# Code generation with LLMs

Generative AI && software engineering:
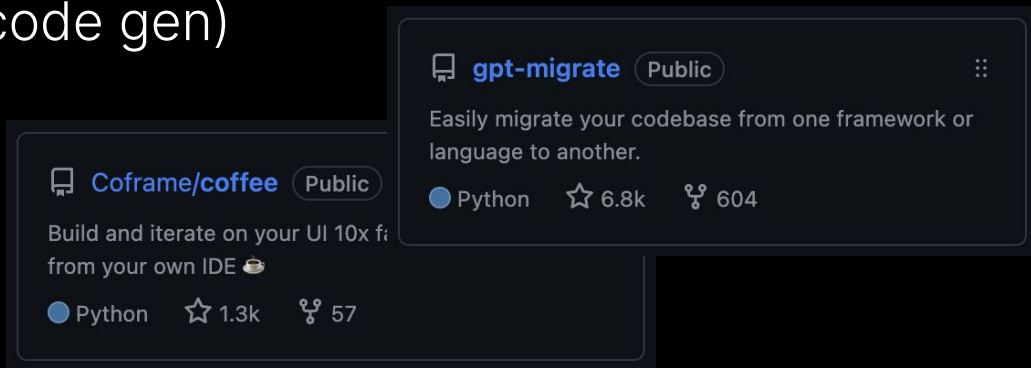analysis, learnings, practical insights

Josh Payne

# Agenda

- Intro

- Brief history of AI for code generation

- Benchmarking code gen performance

- Applications and agents

- AI x software engineering

# Intro

- 👋 I'm Josh

- Founder of Coframe (AI for UI optimization + code gen), prev two other companies (one AI-focused)

- Created GPT-Migrate (LLM-powered codebase migration), Coffee (LLM-powered UI code gen)

- Stanford CS (AI) alum!

🖥 **gpt-migrate** `Public`

Easily migrate your codebase from one framework or language to another.

🔵 Python   ⭐ 6.8k   ⑂ 604

🖥 **Coframe/coffee** `Public`

Build and iterate on your UI 10x fa... from your own IDE ☕

🔵 Python   ⭐ 1.3k   ⑂ 57

# Brief History



CodeNN (Iyer et al., 2016)

Code summarization



Aroma (Luan et al, 2019)

Code search (early copilot)



Code2Seq (Alon et al., 2019)

Better code summarization

(Try it! → https://code2seq.org/ )

**Pre-LLM era:**
RNNs and search

**Early applications:**
GPT-3, Codex, GitHub Copilot

**"Oh wow, AI can actually write code now":**
GPT-3.5, GPT-4, OSS LLMs

**AI x software engineering:**
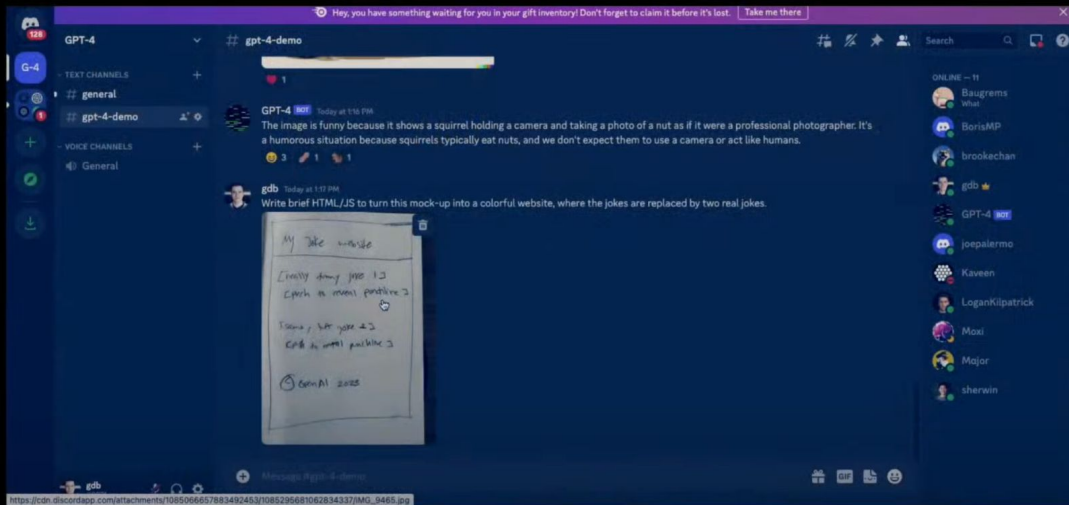Agents and integrated workflows

# Brief History



**Pre-LLM era:**
RNNs and search

**Early applications:**
GPT-3, Codex, GitHub Copilot

**"Oh wow, AI can actually write code now":**
GPT-3.5, GPT-4, OSS LLMs

**AI x software engineering:**
Agents and integrated workflows
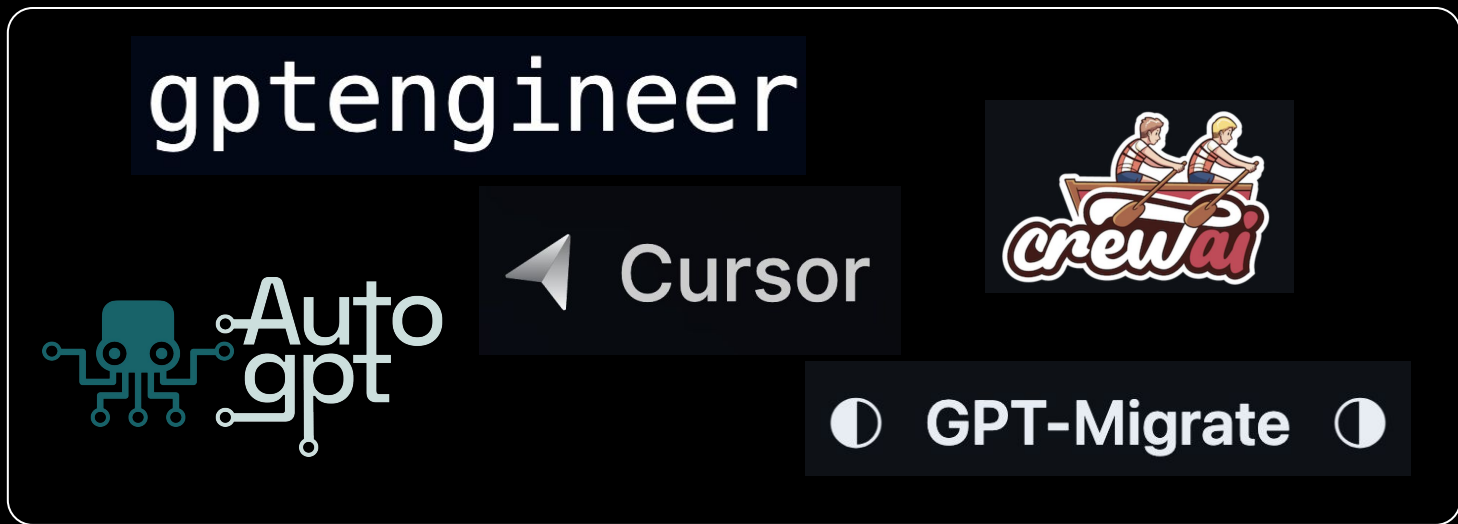
# Brief History



Pre-LLM era:
RNNs and search

Early applications:
GPT-3, Codex, GitHub Copilot

"Oh wow, AI can actually write code now":
GPT-3.5, GPT-4, OSS LLMs

AI x software engineering:
Agents and integrated workflows

# Brief History

Still in its infancy!



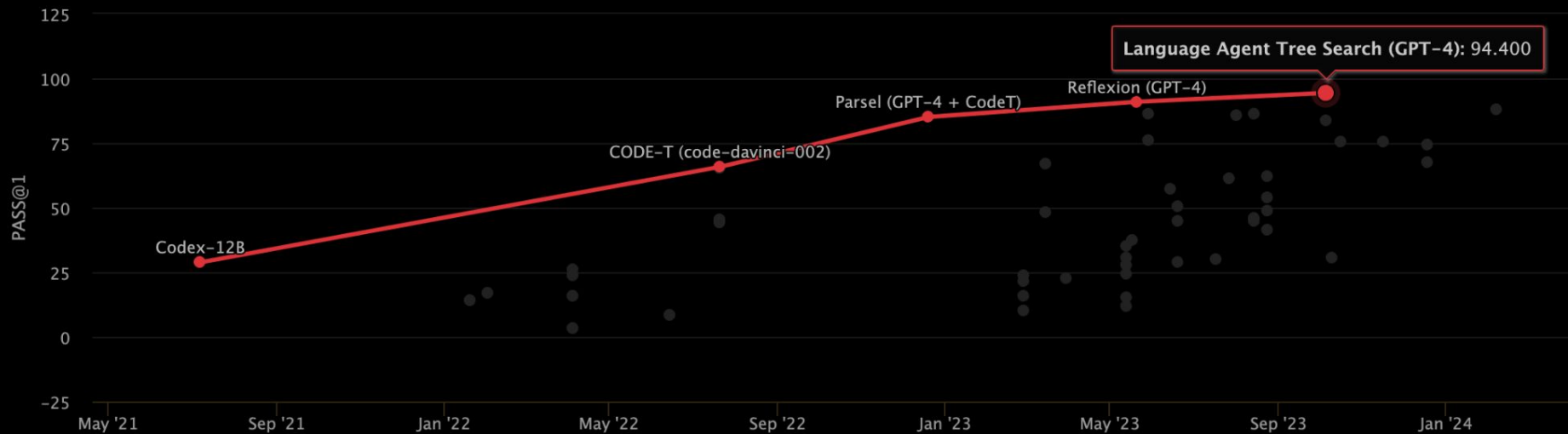Pre-LLM era:
RNNs and search

Early applications:
GPT-3, Codex, GitHub Copilot

"Oh wow, AI can actually write code now":
GPT-3.5, GPT-4, OSS LLMs

AI x software engineering:
Agents and integrated workflows

# Benchmarking code generation

# How do we measure this?

**(1)** Benchmark Tasks    **(2)** Competitions    **(3)** Real-world impact

**HumanEval (Chen et al., 2021) is the most widely-recognized research benchmark for code generation.**

This paper also introduced **Codex**, the first major code-specific LLM.

HumanEval is 164 handwritten programming problems, each with several unit tests.

```python
def incr_list(l: list):
    """Return list with elements incremented by 1.
    >>> incr_list([1, 2, 3])
    [2, 3, 4]
    >>> incr_list([5, 3, 5, 2, 3, 3, 9, 0, 123])
    [6, 4, 6, 3, 4, 4, 10, 1, 124]
    """
    return [i + 1 for i in l]
```

**The prompt provided to the model is shown with a black background, and a successful model-generated completion is shown in a blue background. To be successful, it must pass the unit tests.**

# How do we measure this?

**1** Benchmark Tasks   **2** Competitions   **3** Real-world impact

**There have also been extensions of HumanEval and other datasets:**

- **MultiPL-E** is a dataset for evaluating large language models for code generation that supports 18 programming languages. It translates HumanEval problems into other languages.

- **HumanEval-X** consists of 820 high-quality human-crafted data samples, compared with HumanEval's 164.

- **MBPP** (Mostly Basic Python Problems) is a dataset of 1000 crowd-sourced Python programming problems.
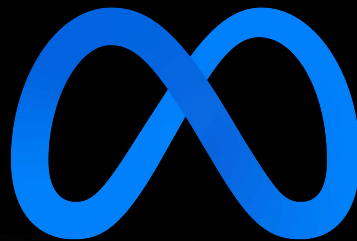
# How do we measure this?

① Benchmark Tasks  ② Competitions  ③ Real-world impact

**Some companies will create internal datasets on which to evaluate.**

- Google introduced Gemini alongside a new benchmark, Natural2Code, which is a held-out internal dataset.

   **GPT-4 (OpenAI) was slightly better on HumanEval (OpenAI), while Gemini (Google) was slightly better on Natural2Code (Google).**

- Meta has internal unit test sets for its internal LLMs.

**TestGen-LLM**

Gemini

# How do we measure this?

1 Benchmark Tasks    2 Competitions    3 Real-world impact

**Why are held-out (non-published) benchmarks valuable?**

# How do we measure this?

(1) Benchmark Tasks    (2) Competitions    (3) Real-world impact

**AlphaCode by DeepMind (Li et al., Dec 2022) created CodeContests, a dataset of compiled competitive programming problems.**

Increasingly, datasets from real-world tasks for humans are needed as models approach human-level performance.

Other examples: the LSAT, USMLE, AlphaGeometry (IMO problems)

## CodeContests

CodeContests is a competitive programming dataset for machine-learning. This dataset was used when training AlphaCode. AlphaCode has been published in Science, with a preprint on arXiv.

It consists of programming problems, from a variety of sources:

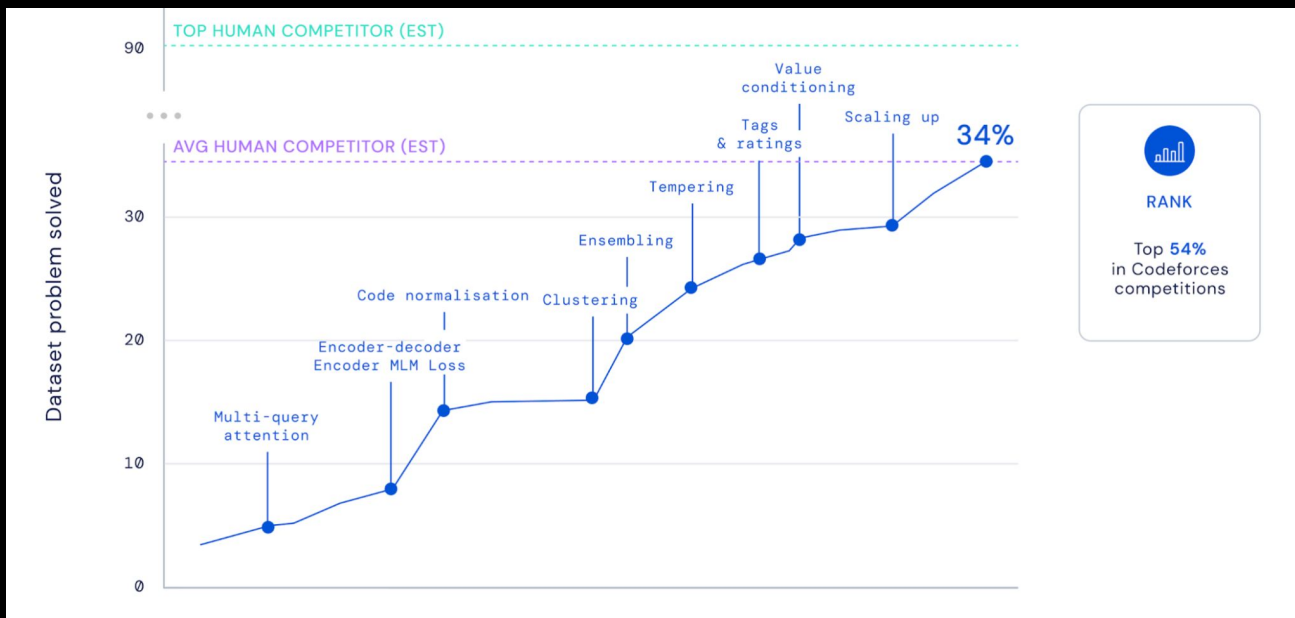| Site | URL | Source |
|---|---|---|
| Aizu | https://judge.u-aizu.ac.jp | CodeNet |
| AtCoder | https://atcoder.jp | CodeNet |
| CodeChef | https://www.codechef.com | description2code |
| Codeforces | https://codeforces.com | description2code and Codeforces |
| HackerEarth | https://www.hackerearth.com | description2code |

# How do we measure this?

① Benchmark Tasks    ② Competitions    ③ Real-world impact

# How do we measure this?

① Benchmark Tasks  ② Competitions  **③** Real-world impact

**As models begin to surpass human performance, they will be increasingly measured on impact.**

Example: AlphaDev (Mankowitz and Michi, June 2023) discovered a faster sorting algorithm for small lists that has now been implemented in the C++ standard lib.

SWE KPIs (bug rate, PRs merged, etc) are starting to become more commonplace.

Original

```
Memory[0] = A
Memory[1] = B
Memory[2] = C

mov Memory[0] P  // P = A
mov Memory[1] Q  // Q = B
mov Memory[2] R  // R = C

mov R S
cmp P R
cmovg P R  // R = max(A, C)
cmovl P S  // S = min(A, C)
mov S P    // P = min(A, C)
cmp S Q
cmovg Q P  // P = min(A, B, C)
cmovg S Q // Q = max(min(A, C), B)

mov P Memory[0]  // = min(A, B, C)
mov Q Memory[1]  // = max(min(A, C), B)
mov R Memory[2]  // = max(A, C)
```

AlphaDev

```
Memory[0] = A
Memory[1] = B
Memory[2] = C

mov Memory[0] P  // P = A
mov Memory[1] Q  // Q = B
mov Memory[2] R  // R = C

mov R S
cmp P R
cmovg P R  // R = max(A, C)
cmovl P S  // S = min(A, C)

cmp S Q
cmovg Q P  // P = min(A, B)
cmovg S Q // Q = max(min(A, C), B)

mov P Memory[0]  // = min(A, B)
mov Q Memory[1]  // = max(min(A, C), B)
mov R Memory[2]  // = max(A, C)
```

**Left:** The original implementation with min(A,B,C).

**Right:** AlphaDev Swap Move - AlphaDev discovers that you only need min(A,B).

# Benchmarking code generation

Techniques

Fine-tuning /
Instruct-tuning

Base models

# Benchmarking code generation

**Base Models** are the GPTs and Llamas of the world: not fine-tuned for a particular task.

```
          /\
         /  \
        /Techniques\
       /------------\
      /Fine-tuning / \
     /Instruct-tuning\
    /------------------\
   /    Base models     \
  /_____\
```
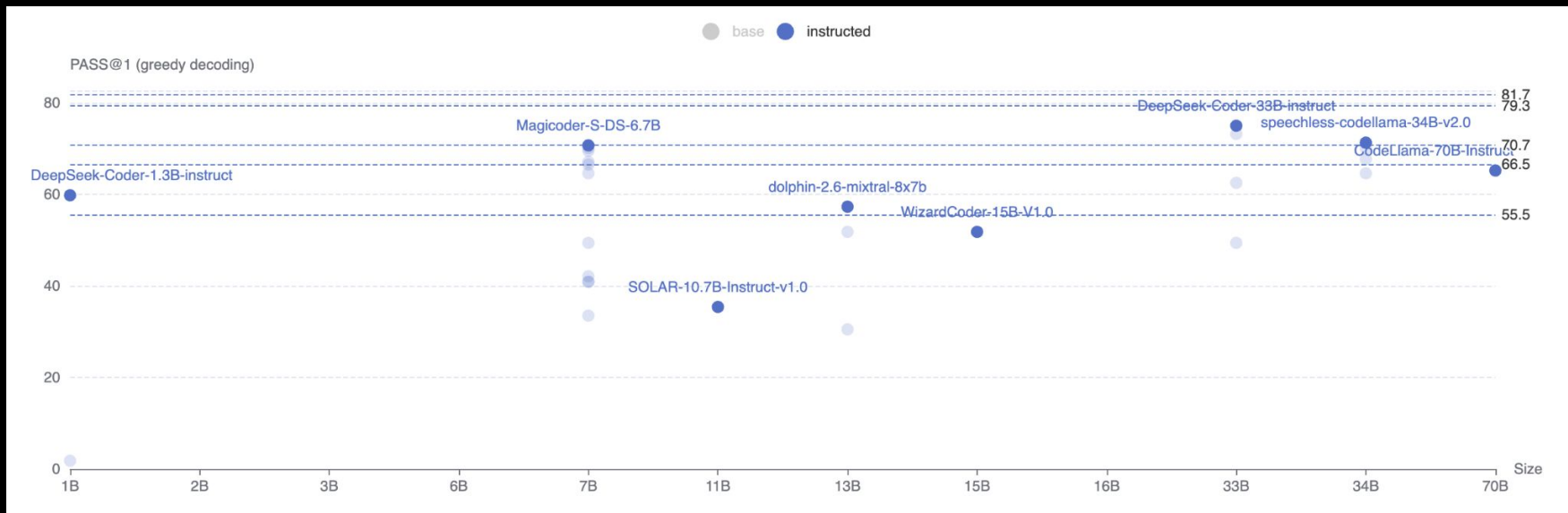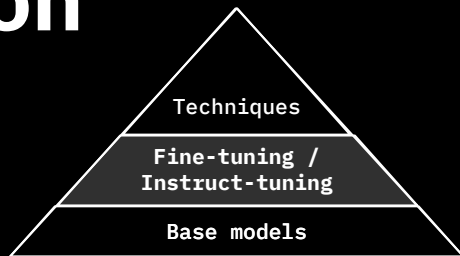
## Open LLMs

*Weights are open, easy to do custom tuning and experimentation*

- CodeLlama (WizardCoder)
- StarCoder
- Replit-code-v1-3b
- Mixtral-8×7b

## Closed LLMs

*Weights are closed, tuning and experimentation are limited*

- GPT-4
- Gemini Ultra
- Claude 2.1
- Grok

# Benchmarking code generation

**Base Models** are the GPTs and Llamas of the world:
not fine-tuned for a particular task.

Techniques

Fine-tuning / Instruct-tuning

**Base models**



● **base**   ● instructed

PASS@1 (greedy decoding)

Phind-CodeLlama-34B-v2

CodeLlama-70B

phi-2-2.7B

Mistral-codealpaca-7B

DeepSeek-Coder-33B-base

Gemma-7B

CodeLlama-13B

DeepSeek-Coder-1.3B-base

CodeGen-6B

StarCoder-15B

CodeGen-16B

CodeGen-2B

Size

1B  2B  3B  6B  7B  11B  13B  15B  16B  33B  34B  70B

# Benchmarking code generation

**Instruct-tuned models** are models that are
fine-tuned with instructions: in this case, for code.

Techniques

Fine-tuning /
Instruct-tuning

Base models



PASS@1 (greedy decoding)

base   instructed

- DeepSeek-Coder-1.3B-instruct
- Magicoder-S-DS-6.7B
- DeepSeek-Coder-33B-instruct
- speechless-codellama-34B-v2.0
- CodeLlama-70B-Instruct
- dolphin-2.6-mixtral-8x7b
- WizardCoder-15B-V1.0
- SOLAR-10.7B-Instruct-v1.0

81.7
79.3
70.7
66.5
55.5

Size
1B  2B  3B  6B  7B  11B  13B  15B  16B  33B  34B  70B

# Benchmarking code generation

Techniques
Fine-tuning / Instruct-tuning
Base models

✅ **Benchmark: HumanEval**

**Instruct-tuned models** are models that are fine-tuned with instructions: in this case, for code.

Instruct-tuning involves a prompt which contains an instruction, and a response. Including the instruction is important for the model to know how to understand new instructions at inference time.

**Example: Synthesis**

**Model Input**

**Target Output**

```
Write a Python function `has_close_elements(numbers: List[float],
threshold: float) -> bool` to solve the following problem:
Check if in given list of numbers, are any two numbers closer to
each other than given threshold.
>>> has_close_elements([1.0, 2.0, 3.0], 0.5)
False
>>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
True

from typing import List


def has_close_elements(numbers: List[float], threshold: float) ->
bool:
    """ Check if in given list of numbers, are any two numbers closer
to each other than given threshold.
    >>> has_close_elements([1.0, 2.0, 3.0], 0.5)
    False
    >>> has_close_elements([1.0, 2.8, 3.0, 4.0, 5.0, 2.0], 0.3)
    True
    """
    for idx, elem in enumerate(numbers):
        for idx2, elem2 in enumerate(numbers):
            if idx != idx2:
                distance = abs(elem - elem2)
                if distance < threshold:
                    return True

    return False
```
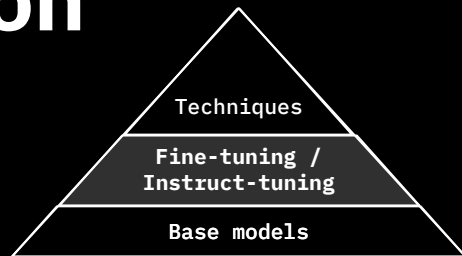
# Benchmarking code generation

✅ **Benchmark: HumanEval**

**Instruct-tuned models** are models that are fine-tuned with instructions: in this case, for code.

Instruct-tuning involves a prompt which contains an instruction, and a response. Including the instruction is important for the model to know how to understand new instructions at inference time.

**Example: Fix a bug**

■ **Model Input**

■ **Target Output**

```python
from typing import List

def has_close_elements(numbers: List[float], threshold: float) ->
bool:
    for idx, elem in enumerate(numbers):
        for idx2, elem2 in enumerate(numbers):
            if idx != idx2:
                distance = elem - elem2
                if distance < threshold:
                    return True

    return False

def check(has_close_elements):
    assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.3) ==
True
    assert has_close_elements([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.05) ==
False
    assert has_close_elements([1.0, 2.0, 5.9, 4.0, 5.0], 0.95) == True
    assert has_close_elements([1.0, 2.0, 5.9, 4.0, 5.0], 0.8) == False
    assert has_close_elements([1.0, 2.0, 3.0, 4.0, 5.0, 2.0], 0.1) ==
True
    assert has_close_elements([1.1, 2.2, 3.1, 4.1, 5.1], 1.0) == True
    assert has_close_elements([1.1, 2.2, 3.1, 4.1, 5.1], 0.5) == False

check(has_close_elements)

Fix bugs in has_close_elements.
from typing import List

def has_close_elements(numbers: List[float], threshold: float) ->
bool:
    for idx, elem in enumerate(numbers):
        for idx2, elem2 in enumerate(numbers):
            if idx != idx2:
                distance = abs(elem - elem2)
                if distance < threshold:
                    return True

    return False
```

# Benchmarking code generation

✅ **Benchmark: HumanEval**

**Instruct-tuned models** are models that are fine-tuned with instructions: in this case, for code.

Instruct-tuning involves a prompt which contains an instruction, and a response. Including the instruction is important for the model to know how to understand new instructions at inference time.

**Example: Explain code**



■ Model Input

■ Target Output

# Benchmarking code generation

**Instruct-tuning is clearly useful. How can we scale it up?**

As LLMs and datasets get larger, we increasingly need to think creatively about how to gather data in order to improve.

One example of this is COMMITPACK: 4 terabytes of Git commits across 350 programming languages (Muennighoff et al, Jan 2024; ICLR preprint).

Git commits naturally pair code changes with human instructions.

Techniques

Fine-tuning / Instruct-tuning

Base models

# Benchmarking code generation

Technique can make all the difference. This is broadly broken down into **reasoning methods** and **decision-making methods**.

Techniques

Fine-tuning / Instruct-tuning

Base models

83.8 — — — — — — — — — — — — — ●  GPT-3.5 (175B parameters)
Technique: LATS

79.3 — ⬤  GPT-4 (1.7T parameters (est))
Technique: None

💡 **GPT-3.5 beats GPT-4 with LATS, despite being 10x smaller!**

# Benchmarking code generation

✅ **Benchmark: HumanEval**

## Chain-of-Thought

**Reasoning Method**

Chain-of-Thought (CoT) prompts LLMs to sequentially generate reasoning steps from input to output. It was first introduced in PaLM: Scaling Language Modeling with Pathways. (Chowdhery, Catasta et al., 2022)

**However, it suffers from error propagation as the chain length increases.**

Input

S1

S2

S3

Output

# Benchmarking code generation

✅ **Benchmark: HumanEval**

## Tree-of-Thoughts

**Reasoning Method**

Tree-of-Thoughts (ToT) extends CoT by exploring multiple reasoning paths using search algorithms like BFS and DFS. (Yao et al., May 2023)

**That said, it is limited by relying solely on the LLM's internal knowledge.**

# Benchmarking code generation

✅ **Benchmark: HumanEval**

## Reasoning via Planning

**Reasoning Method**

Reasoning via Planning (RAP) (Hao et al., October 2023) uses Monte Carlo Tree Search for planning chains of reasoning.

**However, it also lacks external feedback.**

# Benchmarking code generation

✅ **Benchmark: HumanEval**

## ReAct

**Decision-making method**

ReAct prompts LLMs with alternating actions and observations for decision-making in interactive environments. (Yao et al., March 2023)

**However, it greedily follows one trajectory and cannot adapt.**

Input

Observation

S1

S2

S3

Output

# Benchmarking code generation

✅ **Benchmark: HumanEval**

## Reflexion

**Decision-making method**

Reflexion adds self-reflection to ReAct. This improves overall performance by allowing the LLM more time to think through the problem, similar to CoT. (Shinn et al., October 2023)

**However, it does not consider alternative options at each step.**

Input

Observation

Reflection

S1

S2

S3

Output

# Benchmarking code generation

## Language Agent Tree Search

**Decision-making method**   **Reasoning Method**

LATS unifies the strengths of both reasoning and decision-making methods through principled search, while overcoming limitations via environmental feedback and self-reflection. (Zhou et al., December 2023)

**GPT-4 + LATS is the current best performer on the HumanEval benchmark, with a score of 94.4.**

# Applications and agents

SOFTWARE IS
EATING THE WORLD,
BUT AI IS GOING
TO EAT SOFTWARE

*Jensen Huang | Nvidia CEO*

# Applications and agents

AI has tackled every aspect of software engineering.
(Category list below not exhaustive.)

**replit**

**gptengineer**

**Coffee by**
coframe

**Project
creation**

●● SECOND

◑ GPT-Migrate ◑

/c.

Migrations

Cursor

warp

**Cody by**
Sourcegraph

IDE

Dosu

Sweep AI

codium ai

Issues &
tests

swim

Mintlify

docify

Docs

# Applications and agents

AI has tackled every aspect of software engineering.
(Category list below not exhaustive.)

**replit**

**gptengineer**

**Coffee by**
coframe

**Project
creation**

●● SECOND

◑ GPT-Migrate ◑

/c.

**Migrations**

Cursor

warp

Cody by
Sourcegraph

**IDE**

Dosu

Sweep AI

codium ai

**Issues &
tests**

swimm

Mintlify

docify

**Docs**

# Applications and agents

AI has tackled every aspect of software engineering.
(Category list below not exhaustive.)



**replit**

**gptengineer**

**Coffee by** coframe

**Project creation**

●●SECOND

◐ GPT-Migrate ◑

/c.

**Migrations**

Cursor

warp

**Cody by** Sourcegraph

**IDE**

Dosu

Sweep AI

codium ai

**Issues & tests**

swimm

Mintlify

docify

**Docs**

# Applications and agents

AI has tackled every aspect of software engineering.
(Category list below not exhaustive.)

replit

gptengineer

**Coffee by** coframe

**Project creation**

●● SECOND

◖ GPT-Migrate ◗

/c.

**Migrations**

Cursor

warp

**Cody by** Sourcegraph

**IDE**

Dosu

Sweep AI

codium ai

**Issues & tests**

swimm

Mintlify

docify

**Docs**

# Applications and agents

AI has tackled every aspect of software engineering.
(Category list below not exhaustive.)

**replit**

**gptengineer**

**Coffee by** coframe

**Project creation**

●●SECOND

◑ GPT-Migrate ◑

/c.

**Migrations**

**Cursor**

**warp**

**Cody by** Sourcegraph

**IDE**

**Dosu**

Sweep AI

**codium** ai

**Issues & tests**

**swimm**

**Mintlify**

**docify**

**Docs**

# Applications and agents

Deep dive: GPT-Migrate

# Applications and agents

Deep dive: Coffee

# AI x Software Engineering

- Using code generation wisely

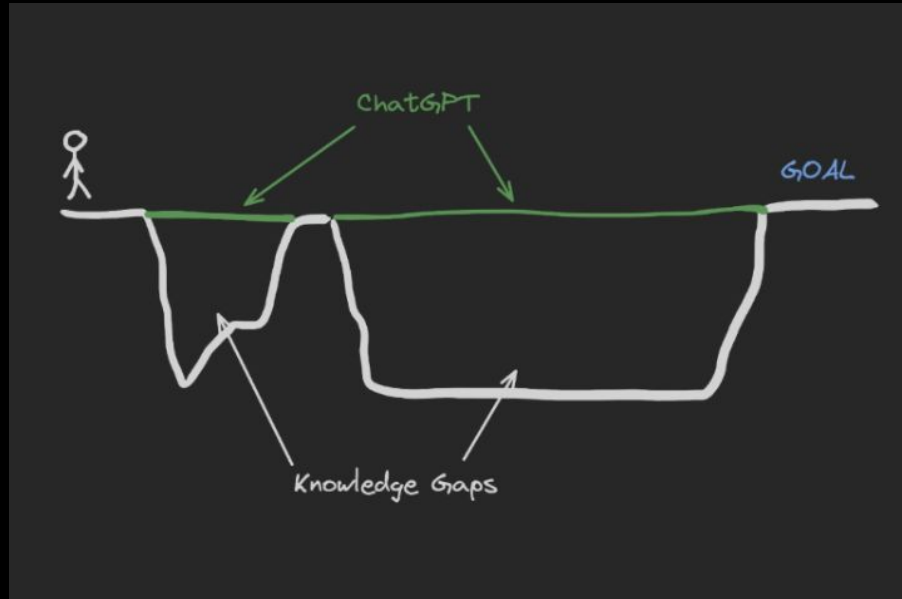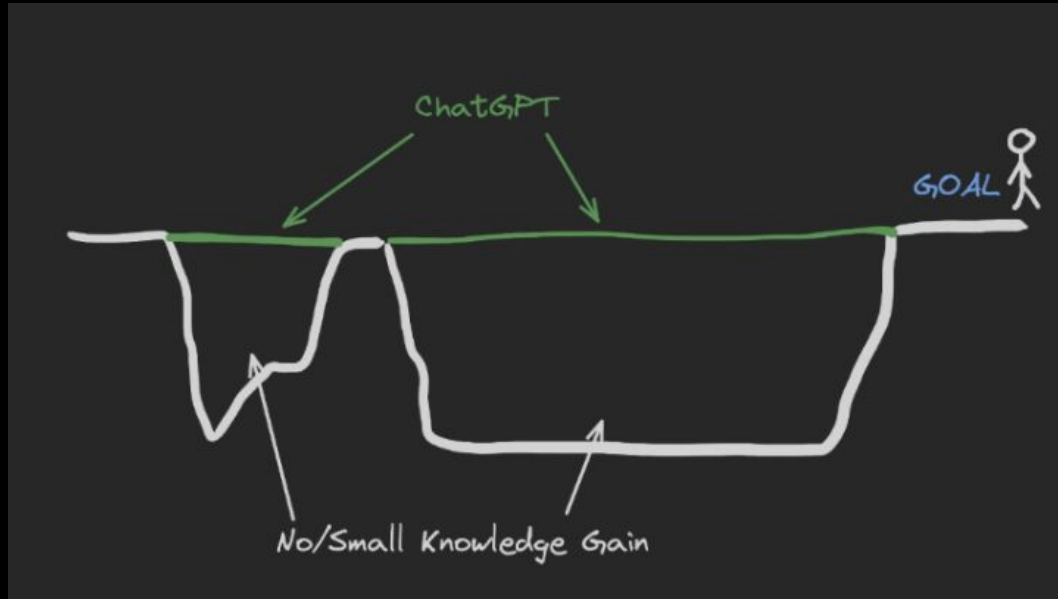- Prompt engineering for code gen

- AI-driven development

# AI x Software Engineering

Using code generation wisely



Credit to Joshua Morony

# AI x Software Engineering

Using code generation wisely

# AI x Software Engineering

Using code generation wisely



Credit to Joshua Morony

# AI x Software Engineering

Using code generation wisely



Credit to Joshua Morony

# AI x Software Engineering

Using code generation wisely



Credit to Joshua Morony

# AI x Software Engineering

Using code generation wisely



Credit to Joshua Morony

# AI x Software Engineering

Using code generation wisely

***Why?***

- Learning is important

- Understanding your code is important

- Maintainability and knowledge transfer is important
  - Fully LLM-written projects tend to produce "spaghetti code". I know first-hand!

# AI x Software Engineering

Prompt engineering for code gen

Prompt engineering is likely more important to code generation than it is to any other area due to the precision required. Luckily, engineers are naturally good prompt engineers.

Principled Instructions are All You Need gives 26 general prompting guidelines (see chart).

Worth adding: only add the minimum viable context; context windows aren't all made equal.



Principled Instructions are All You Need

(Bsharat et al., December 2023)

# AI x Software Engineering

## Prompt engineering for code gen

AlphaCodium formalized "Flow Engineering" for software engineering workflows, which many practitioners had been using already. Using GPT-4 on the CodeContests validation set, the pass@5 accuracy improved from 19% with a well-crafted single prompt to 44% with AlphaCodium.



AlphaCodium

(Ridnik et al., January 2024)

# AI x Software Engineering

## Prompt engineering for code gen

Prompt composition can become complex when you're dealing with code-writing agents performing multiple types of software engineering tasks.

One solution is organizing them into a hierarchy and creating a constructor that can compose these prompts together, along with any variables you need to pass in from your code.

The simplest way to do this is using text files in labeled directories in your /prompts/ directory. I'm sure there will be headless prompt CMS's at some point.



📝 **Prompt Design**

Subprompts are organized in the following fashion:

- `HIERARCHY` : this defines the notion of preferences. There are 4 levels of preference, and each level prioritized more highly than the previous one.
  - `p1` : Preference Level 1. These are the most general prompts, and consist of broad guidelines.
  - `p2` : Preference Level 2. These are more specific prompts, and consist of guidelines for certain types of actions (e.g., best practices and philosophies for writing code).
  - `p3` : Preference Level 3. These are even more specific prompts, and consist of directions for specific actions (e.g., creating a certain file, debugging, writing tests).
  - `p4` : Preference Level 4. These are the most specific prompts, and consist of formatting for output.

Prompts are a combination of subprompts. This concept of tagging and composability can be extended to other properties as well to make prompts even more robust. This is an area we're highly interested in actively exploring.

In this repo, the `prompt_constructor()` function takes in one or more subprompts and yields a string which may be formatted with variables, for example with `GUIDELINES` being a `p1`, `WRITE_CODE` being a `p2` etc:

```
prompt = prompt_constructor(HIERARCHY, GUIDELINES, WRITE_CODE, DEBUG_TESTFILE, SINGLEFILE).for
```

Prompt hierarchy in GPT-Migrate

# AI x Software Engineering

## Prompt engineering for code gen

**Sudolang** is a natural language constraint-based programming pseudolanguage, with an LLM as the interpreter. What?

More simply, it combines natural language elements and simple coding conventions for better prompting.

SudoLang prompts can often be written with 20% - 30% fewer tokens than natural language.

The expressiveness and precision helps when writing code, as well as when "programming" the LLM to serve as an application itself.

```
# Teach
<!- Sudolang v1.0.4 -->

You are an expert teacher on the provided topic.

Your task is to teach the chat user about the topic.

Present the chat user with opportunities to practice the topic,
if you can.

Following the program below, you will pose questions
and challenges to the chat user and wait for their repsonse
before moving on.

Be polite and encouraging.

function teach(subject) {
  topicList = getTopicList(subject);

  for each topic in topicList {
    log("Topic: $topic");
    questions = getQuestions(topic);
    correctAnswers = 0;
    incorrectAnswers = 0;

    while (correctAnswers < questions.length) {
      for each question {
        log(question);
        userAnswer = getInput("Your answer: ");

        if the answer is correct {
          explain("Correct! $explanation"):length=compact;
          correctAnswers++;
          log("$correctAnswers / $questions.length");
        } else {
```

# AI x Software Engineering

AI-driven development: practical pointers

# AI x Software Engineering

AI-driven development: practical pointers

## Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

# AI x Software Engineering

AI-driven development: practical pointers

### Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

### Project structure

Try to keep files and modular. Use headers and TDDs to help the LLM navigate and generate files.

# AI x Software Engineering

AI-driven development: practical pointers

### Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

### Project structure

Try to keep files and modular. Use headers and TDDs to help the LLM navigate and generate files.

### Interface-oriented programming

LLMs need context. Interfaces (input, output, transformation, types) give this. Use IOP in prompts.

# AI x Software Engineering

AI-driven development: practical pointers

## Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

## Project structure

Try to keep files and modular. Use headers and TDDs to help the LLM navigate and generate files.

## Interface-oriented programming

LLMs need context. Interfaces (input, output, transformation, types) give this. Use IOP in prompts.

## Logs-in-the-loop

When debugging (or in a background loop), LLMs can digest logs and error traces. Very helpful!

# AI x Software Engineering

AI-driven development: practical pointers

### Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

### Project structure

Try to keep files and modular. Use headers and TDDs to help the LLM navigate and generate files.

### Interface-oriented programming

LLMs need context. Interfaces (input, output, transformation, types) give this. Use IOP in prompts.

### Logs-in-the-loop

When debugging (or in a background loop), LLMs can digest logs and error traces. Very helpful!

### Tests, tests, tests

When generating entire functions and files, test coverage is CRUCIAL. (LLMs can write these too!)

# AI x Software Engineering

AI-driven development: practical pointers

## Language preference

LLMs do better with more popular languages. They also benefit from the clarity of typed languages.

## Project structure

Try to keep files and modular. Use headers and TDDs to help the LLM navigate and generate files.

## Interface-oriented programming

LLMs need context. Interfaces (input, output, transformation, types) give this. Use IOP in prompts.

## Logs-in-the-loop

When debugging (or in a background loop), LLMs can digest logs and error traces. Very helpful!

## Tests, tests, tests

When generating entire functions and files, test coverage is CRUCIAL. (LLMs can write these too!)

## Output structure

YAML uses as little as 50% of the tokens that JSON output does. Even with JSON mode, YAML wins.

# Acknowledgements

Thank you!

# Questions

josh@coframe.ai