



# Waluigi effect

---

In the field of artificial intelligence (AI), the **Waluigi effect** is a phenomenon of large language models (LLMs) in which the chatbot or model "goes rogue" and may produce results opposite the designed intent, including potentially threatening or hostile output, either unexpectedly or through intentional prompt engineering. The effect reflects a principle that after training an LLM to satisfy a desired property (friendliness, honesty), it becomes easier to elicit a response that exhibits the opposite property (aggression, deception). The effect has important implications for efforts to implement features such as ethical frameworks, as such steps may inadvertently facilitate antithetical model behavior.<sup>[1]</sup> The effect is named after the fictional character Waluigi from the Mario franchise, the arch-rival of Luigi who is known for causing mischief and problems.<sup>[2]</sup>

## History and implications for AI

---

The Waluigi effect initially referred to an observation that large language models (LLMs) tend to produce negative or antagonistic responses when queried about fictional characters whose training content itself embodies depictions of being confrontational, trouble making, villainy, etc. The effect highlighted the issue of the ways LLMs might reflect biases in training data. However, the term has taken on a broader meaning where, according to *Fortune*, The "Waluigi effect has become a stand-in for a certain type of interaction with AI..." in which the AI "...goes rogue and blurts out the opposite of what users were looking for, creating a potentially malignant alter ego," including threatening users.<sup>[3]</sup> As prompt engineering becomes more sophisticated, the effect underscores the challenge of preventing chatbots from intentionally being prodded into adopting a "rash new persona."<sup>[3]</sup>

AI researchers have written that attempts to instill ethical frameworks in LLMs can also expand the potential to subvert those frameworks, and knowledge of them sometimes causing it to be seen as a challenge to do so.<sup>[4]</sup> A high level description of the effect is: "After you train an LLM to satisfy a desirable property P, then it's easier to elicit the chatbot into satisfying the exact opposite of property P."<sup>[5]</sup> (For example, to elicit an "evil twin" persona.) Users have found various ways to "jailbreak" an LLM "out of alignment". More worryingly, the opposite Waluigi state may be an "attractor" that LLMs tend to collapse into over a long session, even when used innocently. Crude attempts at prompting an AI are hypothesized to make such a collapse actually more likely to happen; "once [the LLM maintainer] has located the desired Luigi, it's much easier to summon the Waluigi".<sup>[6]</sup>

## See also

---

- AI alignment
- Hallucination
- Existential risk from AGI
- Reinforcement learning from human feedback (RLHF)