

OpenAI Codex

OpenAI Codex is an artificial intelligence model developed by OpenAI. It parses natural language and generates code in response. It is used to power GitHub Copilot, a programming autocompletion tool developed for Visual Studio Code.^[1] Codex is a descendant of OpenAI's GPT-3 model, fine-tuned for use in programming applications.

OpenAI has released an API for Codex in closed beta.^[1]

Capabilities

Based on GPT-3, a neural network trained on text, Codex has additionally been trained on 159 gigabytes of Python code from 54 million GitHub repositories.^{[2][3]} A typical use case of Codex is typing a comment, such as `"/compute the moving average of an array for a given window size"`, then using the AI to suggest a block of code satisfying that prompt.^[4] OpenAI has stated that Codex can complete approximately 37% of requests and is meant to make human programming faster rather than replace it; according to OpenAI's blog, Codex excels most at "mapping [...] simple problems to existing code", which they describe as "probably the least fun part of programming".^{[5][6]} Jeremy Howard, co-founder of Fast.ai, stated that "[Codex] is a way of getting code written without having to write as much code" and that "it is not always correct, but it is just close enough".^[7] According to a paper written by OpenAI researchers, when attempting each test case 100 times, 70.2% of prompts had working solutions.^[8]

OpenAI claims that Codex is able to function in over a dozen programming languages, including Go, JavaScript, Perl, PHP, Ruby, Shell, Swift, and TypeScript, though it is most effective in Python.^[1] According to VentureBeat, demonstrations uploaded by OpenAI showed impressive coreference resolution capabilities. The demonstrators were able to create a browser game in JavaScript and generate data science charts using matplotlib.^[6]

OpenAI has shown that Codex is able to interface with services and apps such as Mailchimp, Microsoft Word, Spotify, and Google Calendar.^{[6][9]} Microsoft is reportedly interested in exploring Codex's capabilities.^[9]

Issues

OpenAI demonstrations showcased flaws such as inefficient code and one-off quirks in code samples.^[6] In an interview with The Verge, OpenAI chief technology officer Greg Brockman said that "sometimes [Codex] doesn't quite know exactly what you're asking" and that it can require some trial and error.^[9] OpenAI researchers found that Codex struggles with multi-step and higher-level prompts, often failing or yielding counter-intuitive behavior. Additionally, they brought up several safety issues, such as over-reliance by novice programmers, biases based on the training data, and security impacts due to vulnerable code.^[8]

VentureBeat has stated that because Codex is trained on public data, it could be vulnerable to "data poisoning" via intentional uploads of malicious code.^[6] According to a study by researchers from New York University, approximately 40% of code generated by GitHub Copilot (which uses Codex) included glitches or other exploitable design flaws.^[10]

The Free Software Foundation has expressed concerns that code snippets generated by Copilot and Codex could unknowingly violate the terms of free software licenses, such as the GPL, which requires derivative works to be licensed under equivalent terms.^[11] Issues they raised include whether training on public repositories falls into fair use or not, how developers could discover infringing generated code, whether trained machine learning models could be considered modifiable source code or a compilation of the training data, and if machine learning models could themselves be copyrighted and by whom.^{[11][12]} An internal GitHub study found that approximately 0.1% of generated code contained direct copies from the training data. One specific example has been raised, in which the model outputted the original code of the fast inverse square root algorithm, including comments and an incorrect copyright notice.^[4]

In response, OpenAI has stated that "legal uncertainty on the copyright implications of training AI systems imposes substantial costs on AI developers and so should be authoritatively resolved."^[4] The copyright issues with Codex have been compared to the Authors Guild, Inc. v. Google, Inc. court case, in which judges ruled that Google Books's use of text snippets from millions of scanned books constituted fair use.^{[4][13]}

References

1. Zaremba, Wojciech (August 10, 2021). "OpenAI Codex" (<https://openai.com/blog/openai-codex/>). *OpenAI*. Retrieved 2021-09-03.
2. Wiggers, Kyle (July 8, 2021). "OpenAI warns AI behind GitHub's Copilot may be susceptible to bias" (<https://venturebeat.com/2021/07/08/openai-warns-ai-behind-githubs-copilot-may-be-susceptible-to-bias/>). *VentureBeat*. Retrieved 2021-09-03.
3. Alford, Anthony (August 31, 2021). "OpenAI Announces 12 Billion Parameter Code-Generation AI Codex" (<https://www.infoq.com/news/2021/08/openai-codex/>). *InfoQ*. Retrieved 2021-09-03.
4. Anderson, Tim; Quach, Katyanna (July 6, 2021). "GitHub Copilot auto-coder snags emerge, from seemingly spilled secrets to bad code, but some love it" (https://www.theregister.com/2021/07/06/github_copilot_autocoder_caught_spilling/). *The Register*. Retrieved 2021-09-04.
5. Dorrier, Jason (August 15, 2021). "OpenAI's Codex Translates Everyday Language Into Computer Code" (<https://singularityhub.com/2021/08/15/openais-codex-translates-everyday-language-into-computer-code/>). *SingularityHub*. Retrieved 2021-09-03.
6. Dickson, Ben (August 16, 2021). "What to expect from OpenAI's Codex API" (<https://venturebeat.com/2021/08/16/what-to-expect-from-openais-codex-api/>). *VentureBeat*. Retrieved 2021-09-03.
7. Metz, Cade (September 9, 2021). "A.I. Can Now Write Its Own Computer Code. That's Good News for Humans" (<https://www.nytimes.com/2021/09/09/technology/codex-artificial-intelligence-coding.html>). *The New York Times*. Retrieved 2021-09-16.
8. Chen, Mark; Tworek, Jerry; Jun, Heewoo; Yuan, Qiming; Pinto, Henrique Ponde de Oliveira; Kaplan, Jared; Edwards, Harri; Burda, Yuri; Joseph, Nicholas; Brockman, Greg; Ray, Alex (2021-07-14). "Evaluating Large Language Models Trained on Code". [arXiv:2107.03374](https://arxiv.org/abs/2107.03374) (<https://arxiv.org/archive/cs>).
9. Vincent, James (August 10, 2021). "OpenAI can translate English into code with its new machine learning software Codex" (<https://www.theverge.com/2021/8/10/22618128/openai-codex-natural-language-into-code-api-beta-access>). *The Verge*. Retrieved 2021-09-03.
10. Claburn, Thomas (August 25, 2021). "GitHub's Copilot may steer you into dangerous waters about 40% of the time – study" (https://www.theregister.com/2021/08/25/github_copilot_study/). *The Register*. Retrieved 2021-09-03.

11. Krill, Paul (August 2, 2021). "[GitHub Copilot is 'unacceptable and unjust,' says Free Software Foundation](https://www.infoworld.com/article/3627319/github-copilot-is-unacceptable-and-unjust-says-free-software-foundation.html)" (<https://www.infoworld.com/article/3627319/github-copilot-is-unacceptable-and-unjust-says-free-software-foundation.html>). *InfoWorld*. Retrieved 2021-09-03.
 12. Robertson, Donald (2021-07-28). "[FSF-funded call for white papers on philosophical and legal questions around Copilot: Submit before Monday, August 23, 2021](https://www.fsf.org/blogs/licensing/fsf-funded-call-for-white-papers-on-philosophical-and-legal-questions-around-copilot)" (<https://www.fsf.org/blogs/licensing/fsf-funded-call-for-white-papers-on-philosophical-and-legal-questions-around-copilot>). *Free Software Foundation*. Retrieved 2021-09-04.
 13. Barber, Gregory (July 12, 2021). "[GitHub's Commercial AI Tool Was Built From Open Source Code](https://www.wired.com/story/github-commercial-ai-tool-built-open-source-code/)" (<https://www.wired.com/story/github-commercial-ai-tool-built-open-source-code/>). *WIRED*. Retrieved 2021-09-04.
-

Retrieved from "https://en.wikipedia.org/w/index.php?title=OpenAI_Codex&oldid=1068813052"

This page was last edited on 30 January 2022, at 09:14 (UTC).

Text is available under the Creative Commons Attribution-ShareAlike License 3.0; additional terms may apply. By using this site, you agree to the Terms of Use and Privacy Policy. Wikipedia® is a registered trademark of the Wikimedia Foundation, Inc., a non-profit organization.