# LLM Code Generation

Willis Reid
April 2024
University of North Carolina
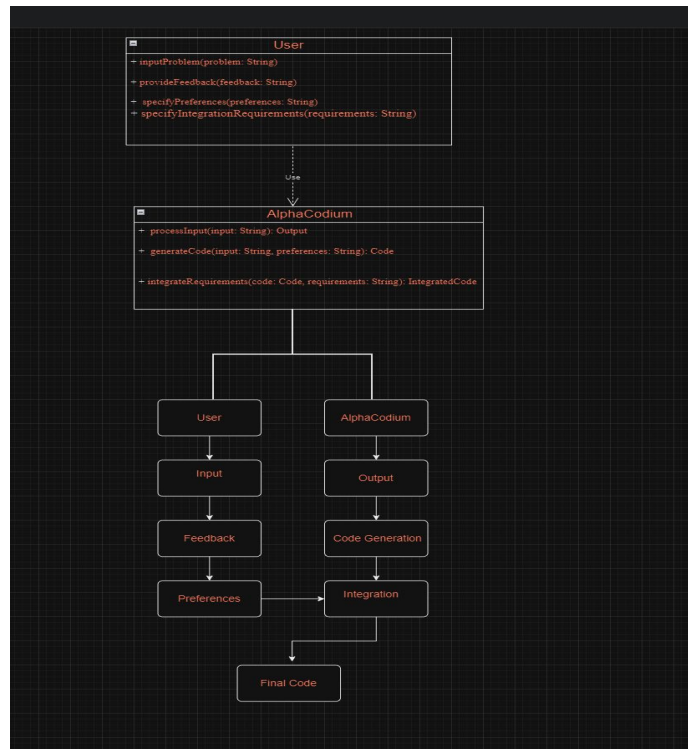
# Table of Content

# Motivation

I chose this topic because I found it interesting and I believe many people will enjoy it too. As I delve deeper and examine the complexity of this topic, my findings continue to astound me. My main reason for choosing this topic is my fascination with it, and I'm eager to share the intriguing thoughts and discoveries I come across.
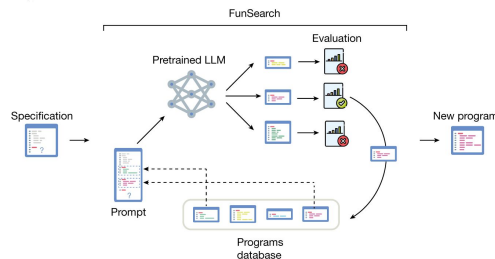
# Existing Related Approaches

- Created a user diagram to illustrate project workflow.

- Executed terminal commands to process and run code efficiently.

- Procured the dataset required for the project from Hugging Face.

# Method # 1

Evaluating Large Language Models Trained on Code.

|  | INTRODUCTORY | INTERVIEW | COMPETITION |
|---|---|---|---|
| GPT-NEO 2.7B RAW PASS@1 | 3.90% | 0.57% | 0.00% |
| GPT-NEO 2.7B RAW PASS@5 | 5.50% | 0.80% | 0.00% |
| 1-SHOT CODEX RAW PASS@1 | 4.14% (4.33%) | 0.14% (0.30%) | 0.02% (0.03%) |
| 1-SHOT CODEX RAW PASS@5 | 9.65% (10.05%) | 0.51% (1.02%) | 0.09% (0.16%) |
| 1-SHOT CODEX RAW PASS@100 | 20.20% (21.57%) | 2.04% (3.99%) | 1.05% (1.73%) |
| 1-SHOT CODEX RAW PASS@1000 | 25.02% (27.77%) | 3.70% (7.94%) | 3.23% (5.85%) |
| 1-SHOT CODEX FILTERED PASS@1 | 22.78% (25.10%) | 2.64% (5.78%) | 3.04% (5.25%) |
| 1-SHOT CODEX FILTERED PASS@5 | 24.52% (27.15%) | 3.23% (7.13%) | 3.08% (5.53%) |

# Method # 2

CodeChain: Improving Code Generation through Modular Self-Revision



Figure 2: An overview of CodeChain: a pretrained LLM is first instructed with chain-of-thought prompting to generate a set of modularized solutions. Generated sub-modules are then extracted from potentially correct solutions and grouped into different semantic clusters. The cluster centroids are selected as representative sub-modules to condition the next self-revision round. The model is instructed to reuse or adapt these modules into its revised solutions.

# Method That Was Duplicate

AlphaCodium



(a) The proposed AlphaCodium flow.



(b) Illustrating the improvement from AlphaCodium.

# Results

| model | set | method | score (pass@5) |
|---|---|---|---|
| DeepSeek -33B [3] | validation | Direct | 7% |
| | | AlphaCodium | **20%** |
| | test | Direct prompt | 12% |
| | | AlphaCodium | **24%** |
| GPT-3.5 | validation | Direct prompt | 15% |
| | | AlphaCodium | **25%** |
| | test | Direct prompt | 8% |
| | | AlphaCodium | **17%** |
| GPT-4 | validation | Direct prompt | 19% |
| | | AlphaCodium | **44%** |
| | test | Direct prompt | 12% |
| | | AlphaCodium | **29%** |

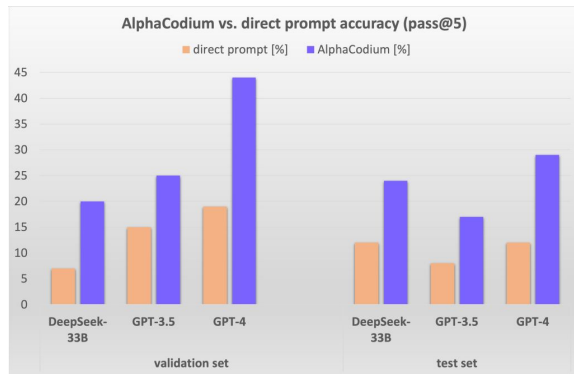| model | set | method | score |
|---|---|---|---|
| GPT-3.5 | validation | AlphaCodium (pass@5) | **25%** |
| | | CodeChain (pass@5) | 17% |
| | test | AlphaCodium (pass@5) | **17%** |
| | | CodeChain (pass@5) | 14% |
| GPT-4 | validation | AlphaCodium (pass@5) | **44%** |
| AlphaCode | | AlphaCode (pass@10@1K) | 17% |
| | | AlphaCode (pass@10@100K) | 24% |
| GPT-4 | test | AlphaCodium (pass@5) | **29%** |
| AlphaCode | | AlphaCode (pass@10@1K) | 16% |
| | | AlphaCode (pass@10@100K) | 28% |

# Results

Compares AlphaCodium's performance with other models like AlphaCode and CodeChain.

Valuation Metric: Utilizes the pass@k metric to quantify the percentage of successfully solved problems using k generated solutions per problem.

AlphaCodium: demonstrates its effectiveness in solving problems compared to single, well-designed direct prompts.

AlphaCode and CodeChain: Represents other models used for comparison in the analysis.

# Future Work

- Identified human errors in coding through dataset.
- Curated dataset showcasing common coding mistakes and associated code snippets.
- Parsed text to detect errors and visualized alongside code outputs.
- Aimed to provide clear understanding of coding mistakes through visualization.

# Conclusion

- Many people fear the impact of machine learning in our field.

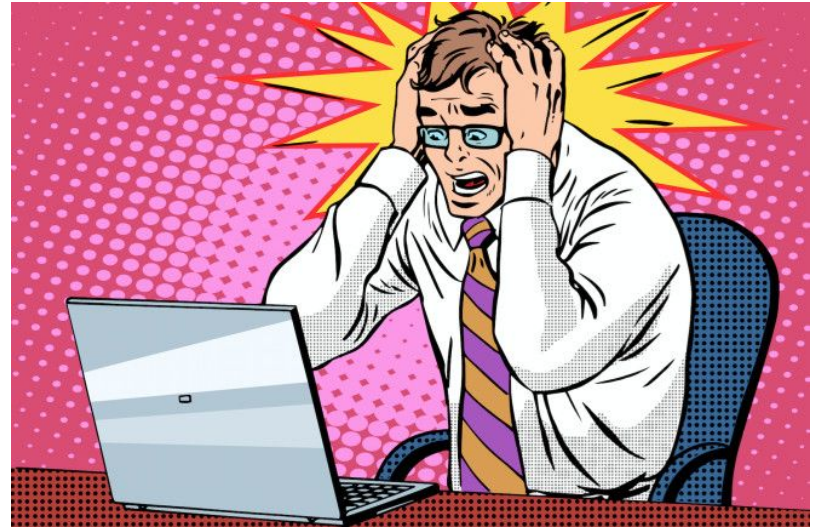- Common concern: Fear of job loss or inability to find work.

- Clarification: Not necessarily true.

- Emphasize the potential for growth and adaptation in the industry.

# Embracing Fear

https://www.youtube.com/watch?v=AgyJv2Qelwk&ab_channel=Fireship

# Resources

Ridnik, Tal, et al. "Code Generation with Alphacodium: From Prompt Engineering to Flow Engineering." *arXiv.Org*, 16 Jan. 2024, arxiv.org/abs/2401.08500. Accessed 23 Feb. 2024. ( https://arxiv.org/abs/2401.08500)

Chen, Mark, et al. "Evaluating Large Language Models Trained on Code." *arXiv.Org*, 14 July 2021, arxiv.org/abs/2107.03374. Accessed 23 Feb. 2024. (https://arxiv.org/abs/2107.03374)

Le, Hung, et al. "CodeChain: Towards Modular Code Generation through Chain of Self-Revisions with Representative Sub-Modules." *arXiv.Org*, 28 Nov. 2023, arxiv.org/abs/2310.08992. Accessed 23 Feb. 2024. ( https://arxiv.org/abs/2310.08992)