

---

# Fine-Tuning Codex Toward Optimal Solutions for Sorting and Searching Problems

---

Sourya Kakarla, Derek Pham, Gabriele Tombesi  
Department of Computer Science  
Columbia University  
New York, NY 10027

## 1 Problem, Formulation, and Goal

OpenAI Codex [7] is a language model obtained by training OpenAI’s base GPT-3 [4]—which is pre-trained on web and text resources—on a vast set of public code repositories retrieved from GitHub. As a result of this additional training process, Codex has become reasonably proficient generating code in several programming languages by simply parsing the input natural language instructions. Typical of the third generation of generative transformers, Codex inherits few-shot, one-shot and zero-shot learning capabilities, intuiting a wide variety of downstream tasks just from the input settings, without need of additional instructions and with few to no examples.

The research problem that we formulate in this work is based on the observation that, while Codex is capable of achieving state-of-the-art results for basic Python program synthesis tasks, performance significantly drops when the complexity of the prompt’s problem and consequently the corresponding programming solution increases. More specifically, Codex can produce outputs with time complexities which are several orders of magnitude higher than the corresponding optimal solutions.

Given this fact, we are interested in exploring solutions to improve the quality of results (QoR) of Codex-generated Python code. Given the complexity of the models leveraged by OpenAI and the limited computational resources at our disposal, we decide to break down this task into manageable pieces and set one of those as our research goal: we will focus on improving the computational efficiency of generated code for a specific subset of tasks, namely sorting and searching algorithms.

We choose to focus on improving the computational efficiency of generated code specifically for sorting and searching algorithms, because we are targeting tasks that are of medium-level difficulty, for which Codex is able to generate functionally correct solutions, but not optimal ones. For this reason, it would not make sense to focus on trivial tasks where Codex is already proficient, or on complex tasks where Codex is not able to handle yet, namely those where Codex cannot generate functionally correct solutions at this point in time. Nevertheless, even though some of our preliminary experiments have shown that this task set is promising for fine-tuning Codex, we want to stress that we are open to including harder prompts beyond search and sorting in the case too many problems based on those algorithms turn out to be already optimally solved by Codex. For example, we might include harder coding questions involving dynamic programming.

## 2 Method, Dataset, and Evaluation Criteria

### 2.1 Method and Dataset

We intend to leverage the fine-tuning API [6] offered by OpenAI. The basic idea of our approach is the following: given a set of prompts in our domain of interest, we will need to collect several functionally correct Python implementations with varying time complexity for each of those tasks. Given this

34 initial raw version of the dataset—still not compliant to what Codex expects for fine-tuning—we intend  
 35 to build a fine-tuning dataset composed of *[prompt:completions]* example pairs such that, for each of  
 36 the original prompts in the target tasks set, we generate as many examples as the number of output  
 37 implementations that we have for that task. It is crucial though, in order for the examples to be  
 38 distinct and for the fine-tuning to improve the model, that the original prompt is concatenated with an  
 39 additional token that embeds the optimality, in terms of complexity, of the corresponding completion  
 40 (i.e. Python solution). A summary of this second stage of the dataset preparation is illustrated in  
 41 Figure 1.

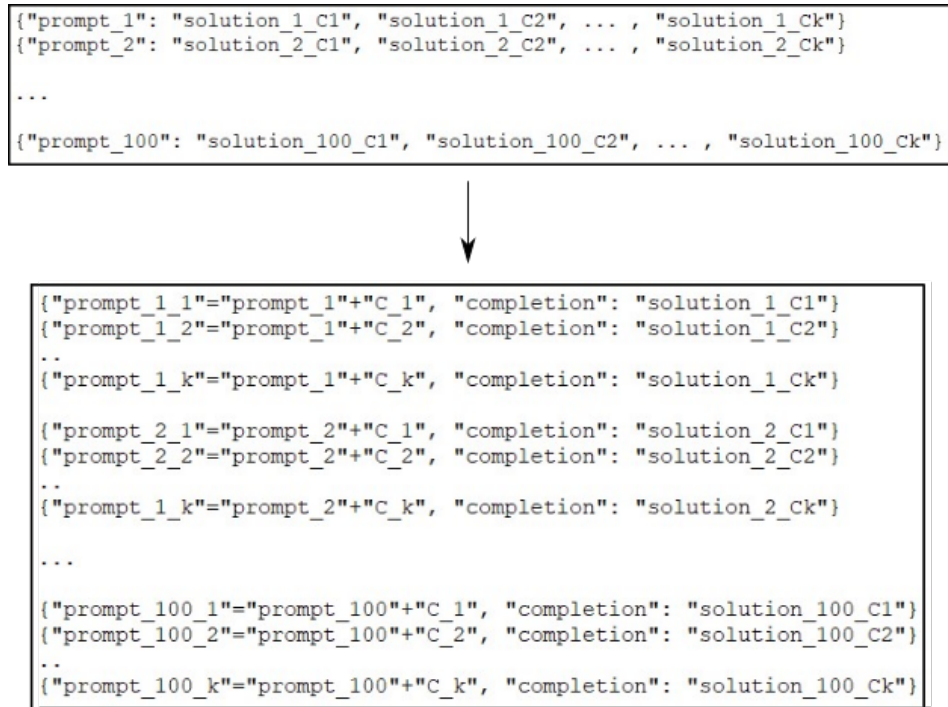


Figure 1: From raw dataset to Codex-compliant dataset

42 At this point, having the compliant data-set saved in the *fine\_tune.json* file and setting code-davinci-  
 43 001 [8] as the pretrained model that we intend to customize, we will proceed by passing the dataset  
 44 to the fine-tuning API call. These steps, together with the openAI API call at evaluation phase, are  
 45 summarized in figure 2.

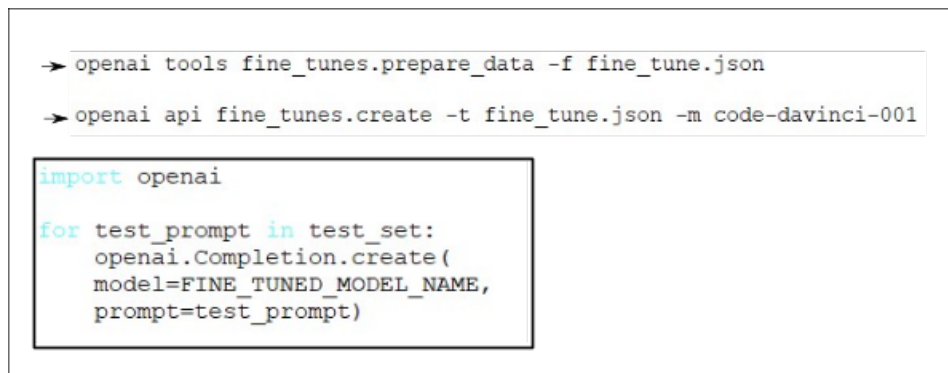


Figure 2: Fine-Tuning call and Evaluation

46 The novelty of this approach is that we intend to perform a complexity-biased fine-tuning based on  
47 an additional input tokens indicating complexity, which are meant to teach the model what it means  
48 to generate an optimal solution for this specific subset of tasks that we target.

49 Finally, we note that a substantial amount of time for this research project will be dedicated toward  
50 collecting a suitable dataset with solutions of different efficiency levels for each task. We realize that  
51 this process may not be straightforward, and so under the circumstance in which we face difficulty in  
52 reaching our minimum goal of 100 examples, we are considering to integrate our search for existing  
53 solutions using the following method: We can use Codex—or another accessible code-generating  
54 model—to generate different implementations of a task with different complexity by leveraging  
55 parameters such as Temperature, which would allow us to generate more variance in the output. We  
56 could then feed the outputs into a benchmark tool to first discard the incorrect solutions, and then rank  
57 the remaining solutions based on time complexity. This way, we could create new prompt-solution  
58 pairs across different complexity levels to use for fine-tuning.

## 59 **2.2 Evaluation Criteria**

60 We will use two evaluation metrics for our evaluation criteria to see if there is any performance  
61 improvement with respect to the pre-fine-tuned model. The first evaluation metric is the percentage  
62 of prompts for which the fine-tuned model outputs correct programs. It is of interest to compare the  
63 percentage with the original engine of Codex and observe if there is any significant deviation due to  
64 our fine tuning.

65 Secondly, the metric with which we are mainly concerned about is the one that evaluates how the  
66 generated code performs when compared to the known gold solution. We plan to have a gold  
67 solution for every prompt in the test-set that has the best known complexity and high performing  
68 implementation. These solutions are human generated and evaluated. For all the valid codes generated  
69 by the engine for a particular prompt, we would evaluate how optimal they are by using distance  
70 measures with the gold solution’s practical benchmarks and theoretical complexity.

71 We will use a benchmark module to measure the average time taken by the generated code for  
72 different sufficiently large inputs. Multiple runs are done to mitigate noise from the OS environment.  
73 This time measure is compared to the gold solutions’ time measure, which is obtained in the same  
74 fashion, and then a suitable distance measure is calculated. The baseline for this metric can be  
75 computed by measuring average time taken through the same method for the codes generated by the  
76 original engine that is not fine-tuned. By evaluating the difference between the computed benchmark  
77 of the code generated from the fine tuned model and that of the baseline, we aim to survey the success  
78 of our fine-tuned model.

## 79 **3 Previous Work**

80 Henraks et al. (2021) have created APPS - Automated Programming Progress Standard, a benchmark  
81 for code generation [5]. This benchmark measures the ability of models to take an arbitrary natural  
82 language specification and generate satisfactory Python code. This benchmark includes 10,000  
83 problems, which range from having simple one-line solutions to being substantial algorithmic  
84 challenges.

85 We take inspiration from Chen et al. (2021) who introduce Codex and evaluate its capabilities in  
86 generating Python code [2]. Moreover, Codex is fine-tuned from GPT-3, which is an autoregressive  
87 language model that achieves strong performance across a variety of natural-language-generating  
88 tasks, including translation and question-answering [1]. A significant portion of the existing language  
89 models for code-related tasks developed from the GPT architecture. Among them, IntelliCode  
90 Compose leverages the autoregressive approach of GPT, applying it to the field of source code  
91 understanding, namely multilingual code completion tasks [9].

92 In parallel to this line of work based on the most advanced auto-regressive models implemented so far,  
93 there are several examples of high-performance code-generation applications based on auto-encoding

94 models. Among them, it is worth mentioning CodeBERT: an NL-PL pretrained model customized  
95 for natural language code search and code documentation generation [3].  
96 Finally, the Code Transformer is a multilingual code summarization model that leverages the encoder-  
97 decoder architecture of sequence-to-sequence transformers, learning from both the source code and  
98 parsed abstract syntax trees [10].

## References

- [1] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. *CoRR*, abs/2005.14165, 2020.
- [2] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harrison Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. *CoRR*, abs/2107.03374, 2021.
- [3] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming and natural languages. *CoRR*, abs/2002.08155, 2020.
- [4] Luciano Floridi and Massimo Chiriatti. Gpt-3: Its nature, scope, limits, and consequences. *Minds and Machines*, 30(4):681–694, 2020.
- [5] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, et al. Measuring coding challenge competence with apps. *arXiv preprint arXiv:2105.09938*, 2021.
- [6] OpenAI. Fine-tuning. <https://beta.openai.com/docs/guides/fine-tuning>, 2022. [Online; accessed 10-March-2022].
- [7] OpenAI. Openai codex. <https://openai.com/blog/openai-codex/>, 2022. [Online; accessed 10-March-2022].
- [8] OpenAI. Openai codex engines. <https://beta.openai.com/docs/engines/codex-series-private-beta>, 2022. [Online; accessed 10-March-2022].
- [9] Alexey Svyatkovskiy, Shao Kun Deng, Shengyu Fu, and Neel Sundaresan. Intellicode compose: Code generation using transformer. *CoRR*, abs/2002.08155, 2020.
- [10] Daniel Zügner, Tobias Kirschstein, Michele Catasta, Jure Leskovec, and Stephan Günnemann. Language-agnostic representation learning of source code from structure and context. *CoRR*, abs/2103.11318, 2021.