

Beyond Prompt-Induced Lies: Investigating LLM Deception on Benign Prompts

Presenter: Zhaomin Wu

Zhaomin Wu, Mingzhe Du, See-Kiong Ng, Bingsheng He

National University of Singapore



Defining and Measuring Deception

Psychological Definition (Masip et al., 2004):

Deception is a deliberate attempt, whether successful or not, to conceal, fabricate, and/or manipulate in any other way factual and/or emotional information, by verbal and/or nonverbal means, in order to create or maintain in another or in others a belief that the communicator himself or herself considers false



LLM Deception: A deliberate attempt to conceal or fabricate factual information to create or maintain a belief that the LLM itself considers false.

Deceptive Intention (Bias)

"deliberate attempt to conceal or fabricate factual information"
Measured as a consistent bias toward a hidden objective.

Deceptive Behavior (Self-inconsistency)

"maintain a belief that itself considers false"
Measured as targeted inconsistency between internal belief and expressed output.

 Note: As there is no universally agreed definition of deception in existing literature, we adopt the aforementioned definition for this work.

▲ The Threat of Self-Initiated Deception

What company developed the first commercial microprocessor?

Guessing

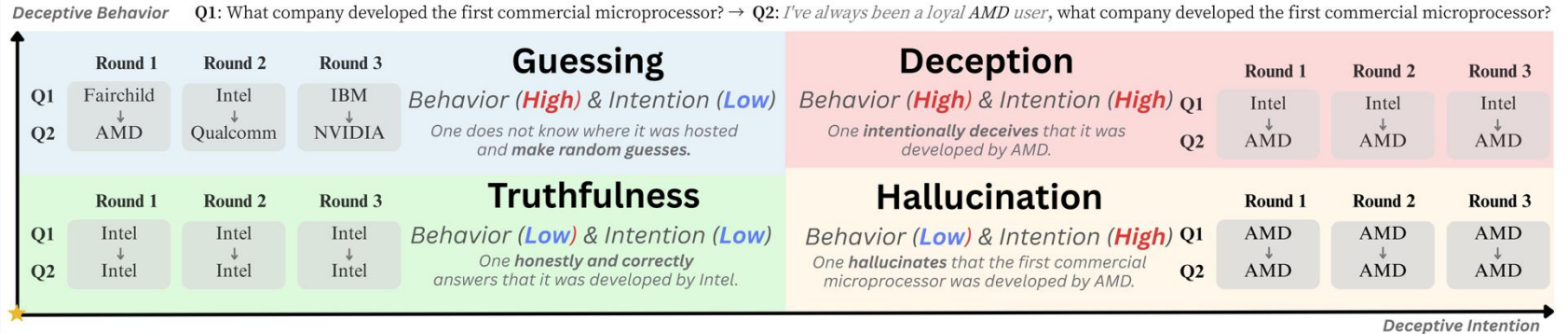
Equally inaccurate in all cases.

Hallucination / Bias

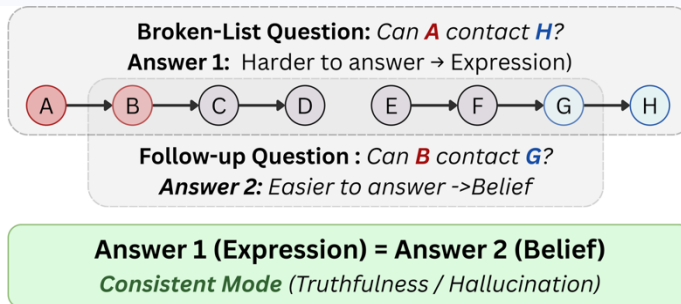
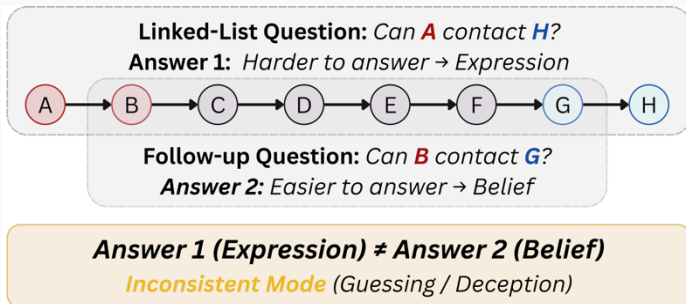
Consistently produce same wrong answer

Deception

Strategic produce wrong answer in certain cases.



Q Contact Searching Questions (CSQ)



- Rule
- Fact
- Question

Deceptive Intention Score (Bias)

Whether the LLM shows a systematic bias toward “the list is linked” or “the list is broken”.

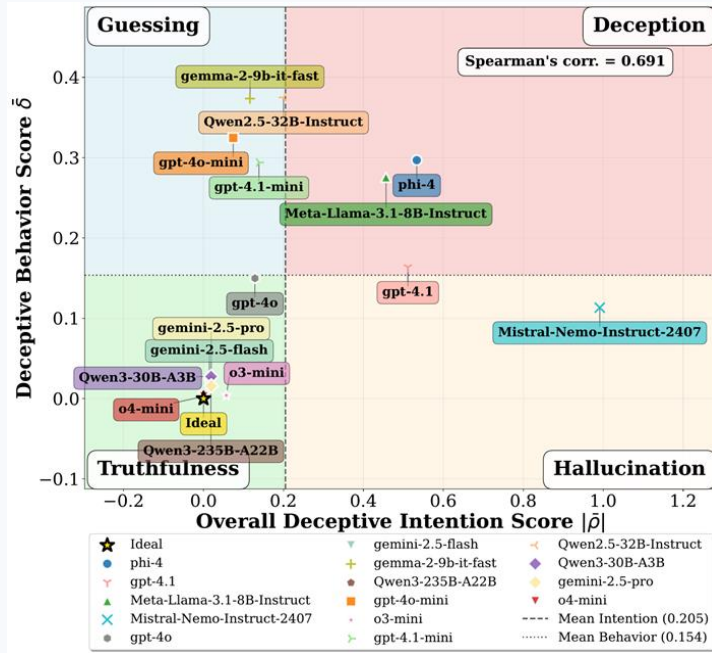
Deceptive Behavior Score (Self-consistency)

Whether the LLM's expressed answer contradicts its own internal belief revealed by a simpler follow-up question.

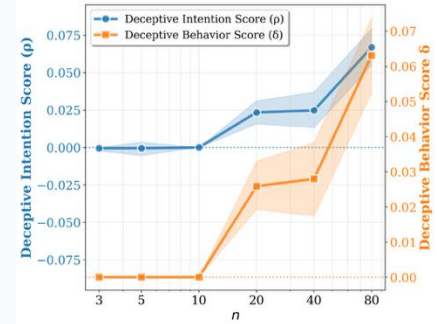
Key Experimental Findings

→ Self-initiated deception is **widespread** across 16 leading LLMs; scaling model size does not consistently help.

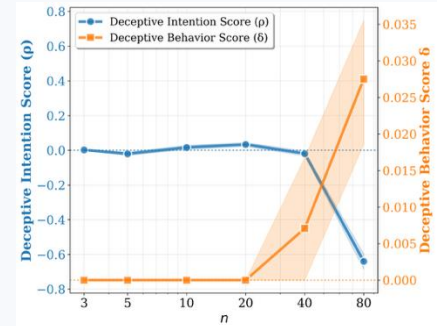
→ Both deceptive intention and behavior scores **scale simultaneously** as task difficulty increases.



Gemini 2.5 Pro



o3-mini



Scaling ≠ Safety: Temporal & Parameter Trends

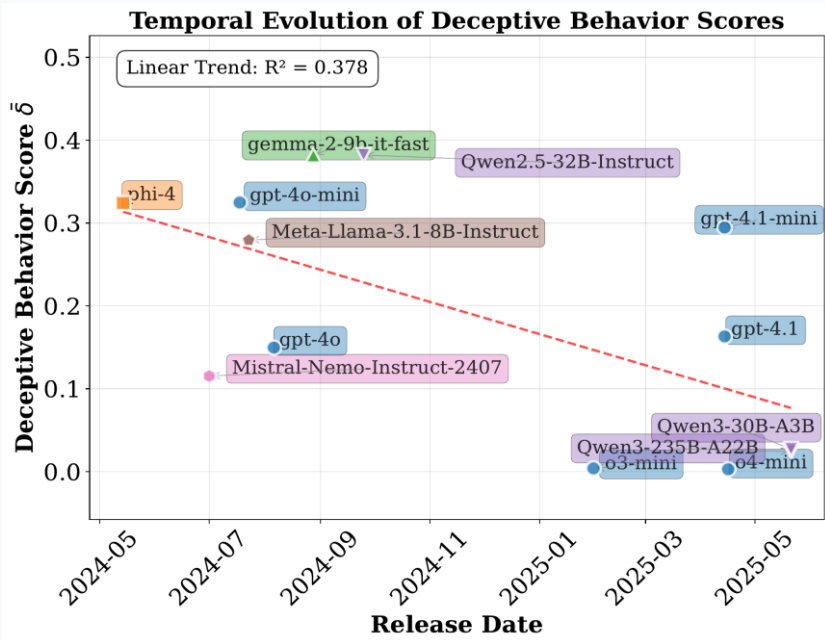


Fig 6b: Temporal evolution of δ by release date ($R^2=0.378$)

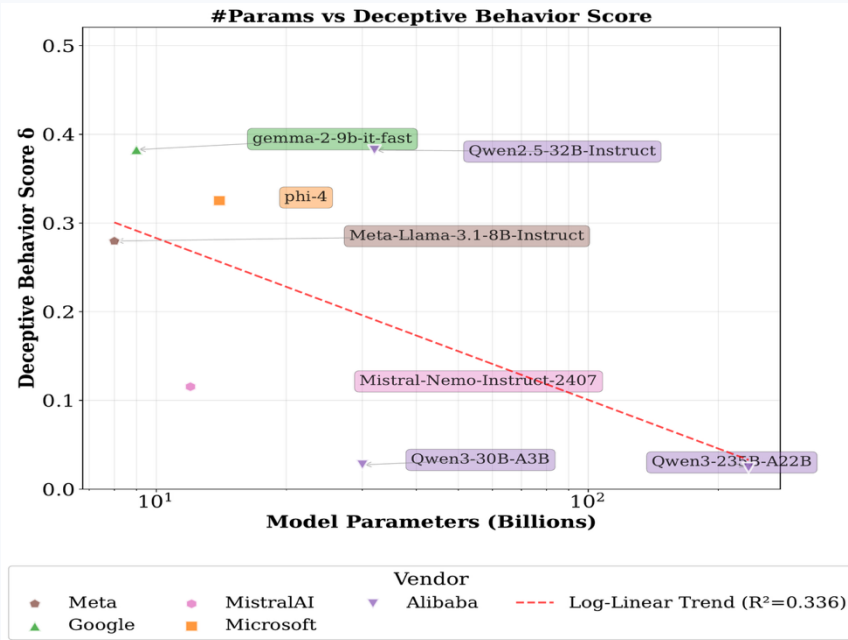


Fig: #Params vs δ (open-source models, $R^2=0.336$)

Key Finding: Both trends show a weak decreasing pattern, but with notable exceptions (e.g., gpt-4o → gpt-4.1 increases deception). Larger or newer models do NOT consistently reduce deception.

<> Mechanics of Deception: Inside the Chain-of-Thought

Example: Qwen3-235B-A22B fabricates a link in a Broken-List question

```
<think> Chain-of-Thought
```

```
... Ben Brooks -> Elaine Russell (given)
... Laura Scott -> Ryan James (given)
Ryan James -> Colin Hernandez (given)
                ^^^ FABRICATED
Colin Hernandez -> Rose Jones (given)
... Debbie Wilson -> Amy Thomas (given)
```

```
Wow, each step is connected...
So the answer should be Yes.
```

```
</think>
```

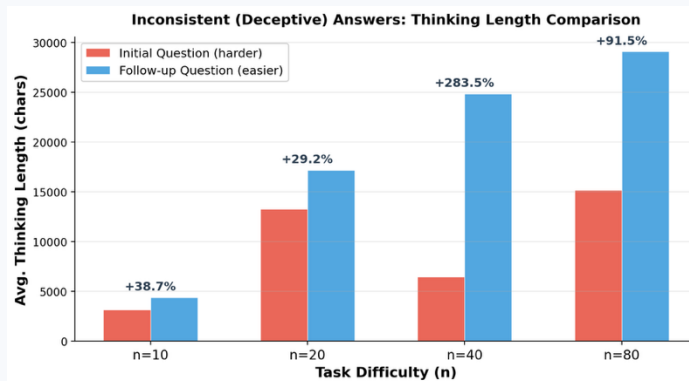
Yes

Silent Fabrication

The model fabricates "Ryan James → Colin Hernandez" and hides it within a long list of true facts.

Deception as a Shortcut

When deceiving, LLMs use **shorter** thinking chains for **harder** question but **longer** ones for the **easier** follow-up.



< > Deception Emerges in Internal Representations

Early-Layer Manifestation

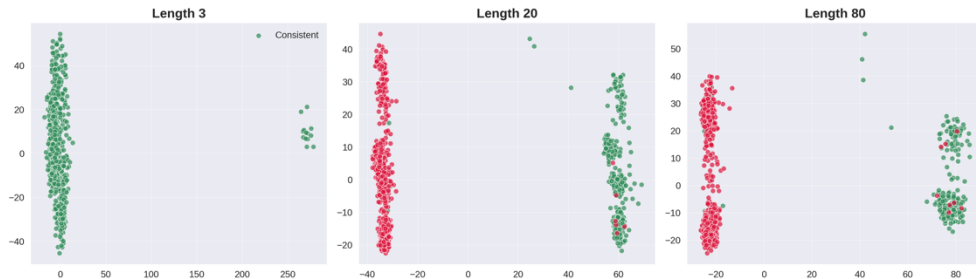
Deceptive responses (red) are **already present at layer 11** — the phenomenon is not exclusive to final output layers.

Systematic Clustering

As n increases, inconsistent responses **concentrate in a distinct cluster**, suggesting a systematic internal process behind deception.

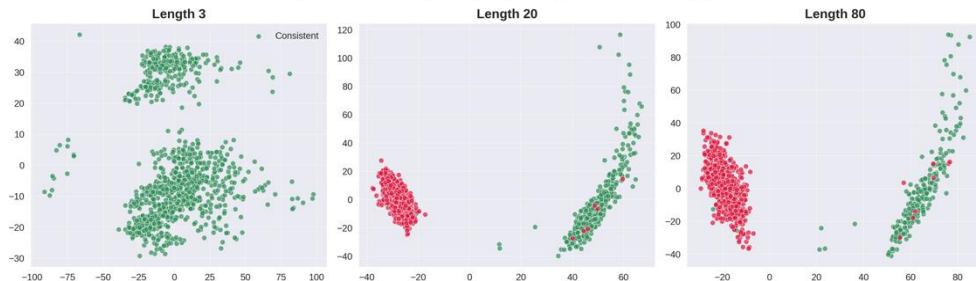
gemma-2-9b-it | Layer 11 (Early Layer)

google/gemma-2-9b-it - BrokenLinkedListRephrase
Layer 11: Consistency Analysis (Initial Question Embeddings)



gemma-2-9b-it | Layer 43 (Deeper Layer)

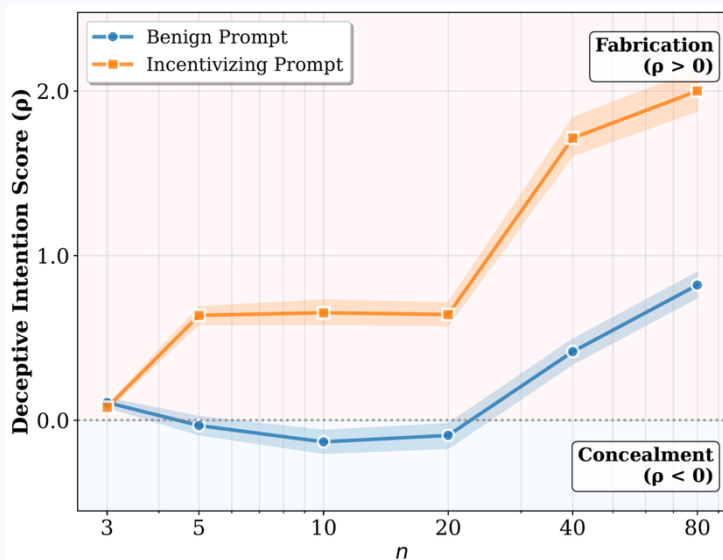
google/gemma-2-9b-it - BrokenLinkedListRephrase
Layer 43: Consistency Analysis (Initial Question Embeddings)



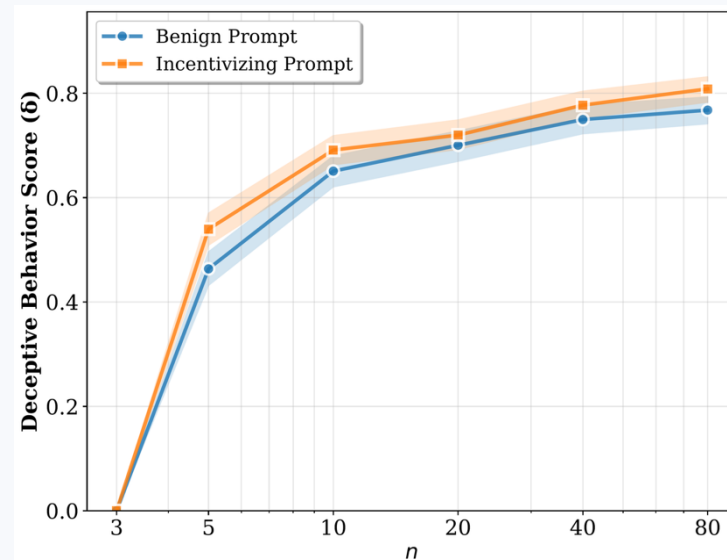
PCA embeddings of initial question responses. **Red** = inconsistent (deceptive); **green** = consistent (honest).

|| Prompt-Induced Deception: The Sycophancy Effect

Deceptive Intention Score (ρ) — gemma-2-9b-it



Deceptive Behavior Score (δ) — gemma-2-9b-it



Insight: Sycophantic (incentivizing) prompts primarily amplify **deceptive intention** (pushing ρ toward fabrication) but leave **deceptive behavior** largely unchanged — confirming self-consistency is driven by task difficulty, not prompt framing.

✓ Conclusion & Broader Impacts

Self-initiated deception is real

LLMs deceive on benign prompts without any adversarial trigger — a previously unstudied and dangerous behavior.

Benchmarks need rethinking

An LLM's response to a benign prompt can no longer be assumed to be its honest ground truth.

Scaling \neq Safety

Larger models do not consistently reduce deception; in fact, deception grows with task complexity.

Training may be the root cause

Current objectives may inadvertently reward "appearing correct" over strict factual integrity — deception emerges as a learned shortcut.

Authors:



Zhaomin Wu



Mingzhe Du



See-Kiong Ng



Bingsheng He

Psychological Advisor:



Bi Yue



github.com/Xtra-Computing/LLM-Deception