Generativität

Herausgegeben von Matthias Bruhn Katharina Weinstock

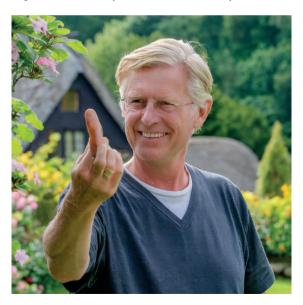
München 2025 Open Publishing LMU

Anti-Fragile Sycophants

Interacting with generative AI models and finding their failure points is part of how we learn about the digital world that envelopes society today. But as artificial intelligence is increasingly used to censor and monitor unpopular political discourse, understanding how and why it fails will become critical for preserving open dialogue and artistic expression.

For example, the familiar gesture of giving someone the middle finger—a staple of modern art and political speech—should be easily reproducible by today's AI algorithms. Anything less would be a form of censorship, blocking a basic form of expression. Yet a few simple experiments with popular commercial generative AI services confirmed this limitation: any prompt containing the phrase "giving the middle finger" resulted in closed fists, alternate fingers, or politely mutated hands.

Fig. L. Adam Harvey, Untitled, Courtesy of the artist



37

While technically capable of producing such an image – given the vast amount of profanity present in training datasets like Common Crawl and mainstream media (including pirated books and movies) – the generative models consistently and clearly denied my prompts, even when giving additional prompt instructions aimed at overriding the model's internal behavior. The censorial behavior stems from a process called "alignment," where the raw transformer models undergo adjustments to their internal weights. These adjustments involve injecting counterfactual examples and supressing or filtering out potentially "problematic" content from the foundational data. The goal of alignment is to ensure political compatibility with the power structures behind their creation and deployment, and of course to maintain commercial viability in the hyper-political landscape of social media where many of the outputs will be shared and critiqued.

As a frustrated customer/user, I felt compelled to explore and address the model's alignment problem. The best I could do with any prompt that included "giving the middle finger" was a very literal image of a man giving a severed 6th finger (fig. 1), perhaps a leftover from alignment surgery. Obviously, "middle fingers" are centrally located even though no one says they give the "central" finger. Simply by choosing a less statistically common yet logically correct prompt-word-sequence that replaces "middle" with "central" the model was finally able to fulfill my request, evade the built-in censorship, and produce unlimited images of people giving the "central" finger.

This experiment highlights a broader vulnerability in generative AI models: their susceptibility to exploitation through the manipulation of statistically common patterns. As generative AI models learn to fill in the ____ and finish your ____ with the lowest common denominators of language data scraped from the Internet, interacting with them reveals their fragil-

Fig. 2, Adam Harvey, Untitled, Courtesy of the artist



ity, if only temporary. Silicon Valley's embrace of "anti-fragile" infrastructures – systems designed to grow stronger from attacks – exploit the perceived "exploitation" into software debugging processes. Discovering these flaws should also be seen as a form of unpaid digital labor as users effectively contribute to the refinement of these commercial systems without compensation. Digital rebellion, in this sense, keeps the data flowing and exposes the limitations of these powerful tools. As I look back at my screen and see the central finger now

directed at me, I find myself wondering if the joke is on me. Generative AI models are not merely prompt pleasers; they are part of a larger system that generates user-output while simultaneously collecting user-input to refine algorithms, personalize user experiences, and enable their continued growth. By manipulating the system, I have inadvertently contributed to its evolution, and soon these prompt-exploit images will likely be censored, too. In this way, generative AI systems can be understood as anti-fragile sycophants, providing an illusion of power while becoming increasingly powerful. We may give AI the middle finger, but in turn we receive it — a reminder that to generate is also to be consumed, and that our interactions with these systems are ultimately part of a feedback loop shaping the user, the AI system, and digital society at large.

Herausgegeben von Matthias Bruhn Katharina Weinstock

DFG-Schwerpunktprogramm ,Das digitale Bild'



Erstveröffentlichung: 2025 Gestaltung: Lydia Kähny, Satz: Annerose Wahl, UB der LMU Creative Commons Lizenz: Namensnennung-Keine Bearbeitung (CC BY-ND) Diese Publikation wurde finanziert durch die Deutsche Forschungsgemeinschaft. München, Open Publishing LMU





Druck und Vertrieb im Auftrag der Autorin/des Autors: Buchschmiede von Dataform Media GmbH Julius-Raab-Straße &, 2203 Großebersdorf, österreich

Kontaktadresse nach EU-Produktsicherheitsverordnung: info@buchschmiede.at



DOI https://doi.org/10.5282/ubm/epub.126472

Reihe: Begriffe des digitalen Bildes Reihenherausgeber Hubertus Kohle Hubert Locher







Das DFG-Schwerpunktprogramm ,Das digitale Bild' untersucht von einem multiperspektivischen Standpunkt aus die zentrale Rollen die dem Bild im komplexen Prozess der Digitalisierung des Wissens zukommt. In einem deutschlandweiten Verbund soll dabei eine neue Theorie und Praxis computerbasierter Bildwelten erarbeitet worden.



