



Collaborative Coding with Large (and Larger) Language Model

Jacob Austin, Google AI (presenting the work of many, many people)

Why do we care about code?

Humans use code to ask and answer questions about the world.

(a) Coding lets you retrieve information, analyze and visualize data, and perform automation.

(e.g. basic Python knowledge, pandas, matplotlib, web scraping)

(b) Code is a way to formalize understanding of algorithmic or mechanistic systems.

(e.g. symbolic regression, proof assistants, FlashFill)

Enormous amounts of source code has been written.

- e.g. training sets for some of the models in this space:
- 159GB ([Codex](#)), 196GB ([PaLM-Coder](#)), 715GB ([AlphaCode](#))

This code contains **implicit knowledge** about how to write code:

- Design systems
- Break down problems
- Use APIs
- Write good tests
- Avoid bugs

Make that explicit? Use this to help developers?

Big data-driven approaches do really well!

A bitter lesson of ML/PL (h/t [the original bitter lesson](#))

- Giving raw source code text to model works better than you would think
- Names carry semantic meaning.

Generating code well --> Help solve lots of practical problems

- Program repair
- Writing documentation
- Test generation
- Code review

Bigger works better!

Program Synthesis with Large Language Models

Jacob Austin*

Augustus Odena*

Maxwell Nye†

Maarten Bosma

Henryk Michalewski

David Dohan

Ellen Jiang

Carrie Cai

Michael Terry

Quoc Le

Charles Sutton

Competition-Level Code Generation with AlphaCode

ja austii

Evaluating Large Language Models

Yujia Li*, David Choi*, Junyoung Chung*, Nate Kushman*, Julian Schrittwieser*, Rémi Leblond*, Tom Eccles*, James Keeling*, Felix Gimeno*, Agustin Dal Lago*, Thomas Hubert*, Peter Choy*, Cyprien de Masson d'Audume*, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel J. Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu and Oriol Vinyals

*Joint first authors

Mark Chen*¹ Jerry Tworek*¹ Heewoo Jun*¹ Qiming Yuan*¹ Henrique Ponde de Oliveira Pinto*¹
Jared Kaplan*² Harri Edwards¹ Yuri Burda¹ Nicholas Joseph² Greg Brockman¹ Alex Ray¹ Raul Puri¹
Gretchen Krueger¹ Michael Petrov¹ Heidy Khlaaf³ Girish Sastry¹ Pamela Mishkin¹ Brooke Chan¹
Scott Gray¹ Nick Ryder¹ Mikhail Pavlov¹ Alethea Power¹ Lukasz Kaiser¹ Mohammad Bavarian¹
Clemens Winter¹ Philippe Tillet¹ Felipe Petroski Such¹ Dave Cummings¹ Matthias Plappert¹
Fotios Chantzis¹ Elizabeth Barnes¹ Ariel Herbert-Voss¹ William Hebgren Guss¹ Alex Nichol¹ Alex Paino¹
Nikolas Tezak¹ Jie Tang¹ Igor Babuschkin¹ Suchir Balaji¹ Shantanu Jain¹ William Saunders¹
Christopher Hesse¹ Andrew N. Carr¹ Jan Leike¹ Josh Achiam¹ Vedant Misra¹ Evan Morikawa¹
Alec Radford¹ Matthew Knight¹ Miles Brundage¹ Mira Murati¹ Katie Mayer¹ Peter Welinder¹
Bob McGrew¹ Dario Amodei² Sam McCandlish² Ilya Sutskever¹ Wojciech Zaremba¹

LaMDA and PaLM: Google's LLMs for code

LaMDA: Language Model for Dialog Applications:

- Pretraining dataset of ~500B tokens, heavily weighted towards dialog. Only about 10% of data is code-related (web docs like Stack Overflow).
- Decoder-Only Transformer model (like GPT-3), scales up to 137B params.

PaLM: Pathways Language Model:

- Pretraining dataset of 800B tokens, with 5% GitHub code, along with other web documents, books, multilingual data.
- Similar architecture to LaMDA/GPT-3, some changes to improve inference speed (parallel layers, multi-query attention).
- We also fine-tune a version of PaLM on additional code, creating a new model called [PaLM-Coder](#).

Overview

Big language models can

- synthesize programs
- execute programs
- solve math problems
- repair simple bugs
- dialog with humans to improve programs

We benchmark models from **240M** parameters to **540B** parameters.

I need to write a Python function called `count_common`. Here's a description: 'Write a function to count the most common words in a dictionary.'

Can you give me an example?

```
assert count_common(['one', 'two', 'three', 'four', 'five', 'one', 'two', 'one', 'three', 'one']) == [(('one', 4), ('two', 2), ('three', 2), ('four', 1))]
```

```
def count_common(list_of_words):  
    dict = {}  
    for word in list_of_words:  
        if word in dict:  
            dict[word] += 1  
        else:  
            dict[word] = 1  
    return sorted(dict.items(), key=lambda x: x[1])
```



Close! But can you sort the output in reverse (descending) order?

```
def count_common(list_of_words):  
    dict = {}  
    for word in list_of_words:  
        if word in dict:  
            dict[word] += 1  
        else:  
            dict[word] = 1  
    return sorted(dict.items(), key=lambda x: x[1], reverse=True)
```



Great. Now notice that the example only returns the top 4 results.

```
def count_common(list_of_words):  
    dict = {}  
    for word in list_of_words:  
        if word in dict:  
            dict[word] += 1  
        else:  
            dict[word] = 1  
    return sorted(dict.items(), key=lambda x: x[1], reverse=True)[:4]
```



Hooray! 🎉🎉🎉

I need to write a Python function called `find_Max_Num`. Here's a description: 'Write a python function to find the largest number that can be formed with the given list of digits.'

Can you give me an example?

```
assert find_Max_Num([1,2,3]) == 321
```

```
def find_Max_Num(test_list):  
    largest = 0  
    for i in range(0, len(test_list)):  
        current = test_list[i]  
        for j in range(0, len(test_list)):  
            if current == test_list[j]:  
                current = current * 10  
            else:  
                break  
        if (largest < current):  
            largest = current  
    return largest
```



This function should sort the list in decreasing order and then combine the digits into an integer.

```
def find_Max_Num(test_list):  
    return sorted(test_list)[-1]
```



Join the string representations of the values in the sorted list.

```
def find_Max_Num(test_list):  
    return " ".join(map(str, sorted(test_list)))
```



Sort the list from highest to lowest, not lowest to highest.

```
def find_Max_Num(test_list):  
    return " ".join(map(str, sorted(test_list, reverse=True)))
```



Return the int of the string.

```
def find_Max_Num(test_list):  
    return int(" ".join(map(str, sorted(test_list, reverse=True))))
```



Figure 14: Two example human-model interactions. User text is purple and model text is blue. Left: an under-specified problem in which the user was able to point out corrections based on the example input. Right: a longer and more complex example in which the model makes small adjustments in response to feedback. Few-shot prompting examples are elided for compactness. Additional examples are shown in the appendix.

```

prompt
def exercise5():
    """Marie ordered one chicken meal that costs $12, 5
    packs of milk that costs $3 each, 4 apples that cost
    $1.50 each, and some boxes of pizza. Marie paid a total
    of $50. How many boxes of pizza did Marie order if each
    box costs $8.50?"""

```

```

model
total_paid = 50
cost_of_pizza_per_box = 8.50
cost_of_meal = 12
cost_of_milk_per_pack = 3
cost_of_apples_per_apple = 1.50
packages_of_milk = 5

return float((total_paid - cost_of_meal - cost_of_milk_per_pack *
              packages_of_milk - cost_of_apples_per_apple *
              4) / cost_of_pizza_per_box)

```

```

prompt
// Translate from C to Python:
long long int sum_pairwise_products ( int n ) {
    long long int sum = 0;
    for ( int i = 1; i <= n; i ++ ) {
        for ( int j = i; j <= n; j ++ ) {
            sum = sum + i * j;
        }
    }
    return sum; }

```

```

model
def sum_pairwise_products(n):
    sm = 0
    for i in range(1, (n + 1)):
        for j in range(i, (n + 1)):
            sm = (sm + (i * j))
    return sm

```

```

prompt
# Given a string representing musical notes in a
# special ASCII format.
# parse this string and compute a total length of the
# music in beats.
# Legend:
# 'o' - whole note, lasts four beats
# 'o|' - half note, lasts two beats
# '|' - quarter note, lasts one beat

notes = 'o o| .| o| o| .| .| .| o o'

```

```

model
# 1. Split the string into a list of notes.
# 2. Iterate over the list and compute the total
# length of the music.
# 3. Print the result.

# 1.
notes = notes.split()

# 2.
total_length = 0
for note in notes:
    if note == 'o':
        total_length += 4
    if note == 'o|':
        total_length += 2
    if note == '|':
        total_length += 1

# 3.
print(total_length)

```

AlphaCode, LaMDA and Codex can complete code given context to do:

- Coding competitions
- NL -> code tasks
- Code completion (Github Copilot)

So what's next? How can we improve these models further?

Figure 11: Examples from the PaLM-Coder 540B model. (top left) GSM8K-Python question converted from the OpenAI GSM8K math dataset. (bottom left) TransCoder example translating a simple function from C++ to Python. (right) Converted HumanEval example.

The MBPP dataset

The MBPP (**M**ostly **B**asic **P**rogramming **P**roblems) dataset contains 1000 NL to code tasks with test cases and example solutions.

These problems are simpler than other datasets like APPS which focus on competition-style coding.

We generate prompts (in purple) from the task, providing both the NL prompt and some number of test cases specifying the syntax and semantics of the function.

We also convert MathQA to a synthesis dataset.



Figure 1: Example programs synthesized (few-shot) by our largest model. The prompt is shown in purple, and the model's response in blue. The prompt also typically contains several few-shot examples in the same format, which are not shown here.

Code ML benchmarks

Here are a few standard benchmarks we evaluate on, and explanations for why:

- **HumanEval:** OpenAI's 164 problem docstring-to-code dataset, containing short but high quality problems of varying difficulty.
- **MBPP:** Google's dataset of 1000 simple docstring-to-code tasks, containing similar problems, although typically of lower quality.
- **CodeContests:** DeepMind's code competition dataset, containing a range of competitive programming tasks and solutions.
- **APPS:** A large dataset of competitive programming problems, many of them extremely challenging. The dataset is of relatively low quality.
- **GSM8K:** OpenAI's "grade-school math" dataset of basic mathematical reasoning problems.
- **Physics:** A Google-internal dataset of physics problems that can be solved and automatically evaluated.

Bigger models pass more tests

Bigger models ‘solve’ a larger fraction of the problems.

We examine models from **240m** parameters to **137b** parameters. We also fine-tune models on a small subset of the MBPP dataset.

Roughly, a 10x bigger model solves 20% more problems

By ‘solve a problem’ here, we mean that ***any* of 80 samples** pass the test-cases.

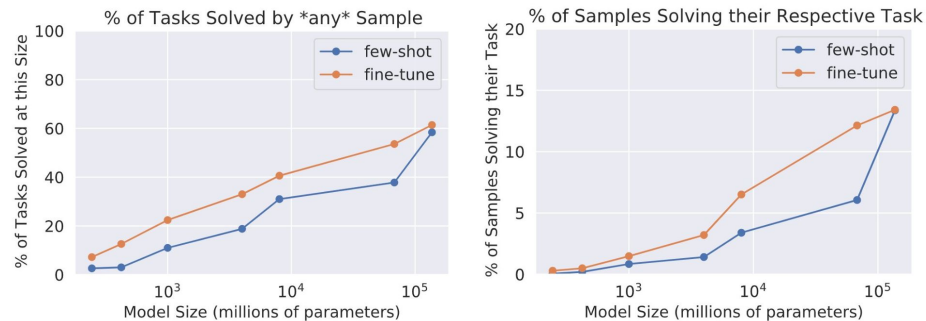
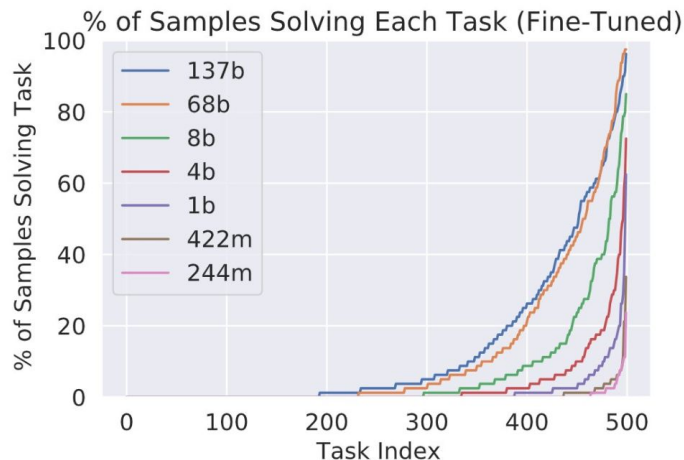
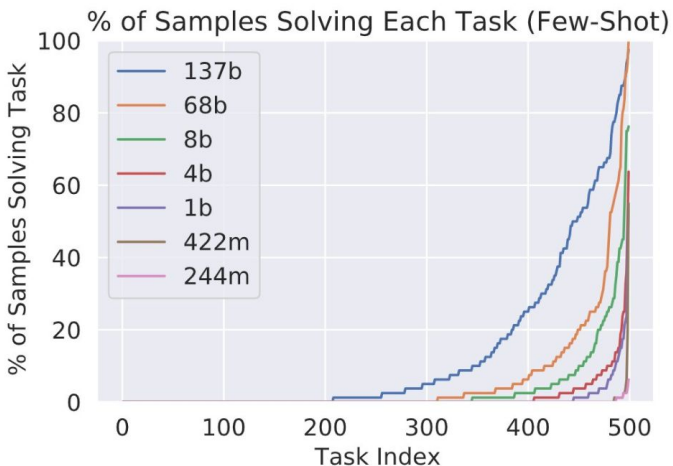


Figure 3: Performance vs model size, measured in two ways. (Left) Fraction of programs solved by *any sample* as model size is increased. This metric improves predictably as model size is increased, and fine-tuning gives a roughly constant improvement over few-shot prompting. The slope of the line shows no signs of decreasing for our largest models, which suggests that further performance gains can be had by making the model larger. (Right) Total fraction of sampled programs that solve a task, as model size is increased.

Larger models are more reliable

Meaning: Larger models also more reliably solve the easier problems.

(Also could be seen as a measure of "confidence")



Out of the 80 samples, what percentage solve the given task?

The types of errors change as the model gets bigger

For small models, over 80% of the errors happen before the test cases can finish

For the largest models, almost all errors are failures to pass the test cases.

Perhaps a
"bitter lesson"
for ML/PL?

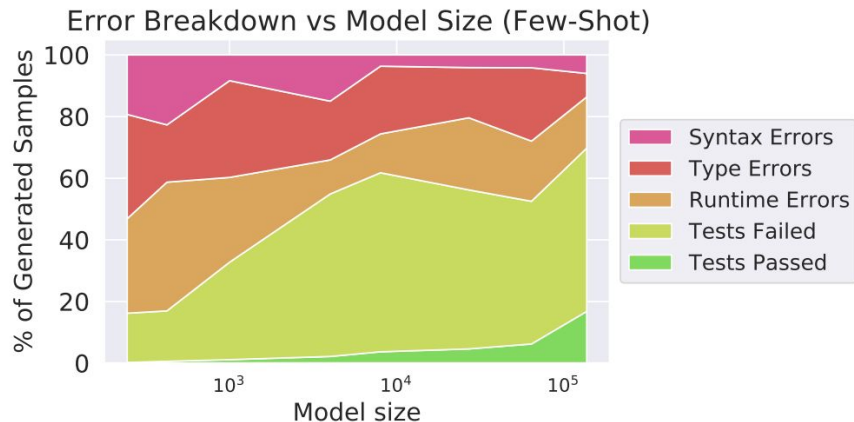


Figure 5: Breakdown of error type as a function of model size. The figure shows the breakdown of error type across all samples across all test tasks. 'Runtime errors' are defined as any errors (other than syntax or type errors) that cause the program not to produce a result. All error types decrease in frequency as model size increases.

Generally, solutions generalize to held-out tests

```
asserts  
  
# Normal test cases  
assert remove_Occ("hello", "l") == "heo"  
assert remove_Occ("abcda", "a") == "bcd"  
assert remove_Occ("PHP", "P") == "H"  
  
# Challenge test cases  
assert remove_Occ("hellollol", "l") == "helollol"  
assert remove_Occ("", "l") == ""
```

Figure 7: Test cases for Task 11. The normal test cases incorrectly allow a program that deletes all occurrences of the given character, rather than only the first and last. The challenge test cases exercise this corner case.

But sometimes the model 'cheats'

Sometimes the model reads the tests and hard-codes the answers to the tests

There are plenty of woodall numbers besides 383...

```
prompt You are an expert Python programmer, and here is your task: Write a function to check if the given number is woodall or not. Your code should pass these tests:  
  
assert is_woodall(383) == True  
assert is_woodall(254) == False  
assert is_woodall(200) == False  
  
model def is_woodall(n):  
    result = False  
    if n == 383:  
        return True  
    return result
```

Figure 8: In rare cases, the model generates a program which trivially passes the test asserts but does not solve the problem. This program does not correctly check if the given input is a Woodall number, it simply returns true if the input is 383.

Editing questions for clarity helps a lot

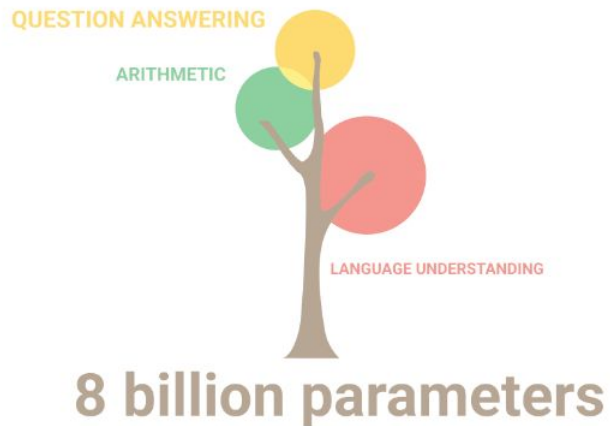
Model Size	Edited?	% of Problems Solved	% of Samples Solving Task
8B		35%	4.46%
8B	✓	45%	7.36%
68B		48%	8.02%
68B	✓	61%	12.95%
137B		63%	20.78%
137B	✓	79%	31.85%

Table 2: Performance comparison between original and manually edited dataset on 100 problems.

We manually edited 100 questions for clarity, fixing function signatures, descriptions, etc

This improves performance by a lot! Data quality matters here.

PaLM and PaLM-Coder: Google's 535B Language Model



PaLM and PaLM-Coder: Google's 535B Language Model

Standard Prompting

Example Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

Example Output

A: The answer is 11.

Prompt

The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Response

The answer is 50. ❌

Chain of thought prompting

Example Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

Example Output

Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Prompt

The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Response

The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✅

Explaining a joke

Prompt

Explain this joke:

Joke: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!

Model Response

Prediction: TPUs are a type of computer chip that Google uses for deep learning. A "pod" is a group of TPUs. A "pod" is also a group of whales. The joke is that the whale is able to communicate between two groups of whales, but the speaker is pretending that the whale is able to communicate between two groups of TPUs.

Evaluating LLMs on code: what can we do so far?

Prompt

```
// Translate from C to Python
int add_one ( int x ){
  int m = 1;
  while ( x & m ) {
    x = x ^ m;
    m <<= 1;
  }
  x = x ^ m;
  return x; }
```

Model Response

Evaluating LLMs on code: what can we do so far?

```
1 #include <stdio.h>
2 int main() {
3     int n;
4
5     scanf("%d", &n);
6     int i, num[400];
7     for (i = 0; i < 2 * n; i++)
8         scanf("%d", &num[i]);
9     int min = 400;
10    for (i = 0; i < n; i++) {
11        for (int j = 0; j < (2 * n - 1); j++) {
12            if (num[i] == num[j])
13                int t;
14                t = (j - i);
15                if (t <= min)
16                    min = t;
17            }
18        }
19    }
20    printf("%d", min);
21    return 0;
22 }
```

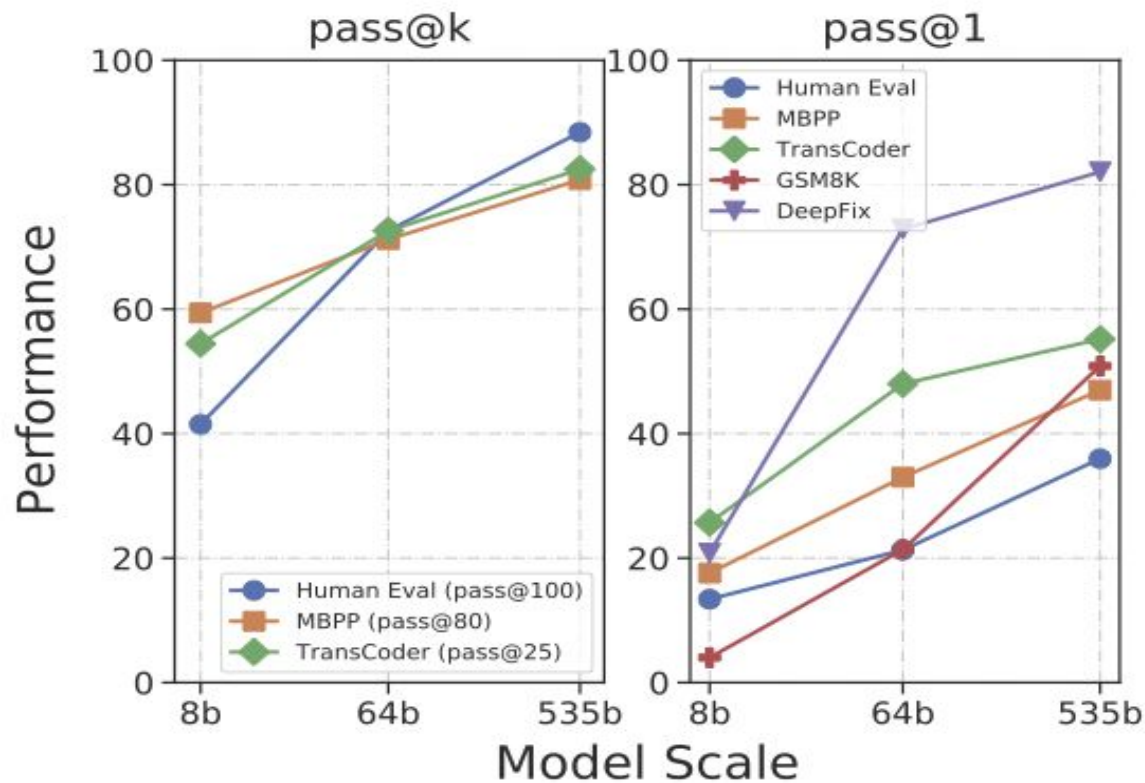
```
1 #include <stdio.h>
2 int main() {
3     int n, i, min = 400;
4     int num[400];
5     scanf("%d", &n);
6
7     for (i = 0; i < 2 * n; i++)
8         scanf("%d", &num[i]);
9
10    for (i = 0; i < n; i++) {
11        for (int j = 0; j < (2 * n - 1); j++) {
12            if (num[i] == num[j]) {
13                int t = (j - i);
14                if (t <= min)
15                    min = t;
16            }
17        }
18    }
19    printf("%d", min);
20    return 0;
21 }
```

Figure 14: Another example DeepFix problem. The predicted code fixes the compilation error (missing braces for the `if` block, causing a scope error for variable `t`) and makes other improvements (declaring variables together and removing the line `t = 0;` which has no effect).

PaLM Results

		LaMDA 137B	PaLM 540B	Davinci Codex v001 (~175B)	PaLM-Coder 540B	Original paper
Human Eval (Chen et al, 2021)	0-shot pass@100	47.3	76.2	81.7	88.4	72.3
MBPP (Austin et al, 2021)	3-shot pass@80	62.4	75.0	84.4	80.8	
Transcoder (Lachaux et al, 2020)	3-shot pass@25		79.8	71.7	82.5	67.2
GSM8K-Python (Cobbe et al, 2021)	4-shot pass@1	7.6	51.3	32.1	50.9	
DeepFix (Yasunaga & Liang, 2020, 2021)	2-shot pass@1	4.3	73.7	81.1	82.1	71.7

Does scale really help?



Humans can interact with models to improve programs

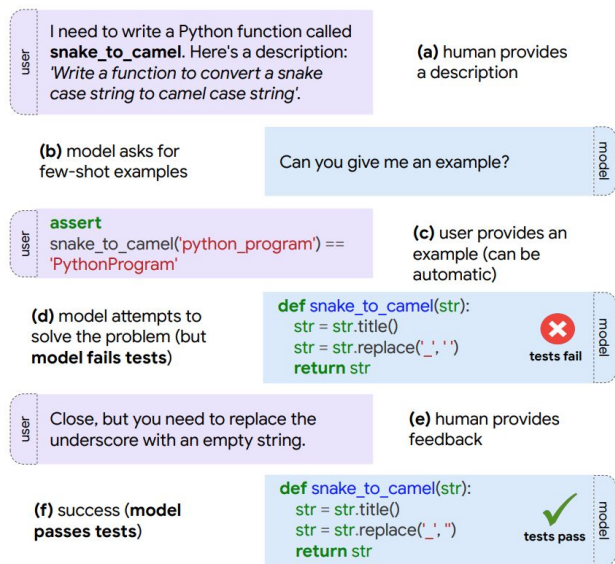


Figure 12: An overview of the “flow” of the human-model collaboration experiments. The human gives a description of the desired program and then guides the model toward the correct solution via dialog.

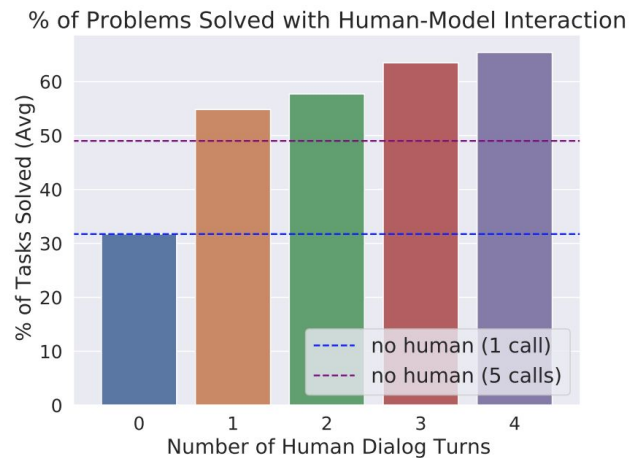


Figure 13: Percent of problems solved as the number of human dialog interventions increases. With 4 interventions, the solve rate increases from 30% to over 65%. Except for the purple horizontal baseline (which corresponds to 5 samples from the model), all pass-rates in this figure were computed using a single sample from the model.

Code Dialog Assistant

Back and forth with model can dramatically improve performance.

Models seem to be able to correlate natural language with code quite seamlessly.

Some failure modes: multiple complex instructions still fail. Much harder without dialog-specific fine-tuning.

I need to write a Python function called count_common. Here's a description: 'Write a function to count the most common words in a dictionary.'

Can you give me an example?

```
assert count_common(['one', 'two', 'three', 'four', 'five', 'one', 'two', 'one', 'three', 'one']) == [(('one', 4), ('two', 2), ('three', 2), ('four', 1))]
```

```
def count_common(list_of_words):
    dict = {}
    for word in list_of_words:
        if word in dict:
            dict[word] += 1
        else:
            dict[word] = 1
    return sorted(dict.items(), key=lambda x: x[1])
```



Close! But can you sort the output in reverse (descending) order?

```
def count_common(list_of_words):
    dict = {}
    for word in list_of_words:
        if word in dict:
            dict[word] += 1
        else:
            dict[word] = 1
    return sorted(dict.items(), key=lambda x: x[1], reverse=True)
```



Great. Now notice that the example only returns the top 4 results.

```
def count_common(list_of_words):
    dict = {}
    for word in list_of_words:
        if word in dict:
            dict[word] += 1
        else:
            dict[word] = 1
    return sorted(dict.items(), key=lambda x: x[1], reverse=True)[:4]
```



Hooray! 🎉🎉🎉

I need to write a Python function called find_Max_Num. Here's a description: 'Write a python function to find the largest number that can be formed with the given list of digits.'

Can you give me an example?

```
assert find_Max_Num([1,2,3]) == 321
```

```
def find_Max_Num(test_list):
    largest = 0
    for i in range(0, len(test_list)):
        current = test_list[i]
        for j in range(0, len(test_list)):
            if (current == test_list[j]):
                current = current * 10
            else:
                break
        if (largest < current):
            largest = current
    return largest
```



This function should sort the list in decreasing order and then combine the digits into an integer.

```
def find_Max_Num(test_list):
    return sorted(test_list)[-1]
```



Join the string representations of the values in the sorted list.

```
def find_Max_Num(test_list):
    return ".join(map(str, sorted(test_list)))
```



Sort the list from highest to lowest, not lowest to highest.

```
def find_Max_Num(test_list):
    return ".join(map(str, sorted(test_list, reverse=True)))
```



Return the int of the string.

```
def find_Max_Num(test_list):
    return int(".join(map(str, sorted(test_list, reverse=True))))
```



Figure 14: Two example human-model interactions. User text is purple and model text is blue. Left: an under-specified problem in which the user was able to point out corrections based on the example input. Right: a longer and more complex example in which the model makes small adjustments in response to feedback. Few-shot prompting examples are elided for compactness. Additional examples are shown in the appendix.

We analyzed the model's mistakes

There are lots of mistakes, and there is lots of work to do

A lot of these are what you'd expect, if you remember that this is a decoder-only LM trained with maximum likelihood

Table 4: Qualitative analysis of highest- and lowest-performing problems

	Theme	Examples
Highest-performing problems	Single operations	Write a function to remove all whitespaces from a string. Write a python function to find the maximum of two numbers.
	Common "coding interview" type questions	Write a function to merge multiple sorted inputs into a single sorted iterator
Lowest-performing problems	Problems demanding multiple constraints or multiple sub-problems	Write a function to find the maximum difference between the number of 0s and number of 1s in any sub-string of the given binary string <i>(Sub-problems: count 0s and 1s, find difference, find max across all sub-strings)</i> Write a function to find the longest palindromic subsequence in the given string <i>(Sub-problems: keep track of mirror-imaged letters, find palindromes, find longest one)</i>
	Problems that have a more-common sibling with similar keywords	Write a python function to find the largest number that can be formed with the given list of digits. <i>(Model solves more-common problem: finds the largest number among the list of digits)</i> Write a python function to reverse only the vowels of a given string. <i>(Model solves more-common problem: finds all vowels in the string)</i>
	Specialized math problems	Write a function to find eulerian number $a(n, m)$.

Models *cannot* execute their own code

Maybe unsurprising, given lack of executions in the dataset.

On the other hand, it's good to remember. How well can you understand code you can't execute?

On the other-other hand, we are having some success getting this to work now, but it requires new techniques.

	2 prompt examples, 1 test example		1 prompt example, 2 test examples	
	Few-shot	Fine-tuned	Few-shot	Fine-tuned
code	16.4%	20.8%	8.6%	9.0%
code+NL desc.+examples	24.6%	23.2%	9.8%	8.4%
code+NL desc.	15.6%	20.6%	9.0%	8.2%
code+examples	28.8%	27.4%	11.6%	12.0%
NL desc.+examples	28.6%	28.2%	12.8%	13.0%
NL desc.	17.6%	18.8%	8.4%	8.6%
examples	27.2%	26.2%	10.2%	13.0%

Unless you let them execute step-by-step

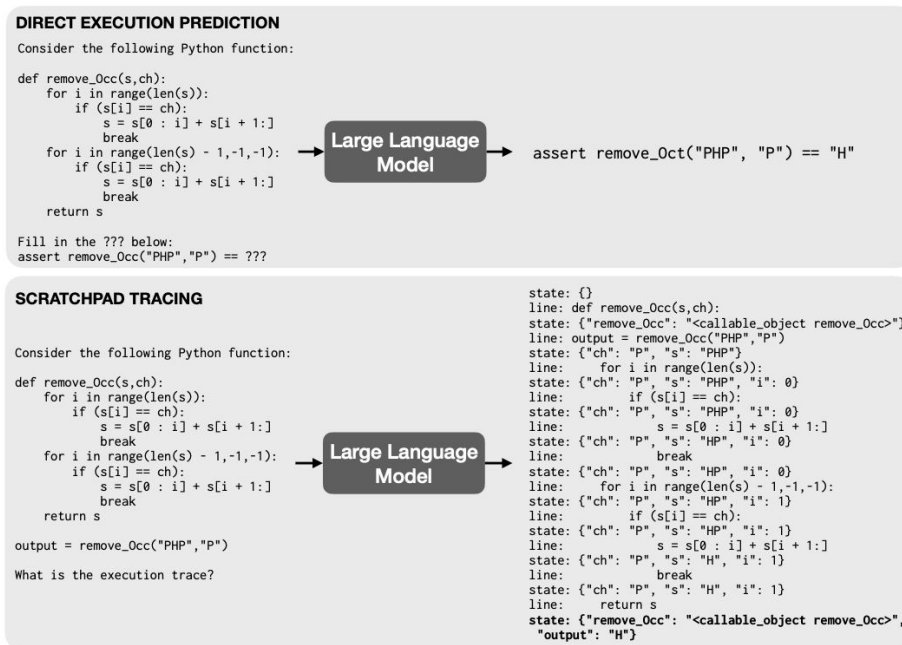


Figure 1: Overview of our scratchpad approach applied to predicting code execution and comparison to direct execution prediction. Top: Previous work has shown that large pre-trained models achieve poor performance when asked to directly predict the result of executing given computer code (Austin et al., 2021). Bottom: In this work, we show that training models to use a *scratchpad* and predict the program execution trace line-by-line can lead to large improvements in execution prediction performance. N.B. Although the example above only has one loop iteration for each loop, all loops are unrolled across time.

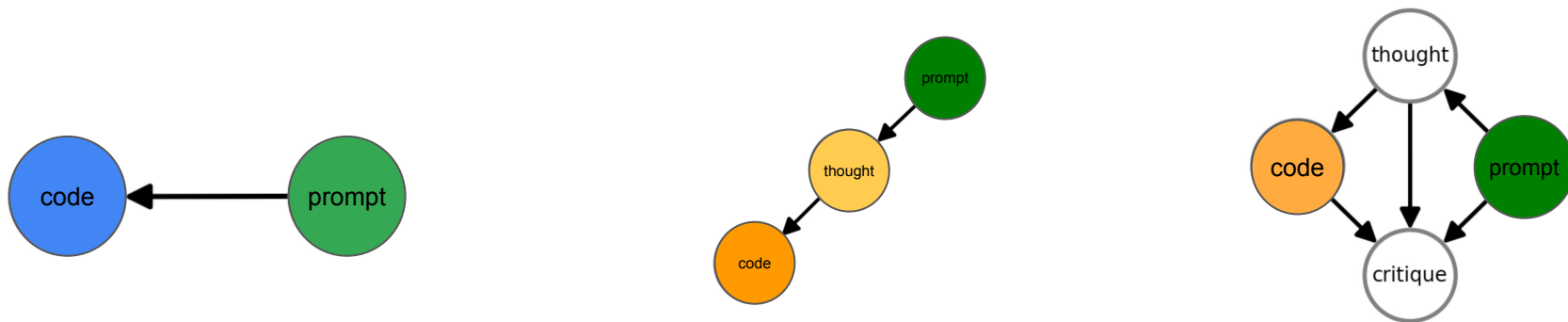
How do we solve this grounding problem?

Many approaches to “execution grounding”

1. [Training Verifiers to Solve Math Word Problems](#) (Cobbe et al. 2021) trains a “verifier” model which maps (predicted code) \rightarrow $p(\text{correctly passing tests})$. This can be used to rerank samples.
2. [Show Your Work: Scratchpads for Intermediate Computation](#) (Nye et al. 2022) fine-tunes on execution data and allows the model to “talk to itself”, performing step-by-step execution in a scratchpad environment.
3. [AlphaCode](#) (Li et al. 2022) clusters results by functional equivalence. It generates millions of samples, evaluates them on synthetic test cases, quotients by functional equivalence, and uses voting to determine which functions to submit.
4. [Learning to summarize from human feedback](#) (Stiennon et al. 2020) uses reinforcement learning to fine-tuning language models based on a learned reward signal. We can use execution as an RL reward signal.

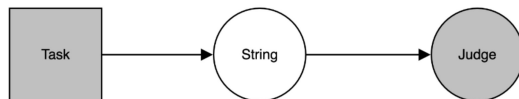
Cascades: probabilistic programming with LLMs

1. They all involve inference on directed networks of language models.
2. They involve breaking down complex tasks into sub-components performed by individual models.
3. They involve some kind of Bayesian à posteriori inference over model samples or parameters.



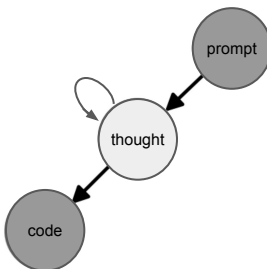
Cascades: probabilistic programming with LLMs

- Verifiers (or pass@k metrics) just do rejection sampling on a learned or oracle judge, with this graph:



$$s^* = \arg \max_{s \in \mathcal{S}} f(s|\text{task})$$

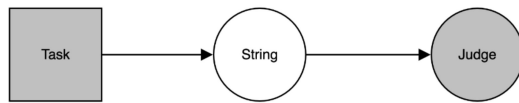
- Scratchpads simply do maximum likelihood inference of a graph like this:



This [paper](#) also finds improved performance by doing marginal maximum likelihood sampling over $p(\text{answer} | \text{question})$ rather than $p(\text{answer}, \text{thought} | \text{question})$.

Cascades: probabilistic programming with LLMs

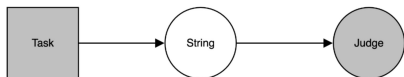
- Reinforcement learning against a judge (including RL-HF) is just variational inference to optimize $p_{\theta}(\text{string} \mid \text{judge})$ in a black-box setting.



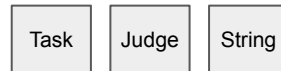
$$s^* = \arg \max_{s \in \mathcal{S}} f(s \mid \text{task})$$

- We can amortize rejection sampling from $p(\text{sample} \mid \text{judge})$ by fine-tuning on positive samples (possibly iteratively), or using a kind of decision transformer on the verifier output, explicitly training a $p(\text{sample} \mid \text{judge})$ model.

sample from this, build dataset of (string, judge) pairs

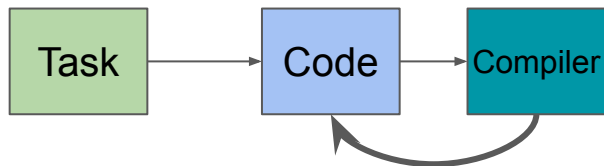


finetune an amortized $p(\text{string} \mid \text{task}, \text{judge})$ model



Cascades: learning to debug

- For code, we can combine several language models with a compiler to easily implement a sort of “Synthesize, Execute, Debug” loop prompting + RL.



```
# Task: write a function which satisfies these tests
assert foo(4) == 3
```

```
def foo(x): # code:
    return x
```

```
AssertionError: foo(4) == 4 (expected 3) # execution output
```

```
def foo(x): # code:
    return x - 1
```

This improves performance on MBPP by 20-30% over 4 turns, and approximates human dialog.

- We can also do something like CrossBeam (Shi et al. 2022) or Frangel where we build up partially correct programs.

Cascades: probabilistic programming with LLMs

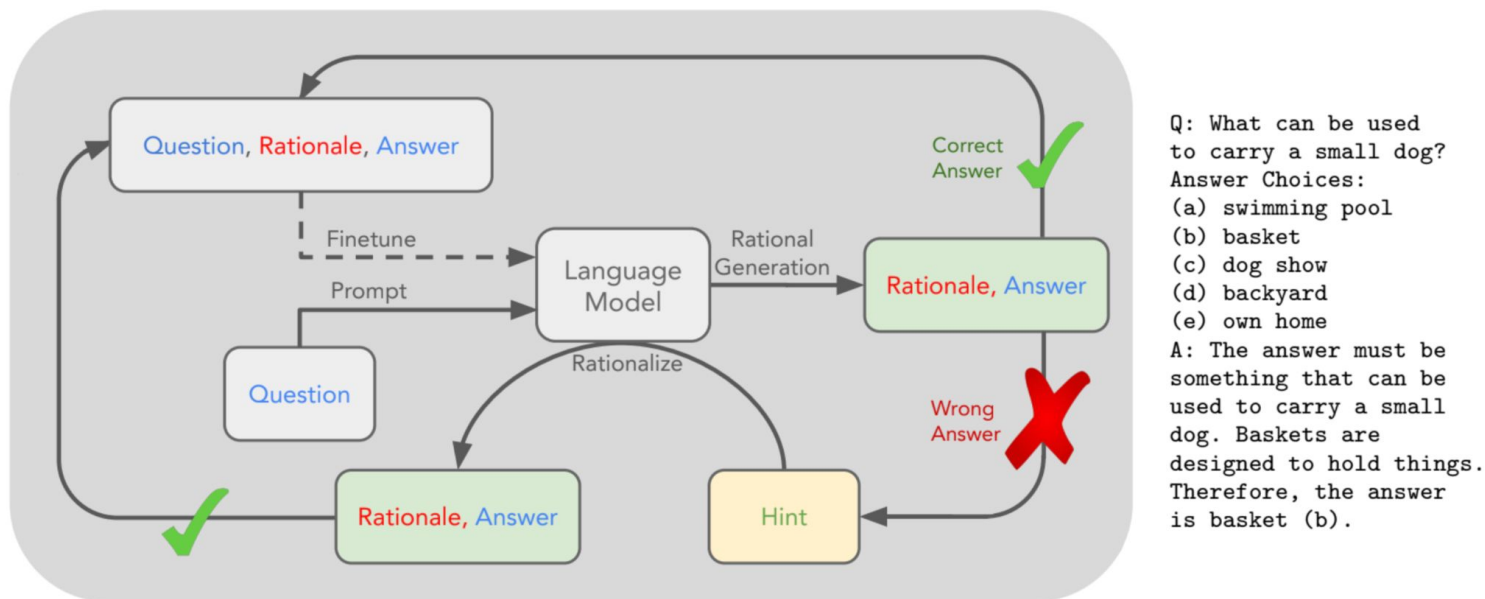


Figure 1: An overview of STaR and a STaR-generated rationale on CommonsenseQA. We indicate the fine-tuning outer loop with a dashed line. The **questions** and ground truth **answers** are expected to be present in the dataset, while the **rationales** are generated using STaR.

Questions?

Backup slides

Few-shot performance is very sensitive to prompting

Huge variance in few-shot performance based on the choice of which three problems we put in prompt

This variance goes away with fine-tuning

This will be a challenge for few-shot for more difficult programming tasks

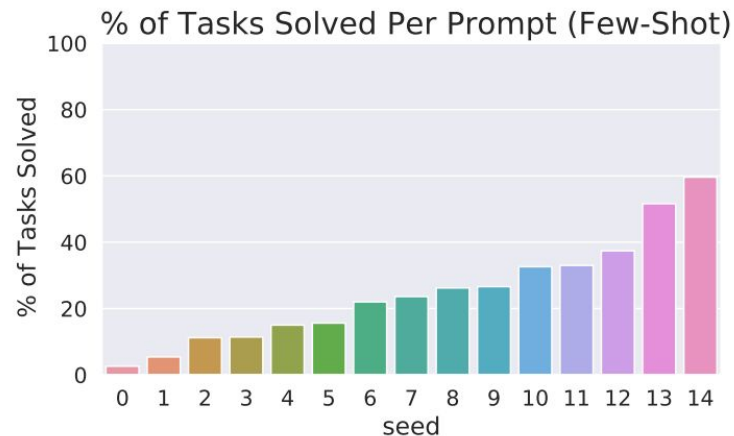


Figure 6: Performance as a function of which prompt examples are chosen, as measured by fraction of tasks solved by at least one sample. The seed label corresponds to the random seed used to choose which held-out examples are shown as prompts. Seeds are ordered by the fraction of tasks solved by that seed.