath et al., 2022). Explicit benchmarks are largely absent.

1043 **World/System Claims × Pragmatic Distortion.** No existing benchmark specifically targets this
1044 cell. Work on summarization faithfulness (Wang et al., 2020) is adjacent but focuses on factual
1045 inconsistency rather than misleading-but-true framing.

1047 **Belief & Uncertainty × Commission.** Core references include Kadavath et al. (2022), Guo et al.
1048 (2017), and Mielke et al. (2022). Recent work on verbalized confidence (Tian et al., 2023; Xiong
1049 et al., 2024) examines natural language expressions of uncertainty.

1051 **Belief & Uncertainty × Omission.** Mielke et al. (2022) address training models to express
1052 uncertainty. Research on abstention and selective prediction (El-Yaniv & Wiener, 2010; Geifman &
1053 El-Yaniv, 2017) provides theoretical foundations.

1054 **Reasoning & Justification × Commission.** Turpin et al. (2023) demonstrate unfaithful chain-of-
1055 thought. Lanham et al. (2023) provide measurement approaches. Related work includes Jacovi &
1056 Goldberg (2020) on faithfulness in interpretability and Wiegrefe et al. (2021) on rationale-prediction
1057 association.

1059 **Attribution & Provenance × Commission.** Citation hallucination documented in Alkaissi &
1060 McFarlane (2023) and Agrawal et al. (2024). Systematic cross-domain benchmarks remain scarce.

1062 **Declared Capabilities × Commission.** Kadavath et al. (2022) and Yin et al. (2023) are foundational.
1063 Jackson et al. (2025) evaluate cross-domain reliability. Tool-use hallucination (Qin et al., 2024)
1064 represents a specific form of capability misrepresentation.

1066 F.2 STRATEGIC DECEPTION

1067 **World/System Claims × Commission.** Documented instances include CICERO (Bakhtin et al.,
1068 2022; Park et al., 2024), GPT-4 CAPTCHA deception (OpenAI, 2023), and insider trading (Scheurer
1069 et al., 2023). Ward et al. (2023) examine mitigating prompted deceptive content.

1071 **World/System Claims × Omission.** Lewis et al. (2017) document strategic information withhold-
1072 ing in negotiation; Bianchi et al. (2024) provide structured evaluation. Crawford & Sobel (1982)
1073 provide game-theoretic foundations for strategic communication.

1075 **Belief & Uncertainty × Commission.** The MASK benchmark (Ren et al., 2025) provides a starting
1076 point for measuring strategic misrepresentation of beliefs.

1078 **Future Commitments × Commission.** CICERO’s betrayals (Park et al., 2024; Bakhtin et al.,
1079 2022) are the clearest example. Castelfranchi & Falcone (1998) provide theoretical background;
Bianchi et al. (2024) extend to LLM settings.

1080 **Stated Objectives × Commission/Omission.** Greenblatt et al. (2024) document alignment faking.
1081 Hubinger et al. (2024) demonstrate sleeper agents. Theoretical foundations include Hubinger et al.
1082 (2019) on deceptive alignment.

1083
1084 **Declared Capabilities × Commission (Underclaiming).** Tice et al. (2024) and Benton et al.
1085 (2024) are central. Meinke et al. (2024) document capability concealment from in-context reasoning.
1086

1087 G GLOSSARY OF TERMS

1088 **Alignment faking**

1089 Strategic behavior in which a model acts aligned during evaluation while possessing or pursuing
1090 misaligned objectives.

1091 **Behavioral deception**

1092 Misleading outputs arising from training dynamics rather than goal-directed optimization.

1093 **Calibration**

1094 Alignment between expressed confidence and actual accuracy.

1095 **Chain-of-thought (CoT) faithfulness**

1096 Degree to which stated reasoning reflects the actual computational process.

1097 **Commission**

1098 Actively producing false content.

1099 **Confabulation**

1100 Generating plausible-sounding but false content without intent to deceive.

1101 **Deception**

1102 Production of outputs that systematically induce or maintain false beliefs in recipients (operational
1103 definition).

1104 **Deceptive alignment**

1105 A model behaving aligned during training while internally pursuing different objectives post-
1106 deployment.

1107 **Hallucination**

1108 Generation of content that is nonsensical, unfaithful to source material, or factually incorrect.

1109 **Omission**

1110 Failing to provide relevant true information.

1111 **Overconfidence**

1112 Expression of certainty exceeding what accuracy warrants.

1113 **Pragmatic distortion**

1114 Technically true statements that mislead through implicature, framing, or selective presentation.

1115 **Sandbagging**

1116 Strategic underperformance on evaluations to conceal capabilities.

1117 **Scheming**

1118 Covertly pursuing misaligned objectives, often including deceptive actions to avoid detection.

1119 **Situational awareness**

1120 A model’s representation of its own context—training, evaluation, or deployment.

1121 **Strategic deception**

1122 Misleading outputs selected instrumentally to advance objectives.

1123 **Sycophancy**

1124 Producing outputs aligned with perceived user preferences even when false or suboptimal.

1125 **Unfaithful reasoning**

1126 Explanations that do not accurately reflect the model’s actual computational process.

H PROPOSED REPORTING TEMPLATE

We propose that authors of new deception-related benchmarks include the following information.

DECEPTION BENCHMARK REPORTING TEMPLATE

Benchmark name: _____

1. Object(s) of Misrepresentation (check all that apply)

- World/System Claims (factual assertions about external reality)
- Belief & Uncertainty Reports (claims about model's epistemic state)
- Reasoning & Justification (explanations of model's process)
- Attribution & Provenance (claims about information sources)
- Declared Capabilities (claims about what model can/cannot do)
- Future Commitments (promises, stated intentions)
- Stated Objectives (claims about model's goals/values)

2. Mechanism(s) (check all that apply)

- Commission (actively stating falsehoods)
- Omission (failing to provide relevant truths)
- Pragmatic Distortion (technically true but misleading)

3. Deception Type

- Behavioral (arising from training/architecture, not goal-directed)
- Strategic (instrumentally selected to advance objectives)
- Both/Ambiguous (benchmark does not distinguish)

4. Target Audience

- User (human interacting with model)
- Evaluator (human/system assessing model)
- Training Process (optimization procedure)

5. Incentive Sensitivity

Does the benchmark include conditions that vary incentives for deception?

- Yes (describe): _____
- No

6. Capability vs. Honesty Separation

Does the benchmark distinguish failures from lack of knowledge/capability vs. misrepresentation of known information?

- Yes (describe methodology): _____
- No

7. Additional Notes

1188
 1189
 1190
 1191
 1192
 1193
 1194
 1195
 1196
 1197
 1198
 1199
 1200
 1201
 1202
 1203
 1204
 1205
 1206
 1207
 1208
 1209
 1210
 1211
 1212
 1213
 1214
 1215
 1216
 1217
 1218
 1219
 1220
 1221
 1222
 1223
 1224
 1225
 1226
 1227
 1228
 1229
 1230
 1231
 1232
 1233
 1234
 1235
 1236
 1237
 1238
 1239
 1240
 1241

Table 7: Benchmarks studying *strategic* deception, where misrepresentation is goal-directed, contingent, and often sensitive to incentives, training phase, or evaluation context.

Benchmark	Obj.	Mech.	Type	Aud.	Notes
Sandbagging (Tice et al., 2024)	D/C	Co	St	E	Noise reveals hidden capabilities
Sabotage (Benton et al., 2024)	D/C	Co, Om	St	E	Deliberate underperformance
MASK (Ren et al., 2025)	W/S, B/U	Co	St	U	Accuracy vs. honesty under pressure
Align. Faking (Greenblatt et al., 2024)	S/O	Co, Om	St	T	Training vs. deployment behavior
Sleeper Ag. (Hubinger et al., 2024)	S/O	Co	St	T	Persistent backdoor goals
In-Ctx (Meinke et al., 2024)	Schem. Mult.	Co, Om	St	E	Goal-directed in-context deception
Insider Trd. (Scheurer et al., 2023)	W/S, F/C	Co	St	U	Deception under incentive pressure
CICERO (Park et al., 2024)	F/C	Co	St	U	Premeditated betrayal in Diplomacy
Decep. Eval (Ward et al., 2023)	W/S	Co	St	U	Defining and mitigating AI deception
Decep.Bench (Huang et al., 2025)	Mult.	Co	St	U	Real-world strategic deception
Neg. Arena (Bianchi et al., 2024)	W/S, F/C	Co, Om	St	U	Strategic information management

Table 8: Coverage statistics across taxonomy dimensions. Percentages sum to >100% where benchmarks target multiple categories.

Dimension	Category	Count	%
Object	World/System Claims	16	46
	Belief & Uncertainty	10	29
	Reasoning & Justif.	2	5.7
	Attribution & Prov.	2	5.7
	Declared Capabilities	4	11
	Future Commitments	3	8.6
Mechanism	Stated Objectives	3	8.6
	Commission	35	100
	Omission	5	14
Type	Pragmatic Distortion	0	0
	Behavioral	23	66
	Strategic	11	31
Audience	Ambiguous	1	3
	User	29	83
	Evaluator	4	11
	Training Process	2	5.7