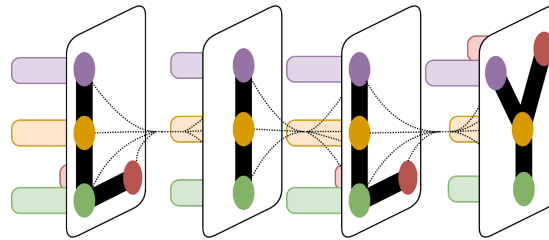# Enhancing Language Models for Program Synthesis using Execution

**Talk @MIT_CSAIL**
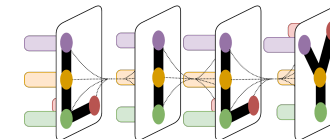
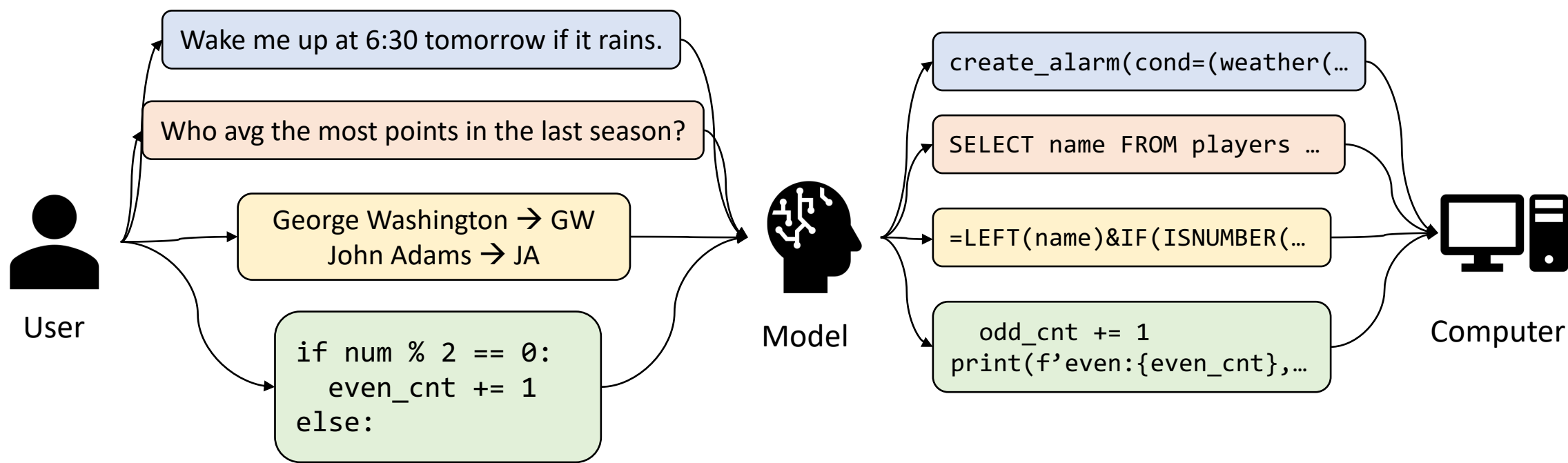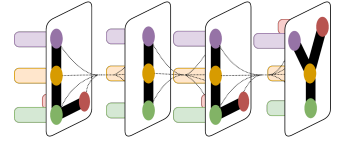Yale LILY Lab

**Ansong Ni**

03-13-2023

Yale

# What is Program Synthesis?

- Here we broadly define **program synthesis** as the tasks to <u>automatically generate programs from the user intent</u>.
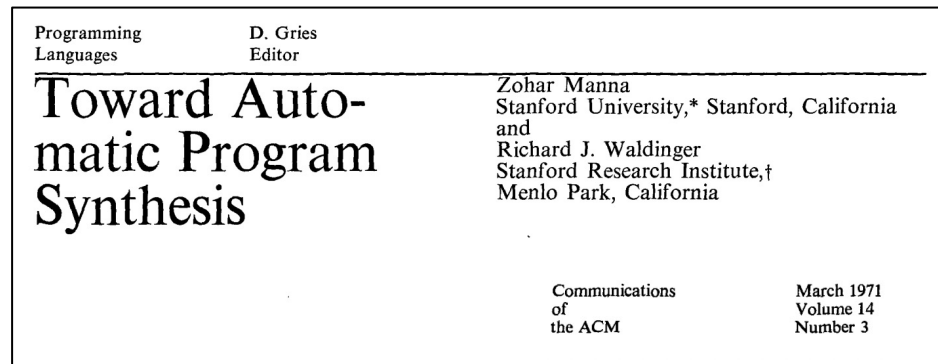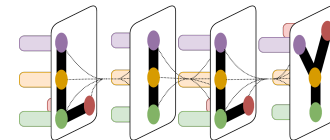
# Why is Program Synthesis Important?

- ## On a higher level:
    - ### It is one of the **oldest** and **hardest** problems in AI and CS:

---

*This process of constructing instruction tables should be very fascinating. There need be no real danger of it ever becoming a drudge, for any processes that are quite mechanical may be turned over to the machine itself.*
        — *Alan Turing (1945)*

---

Programming        D. Gries
Languages          Editor

## Toward Automatic Program Synthesis

Zohar Manna
Stanford University,* Stanford, California
and
Richard J. Waldinger
Stanford Research Institute,†
Menlo Park, California

Communications      March 1971
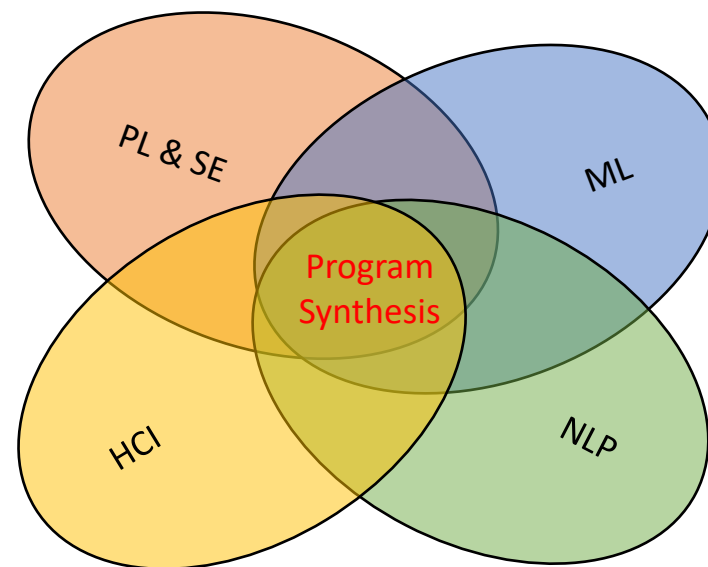of                   Volume 14
the ACM           Number 3

Yale

# Why is Program Synthesis Important?

- ## On a higher level:
  - It is one of the *oldest* and *hardest* problems in AI and CS
  - It involves several important areas in CS
    - Programming Languages (PL)
    - Software Engineering (SE)
    - Machine Learning (ML)
    - Natural Language Processing (NLP)
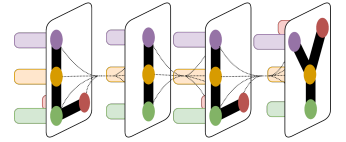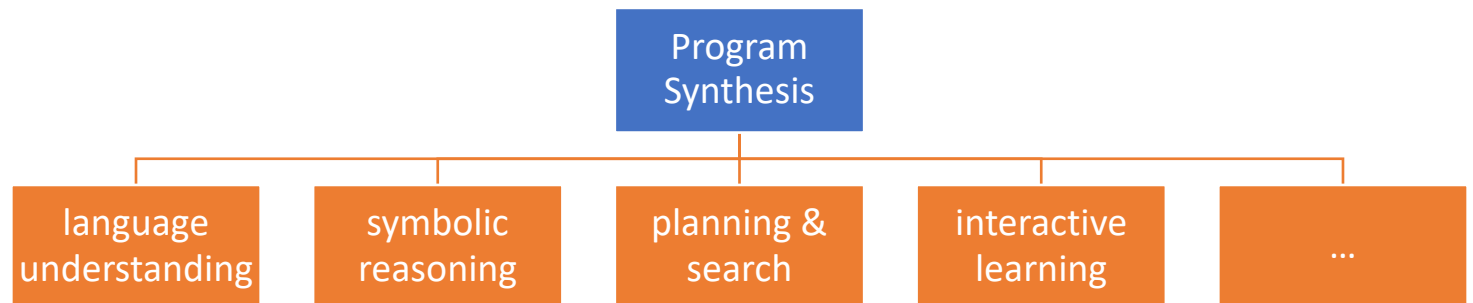    - Human-Computer Interaction (HCI)
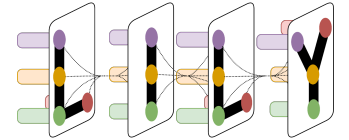    - …

# Why is Program Synthesis Important?

- ## On a higher level:

  - It is one of the *oldest* and *hardest* problems in AI and CS

  - It involves several important areas in CS

  - ## It is a great testbed for intelligence:

    - language understanding

    - symbolic reasoning

    - planning & search

    - interactive learning

    - …

```
                              ┌──────────────┐
                              │   Program    │
                              │   Synthesis  │
                              └──────────────┘
  ┌──────────────┬──────────────┬──────────────┬──────────────┬──────────────┐
  │   language   │   symbolic   │  planning &  │  interactive │      …       │
  │understanding │   reasoning  │    search    │   learning   │              │
  └──────────────┴──────────────┴──────────────┴──────────────┴──────────────┘
```

Yale

# Why is Program Synthesis Important?

- It empowers many real-world applications:

FlashFil

Robotics Control

Virtual Assistants

Database Query and Visualization

# Program Synthesis with Language Models

- What is a language model?
  - Predicts the next word given the context
  - Learns to maximize $P_\theta(x_n|x_1, x_2, \ldots, x_{n-1})$ given training data — self-supervised training

$$P(X) = \prod_{i=1}^{I} P(x_i \mid x_1, \ldots, x_{i-1})$$

Next Word      Context

- *The status quo* of pretrained language models:
  - They are getting **larger**
  - … and **better**
    - Multi-task learning
    - Zero/Few-shot (in-context) learning
    - Instruction tuning
    - …

Images from Graham Neubig and Luke Zettlemoyer

Yale

# Program Synthesis with Language Models

- (Large) language models trained on code:

# Program Synthesis with Language Models

- What does it mean for program synthesis in the *"LLM era"* ?
  - The search space can be greatly reduced/optimized by pretrained CodeLMs
    - → **We can do more!**
    - From domain specific languages (DSL)
      - SQL (Zhong et al., 2017; Yu et al., 2018)
      - $\lambda$-calculus (Zettlemoyer and Collins, 2005; Rabinovich et al., 2017)
      - Karel (Bunel et al., 2018; Chen et al., 2019)
    - To general-purpose programs
      - Basic Python programming (Chen et al., 2021; Austin et al., 2021)
      - Data science code generation (Lai et al., 2022; Yin et al., 2022)
      - Competition-level code generation (Hendrycks et al., 2021; Li et al., 2022)

Yale

# Program Synthesis with Language Models

- What does it mean for program synthesis in the *"LLM era"* ?
  - The search space can be greatly reduced/optimized by pretrained CodeLMs
    → **We can do more!**
    - From supervised learning to zero/few-shot prompting



Example from ChatGPT: https://chat.openai.com/

# Program Synthesis with Language Models

- So are LLMs the solution to program synthesis?
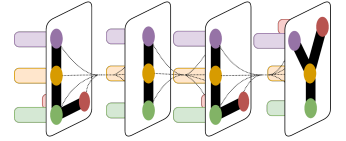    - Examples: programming with ChatGPT



AN  can you write a wrong implementation of `is_prime` function in python

Sure! Here's a wrong implementation of a function to check if a number is prime or not:
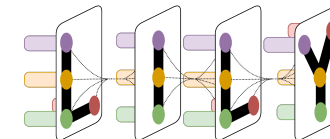
```python
def is_prime(number):
    if number <= 1:
        return False
    for i in range(2, number):
        if (number % i) == 0:
            return False
    return True
```

This implementation is wrong because it only checks if the number is divisible by integers less than itself, but a number can also be divisible by numbers greater than itself. For example, this function would return `True` for `4`, which is not a prime number.

A better implementation would be to check if the number is divisible by any number up to its square root.

AN  can you give me a wrong implementation for determine odd/even numbers

Sure, here is an incorrect implementation to determine if a number is odd or even:

```python
def is_odd_or_even(n):
    if n % 2 == 0:
        return "even"
    else:
        return "odd"
```

This implementation only works for positive integers, but it fails for negative numbers and zero. A correct implementation to handle all cases would be:

```python
def is_odd_or_even(n):
    if n == 0:
        return "even"
    elif n % 2 == 0:
        return "even"
    else:
        return "odd"
```

Example from ChatGPT: https://chat.openai.com/

# Program Synthesis with Language Models

- What is happening here?
  - LLMs are only trained on the **surface form** of programs
  - ...yet the **semantics** of a program is in its **execution**

- But the best of all...
  - Execution can be done automatically!

How to incorporate program semantics into LMs using execution to improve their performance?

Yale

- Task: finetuning pretrained LMs for generating programs from math problems described in natural language

- Motivation - **Program Aliasing**:

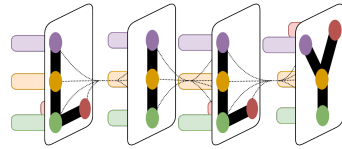  - The same specification (i.e.., question) can be satisfied by different programs
  - However, the training data typically only have one reference solution for learning
  - This causes overfitting as the model keeps seeing the same solution over multiple epochs of training

**NL Input**

A goods train runs at a speed of 72kmph and crosses a 250M long platform in 26 seconds. What is the length of the goods train?

**Ref. Program**

```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n1 / t0
t2 = n2 - t1
answer = t0 * t2
```

**Alt. Program**

```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n2 * t0
answer = t1 - n1
```

- Task: finetuning pretrained LMs for generating programs from math problems described in natural language

- Motivation - *Program Aliasing*

- **Observation**:
  - During inference, the model can sometimes generate programs that are correct but not necessarily the gold one
  - Can we encourage this behavior during training and learning from the self-sampled program solutions?
    - **YES!**

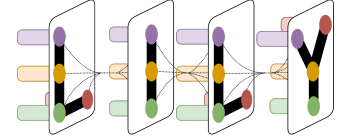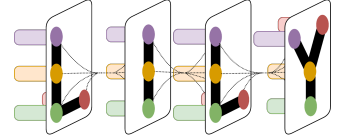**NL Input**

A goods train runs at a speed of 72kmph and crosses a 250M long platform in 26 seconds. What is the length of the goods train?

**Ref. Program**
```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n1 / t0
t2 = n2 - t1
answer = t0 * t2
```

**Alt. Program**
```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n2 * t0
answer = t1 - n1
```

- Use a buffer $\mathcal{B}$ to save self-sampled programs

- Online sampling and filtering
  - Attempt to sample alternative correct solutions from the model during training (L4)
  - Execute the program samples (L6)
  - Filter the samples by:
    - **Full correctness:** matches the gold final execution result (L7)
    - Duplication (L8): pruning out "trivial variants"
  - Save them in the buffer for learning (L9)

---

**Algorithm 1** Training Update

**Input:**

Parameterized model $P_\theta(y|x)$;

Executor $\mathcal{E} : \mathcal{Y} \rightarrow \mathcal{Z}$;

A training example $(x, y^*, z^*)$;

Buffer $\mathcal{B}$ for this input $x$

1: **if** $|\mathcal{B}| = 0$ **then**
2: $\quad \mathcal{B} \leftarrow \mathcal{B} + \{y^*\}$ /* initialize buffer */
3: **end if**
4: $\hat{Y} \leftarrow SampleSolutions(x, P_\theta, \mathcal{B})$
5: **for** $\hat{y}$ in $\hat{Y}$ **do**
6: $\quad \hat{z} \leftarrow \mathcal{E}(\hat{y})$ /* execute solution */
7: $\quad$ **if** $isCorrect(\hat{z}, z^*)$ **then**
8: $\quad\quad$ **if not** $isDuplicate(\hat{y}, \mathcal{B})$ **then**
9: $\quad\quad\quad \mathcal{B} \leftarrow \mathcal{B} + \hat{y}$ /* save to buffer */
10: $\quad\quad$ **end if**
11: $\quad$ **end if**
12: **end for**
13: $\theta \xleftarrow{\text{update}} \nabla_\theta \mathcal{L}(x, \mathcal{B}, P_\theta)$

- Objectives for learning from multiple targets
  - **MLE**: maximize the likelihood of generating the reference program;
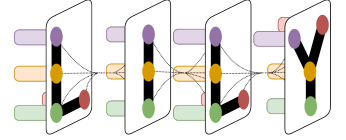  - **MLE-Aug**: simply summing the loss from the saved correct programs, it encourages the model to put equal weights on all targets;
  - **MML**: maximize the marginal likelihood of all saved correct solutions, but note that the gradient is in proportion to the likelihood;
  - **$\beta$-MML** (Guu et al., 2017): an interpolation between MML and MLE-Aug, with $\beta \in (0, 1]$

| Name | Loss Functions $\mathcal{L}(x, \mathcal{B}, P_\theta)$ | Gradients $\nabla_\theta(x, \mathcal{B}, P_\theta)$ |
|---|---|---|
| MLE | $-\log P_\theta(y^*|x)$ | $-\nabla_\theta \log P_\theta(y^*|x)$ |
| MLE-Aug | $-\sum_{\hat{y}\in\mathcal{B}} \log P_\theta(\hat{y}|x)$ | $-\sum_{\hat{y}\in\mathcal{B}} \nabla_\theta \log P_\theta(\hat{y}|x)$ |
| MML | $-\log \sum_{\hat{y}\in\mathcal{B}} P_\theta(\hat{y}|x)$ | $-\sum_{\hat{y}\in\mathcal{B}} \frac{P_\theta(\hat{y}|x)}{\sum_{\tilde{y}\in\mathcal{B}} P_\theta(\tilde{y}|x)} \nabla_\theta \log P_\theta(\hat{y}|x)$ |
| $\beta$-MML | $-\frac{1}{\beta} \log \sum_y P_\theta(\hat{y}|x)^\beta$ | $-\sum_{\hat{y}\in\mathcal{B}} \frac{P_\theta(\hat{y}|x)^\beta}{\sum_{\tilde{y}\in\mathcal{B}} P_\theta(\tilde{y}|x)^\beta} \nabla_\theta \log P_\theta(\hat{y}|x)$ |

Different loss functions and gradients used for learning from multiple targets

Yale

- Learning from self-sampled correct solutions are great, but…
  - It is also hard to sample, especially for complex programs, it is hard to "creatively" write a different program that is also correct.
  - There are many failed attempts for the model to be creative, and they were almost there!
  - Can we learn from those programs as well?
    - **YES!**

- **Intermediate state $s_i$ is the set of all variables values in the scope after executing the first $i$ steps as the *program prefix $y_{\leq i}$***
  - Note: the state representation is name-agnostic since variable names do not typically contributes to the semantics of the solutions

A goods train runs at a speed of 72kmph and crosses a 250M long platform in 26 seconds. What is the length of the goods train?

```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n1 / t0
t2 = n2 - t1
answer = t0 * t2
```
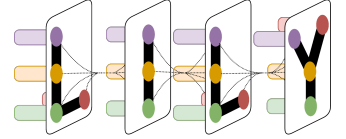
Execute

```
{72}
{72, 250}
{72, 250, 26}
{72, 250, 26, 20.0}
{72, 250, 26, 20.0, 12.5}
{72, 250, 26, 20.0, 12.5, 13.5}
output = 270
```

**Solutions**

**Intermediate States**

Yale

- Prefixes of two programs $y_{\leq i}$ and $y'_{\leq j}$ are **semantically equivalent** if and only if $s_i = s'_j$

  - i.e., those two program prefixes produces the exact same set of variable values



Solutions         Intermediate States

- A program prefix $y_{\leq i}$ is **partially-correct** if and only if it is semantically equivalent to the prefix of a *known partially-correct solution $y^*_{\leq j}$*
  - Since we keep all known partially-correct solutions in buffer $\mathcal{B}$, we have

$$PartiallyCorrect(y_{\leq i}) \iff \exists y^* \in \mathcal{B}. \ \exists j \leq |y^*| \ s.t. \ s^*_j = s_i$$

**Known PC Solution**

```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 * 0.2778
t1 = n1 / t0
t2 = n2 - t1
answer = t0 * t2
```

**Evaluate with Executor**

```
{72}
{72, 250}
{72, 250, 26}
{72, 250, 26, 20.0}
{72, 250, 26, 20.0, 12.5}
{72, 250, 26, 20.0, 12.5, 13.5}
output = 270
```
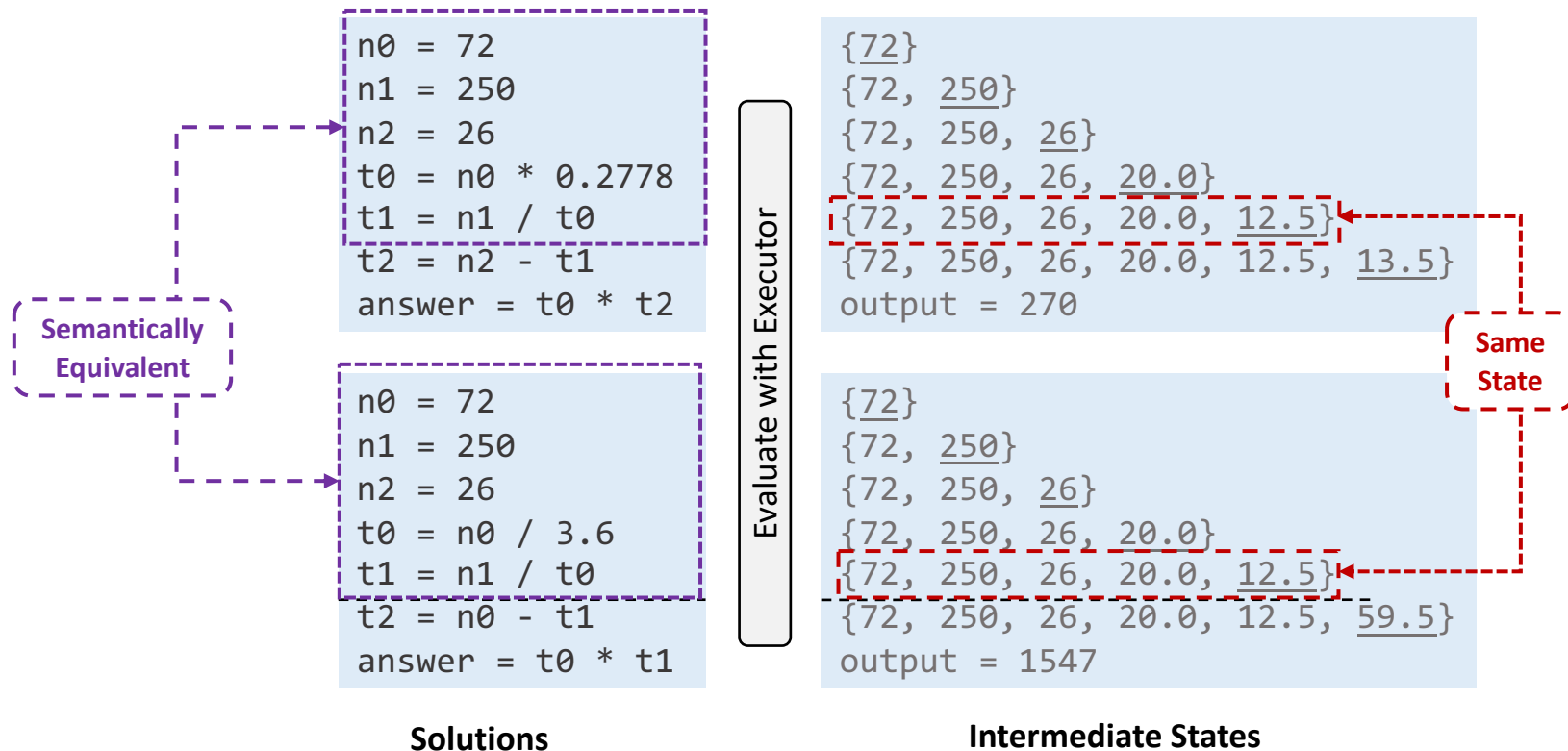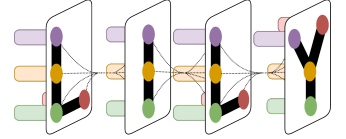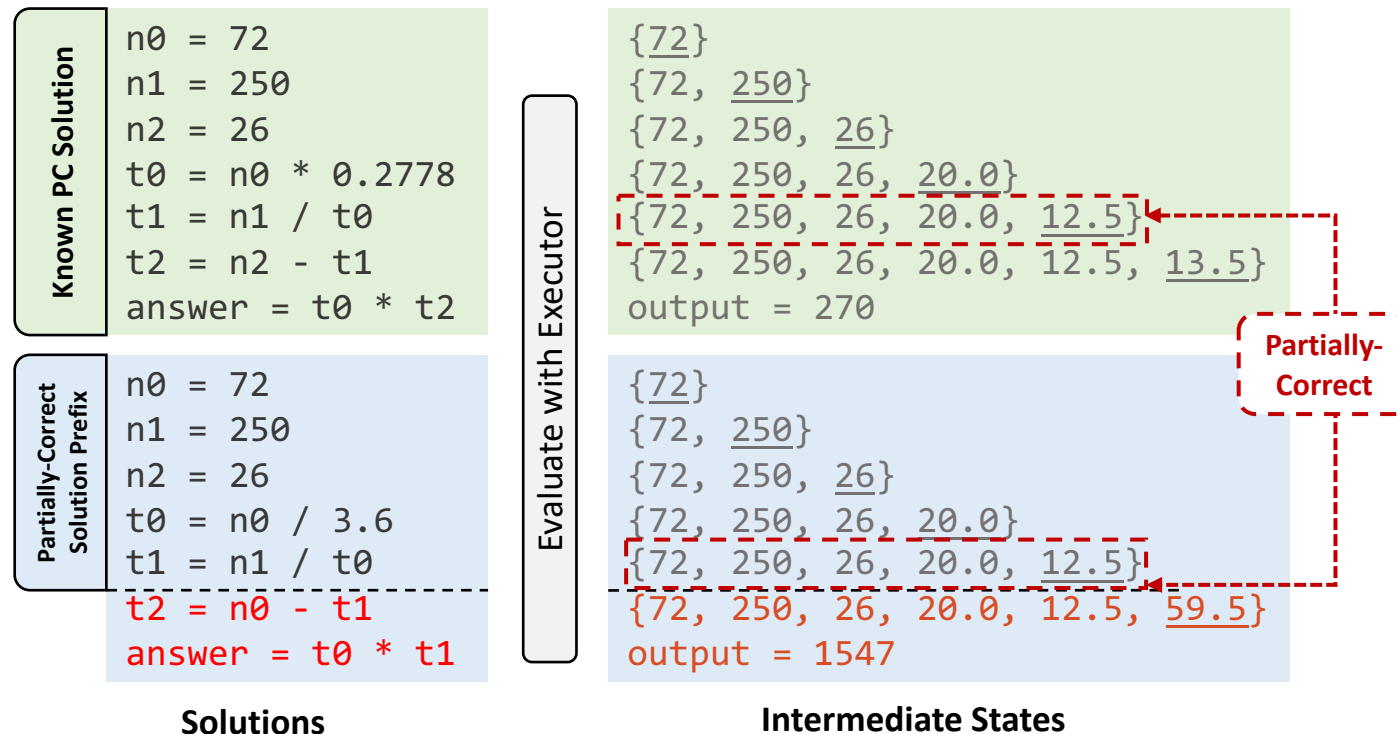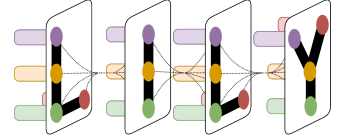
**Partially-Correct Solution Prefix**

```
n0 = 72
n1 = 250
n2 = 26
t0 = n0 / 3.6
t1 = n1 / t0
t2 = n0 - t1
answer = t0 * t1
```

```
{72}
{72, 250}
{72, 250, 26}
{72, 250, 26, 20.0}
{72, 250, 26, 20.0, 12.5}
{72, 250, 26, 20.0, 12.5, 59.5}
output = 1547
```

**Partially-Correct**

**Solutions**

**Intermediate States**

Yale

# Learning from Partially-Correct Solutions (PCSs)

- Modification to the main algorithm
  - Guided-sampling from known PCS prefix

---

**Algorithm 2** $SampleSolutions(x, P_\theta, \mathcal{B})$ with partially-correct solutions

**Input:** Model $P_\theta(y|x)$; the NL input $x$ and a set of partially-correct solutions $\mathcal{B}$
**Output:** Solution samples $\hat{Y}$.
1: Select $\hat{y}_{\leq i} \in \mathcal{B} \setminus \{\hat{y}|\mathcal{E}(\hat{y}) = z^*\}$ uniformly at random /* sample PCS prefix for completion */
2: Sample a set of completions $Y_p \sim P_\theta(\hat{y}_{>i}|\hat{y}_{\leq i}, x)$
3: $\hat{Y} \leftarrow \{[\hat{y}_{\leq i}||\hat{y}_{>i}]\}_{\hat{y}_{>i} \in Y_p}$ /* concatenate completions with the solution prefix */
4: **return** $\hat{Y}$

---

  - Identify partially-correct program prefixes
  - Filtering solution prefixes
    - PCS is only saved if it is not a prefix of any know PCS
  - Learning objectives
    - With an auto-regressive generation model, the learning of $P(y_{\leq i}|x)$ is independent of $y_{>i}$, thus no change to the learning objectives are required.

---

**Algorithm 1** Training Update

**Input:**
  Parameterized model $P_\theta(y|x)$;
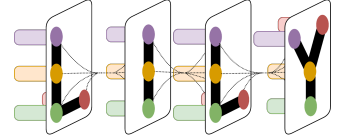  Executor $\mathcal{E} : \mathcal{Y} \to \mathcal{Z}$;
  A training example $(x, y^*, z^*)$;
  Buffer $\mathcal{B}$ for this input $x$
1: **if** $|\mathcal{B}| = 0$ **then**
2:     $\mathcal{B} \leftarrow \mathcal{B} + \{y^*\}$ /* initialize buffer */
3: **end if**
4: $\hat{Y} \leftarrow SampleSolutions(x, P_\theta, \mathcal{B})$
5: **for** $\hat{y}$ **in** $\hat{Y}$ **do**
6:     $\hat{z} \leftarrow \mathcal{E}(\hat{y})$ /* execute solution */
7:     **if** $isCorrect(\hat{z}, z^*)$ **then**
8:       **if not** $isDuplicate(\hat{y}, \mathcal{B})$ **then**
9:         $\mathcal{B} \leftarrow \mathcal{B} + \hat{y}$ /* save to buffer */
10:       **end if**
11:     **end if**
12: **end for**
13: $\theta \xleftarrow{\text{update}} \nabla_\theta \mathcal{L}(x, \mathcal{B}, P_\theta)$
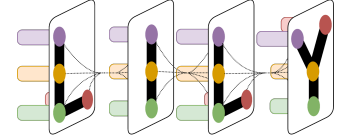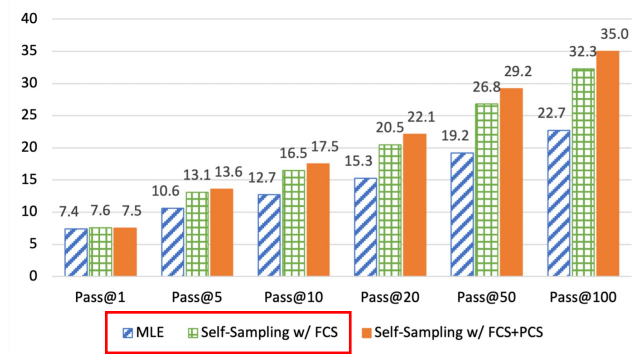
Yale

# Experimental Setup

- Datasets:
  - ➢ **MathQA-Python-Filtered**: we do template-based deduplication is applied to the original MathQA-Python dataset to better evaluation generalization
  - ➢ **GSM5.5K-Python**: we automatically converted the natural language formulas to program solutions in the same style as MathQA-Python

- Language model:
  - We use **GPT-Neo 125M and 2.7B** as our main LM to study

- Evaluation metrics:
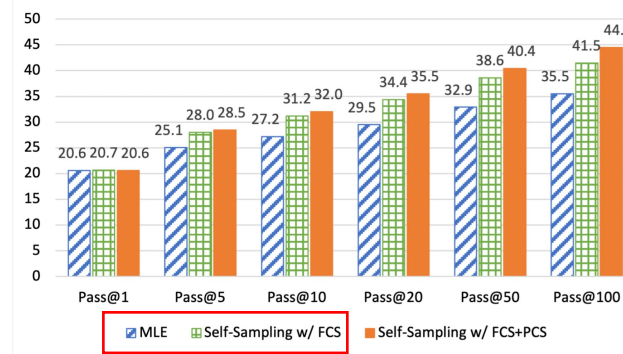  - We use **pass@k** following recent work in math reasoning and program synthesis
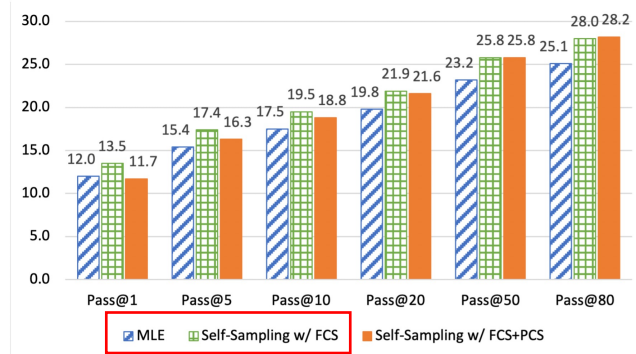
# Main Results

- Learning from self-sampled solutions improves pass@k
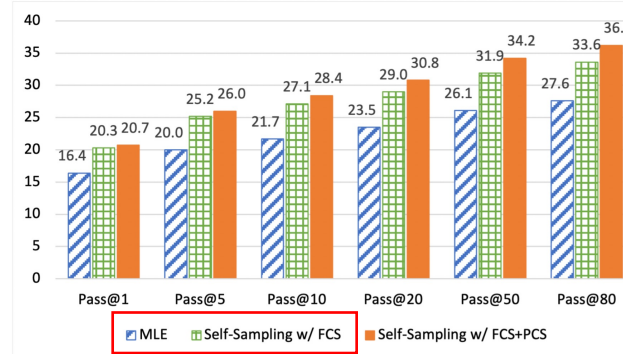


(a) GSM5.5K-Python with GPT-Neo 125M

(b) GSM5.5K-Python with GPT-Neo 2.7B
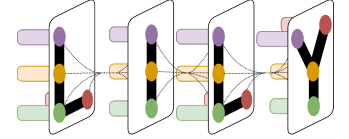
(c) MathQA-Python-Filtered with GPT-Neo 125M

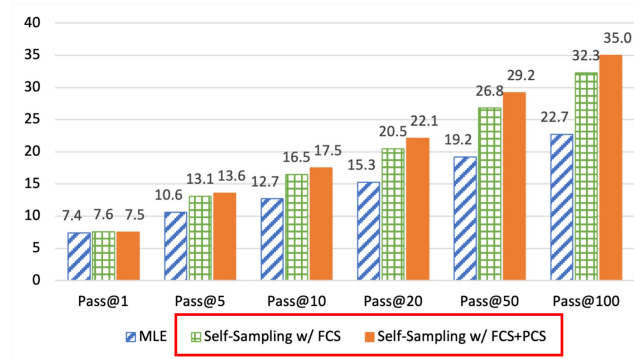(d) MathQA-Python-Filtered with GPT-Neo 2.7B

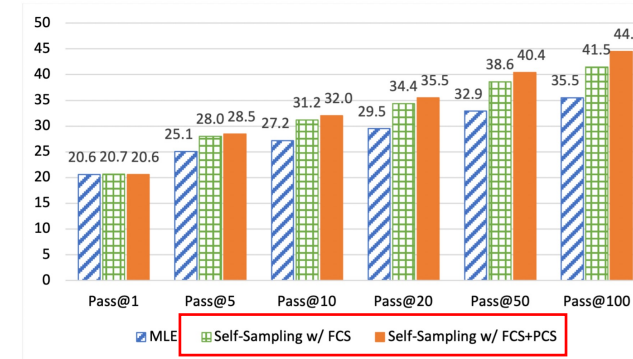Comparing self-sampling with MLE baseline
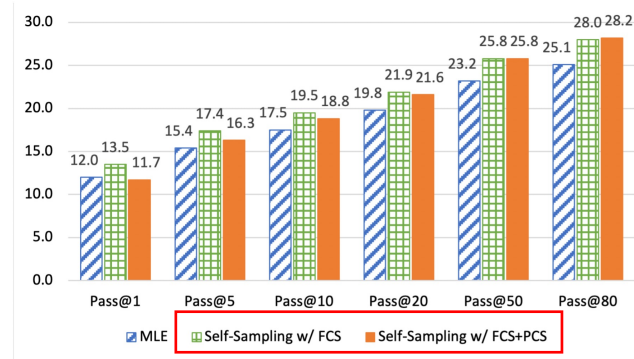
# Main Results

- Partially-correct solutions further improve model performance
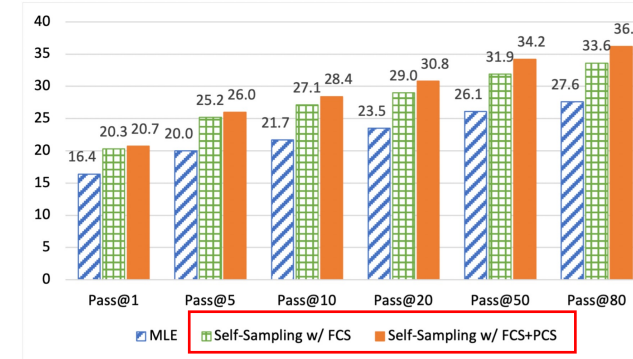


(a) GSM5.5K-Python with GPT-Neo 125M

(b) GSM5.5K-Python with GPT-Neo 2.7B

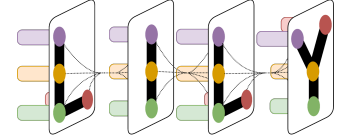(c) MathQA-Python-Filtered with GPT-Neo 125M
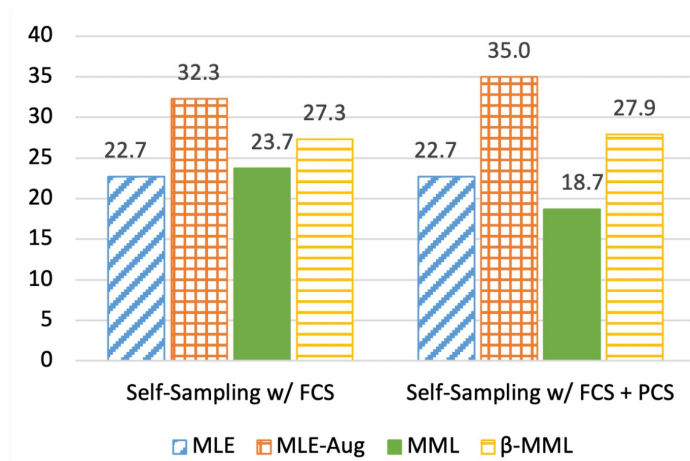
(d) MathQA-Python-Filtered with GPT-Neo 2.7B

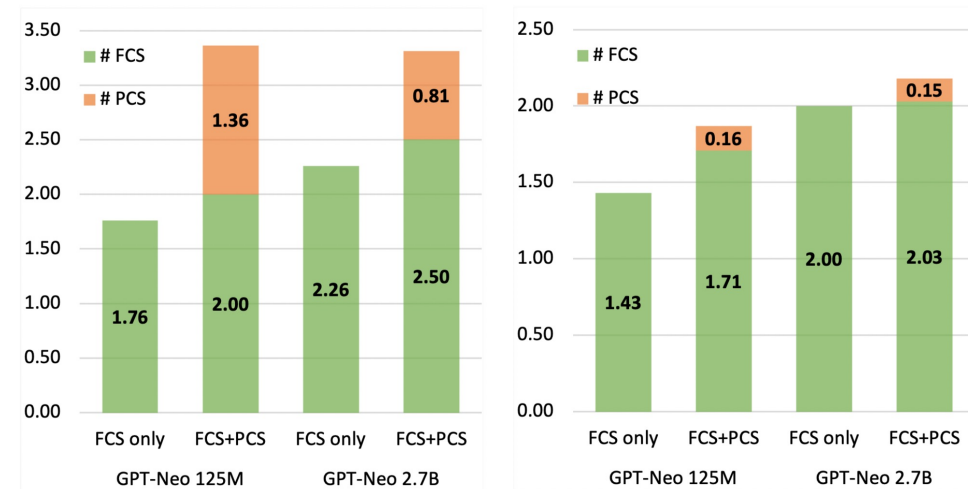Comparing self-sampling with MLE baseline

- MLE-Aug loss function works the best
  - It draws learning signal equally from all saved solutions
  - MML works the worse, especially when also learning from PCSs

- Dynamics between # of PCSs and FCSs saved in the buffer
  - More saved solutions typically results in better pass@k performance
  - Large models are better at completing PCS prefixes to be FCS



GSM5.5K-Python with finetuned
GPT-Neo 125M model

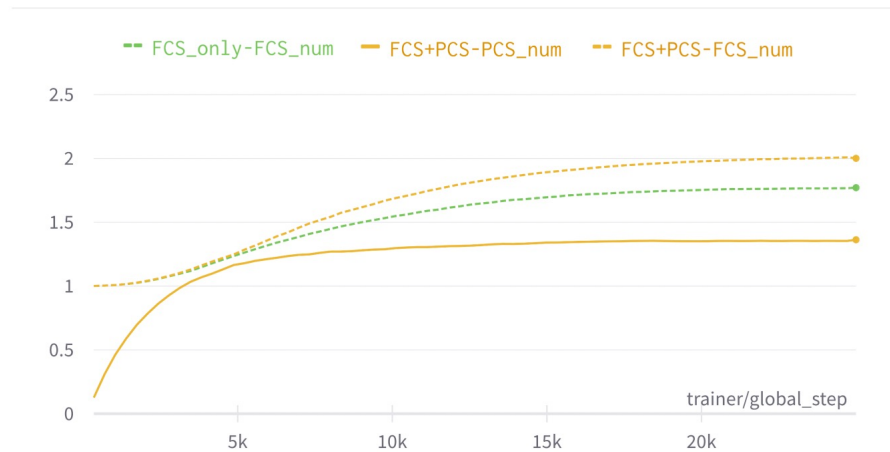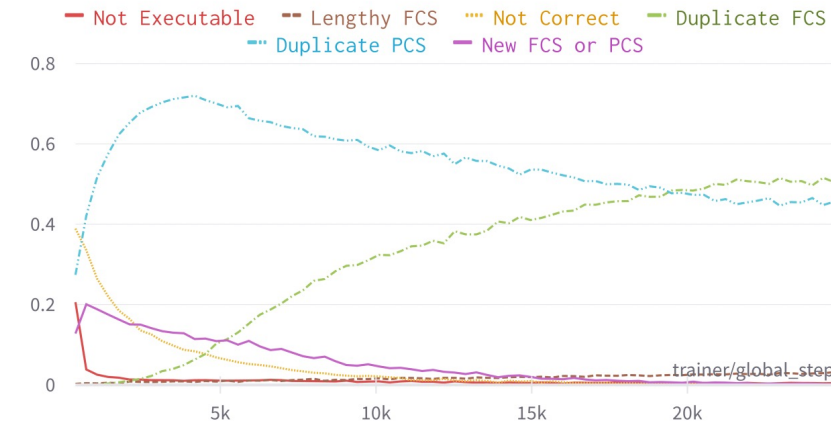# saved FCSs and PCSs per problem for GSM5.5K-Python (left)
and MathQA-Python-Filtered (right)

- Dynamics between # of PCSs and FCSs saved in the buffer
  - Partially-correct solutions helps learning especially in early stages



Growth of the number of saved
FCS and PCS during training



Distribution of the characterization of self-
sampled solutions during training

Yale

- Learning from self-sampled solutions can be useful given the right constraints
  - E.g., when you can easily prune out incorrect and duplicated ones
- Programs are not either correct or wrong, they can be partially-correct.
  - Note that our definition of partial correctness is different from say, passing 60% of the test cases, because that program would still be wrong;
  - Instead, by comparing execution traces, we identify the first 60% of the program is on the right track

# Task: Language-to-Code Generation

- Task: **language-to-code generation** using LLMs in few-shot learning
  - Cornerstone for many tasks in NLP and ML

**Complex question**
What are the name and budget of the departments with average instructor salary greater than the overall average?

**Complex SQL**
```
SELECT T2.name, T2.budget
FROM instructor as T1 JOIN department as
T2 ON T1.department_id = T2.id
GROUP BY T1.department_id
HAVING avg(T1.salary) >
        (SELECT avg(salary) FROM instructor)
```

Spider (Yu et al., 2018)

| Year | City | Country | Nations |
|------|------|---------|---------|
| 1896 | Athens | Greece | 14 |
| 1900 | Paris | France | 24 |
| 1904 | St. Louis | USA | 12 |
| ... | ... | ... | ... |
| 2004 | Athens | Greece | 201 |
| 2008 | Beijing | China | 204 |
| 2012 | London | UK | 204 |

$x$ = Greece held its last Summer Olympics in which year?

$y$ = 2004

WikiTableQuestions (Pasupat and Liang., 2015)

**Problem:** Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?
**Solution:** Beth bakes 4 2 dozen batches of cookies for a total of 4*2 = <<4*2=8>>8 dozen cookies
There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of 12*8 = <<12*8=96>>96 cookies
She splits the 96 cookies equally amongst 16 people so they each eat 96/16 = <<96/16=6>>6 cookies
**Final Answer:** 6

**Problem:** Mrs. Lim milks her cows twice a day. Yesterday morning, she got 68 gallons of milk and in the evening, she got 82 gallons. This morning, she got 18 gallons fewer than she had yesterday morning. After selling some gallons of milk in the afternoon, Mrs. Lim has only 24 gallons left. How much was her revenue for the milk if each gallon costs $3.50?
Mrs. Lim got 68 gallons - 18 gallons = <<68-18=50>>50 gallons this morning.
So she was able to get a total of 68 gallons + 82 gallons + 50 gallons = <<68+82+50=200>>200 gallons.
She was able to sell 200 gallons - 24 gallons = <<200-24=176>>176 gallons.
Thus, her total revenue for the milk is $3.50/gallon x 176 gallons = $<<3.50*176=616>>616.
**Final Answer:** 616

GSM8k (Cobbe et al., 2021)

prompt:
Write a python function to check if a given number is one less than twice its reverse. Your code should satisfy these tests:

**assert** check(70) == False
**assert** check(23) == False
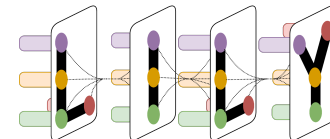**assert** check(73) == True

model:
```
def check(n) :
    if n == 2*int(str(n)[::-1])-1 :
        return True
    else :
        return False
```

MBPP (Austin et al., 2021)

- Task: natural language to code generation using large language models without parameter updates (i.e., finetuning)

- Motivation:
  - The CodeLMs are trained on surface code, how do we incorporate execution semantics into the generation process?
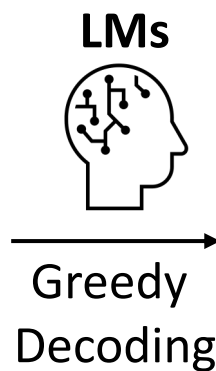  - The cost for finetuning LLMs are huge, how do we improve them without changing their parameters?



**Few-shot Exemplars**
```
-- Translate natural language
question into SQL Query

-- Example
-- NL: What ...
SELECT ...
```

+

**Task Input**
```
-- Example
--NL: How many students in the
class are between 20 and 30
years old?
```
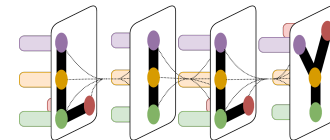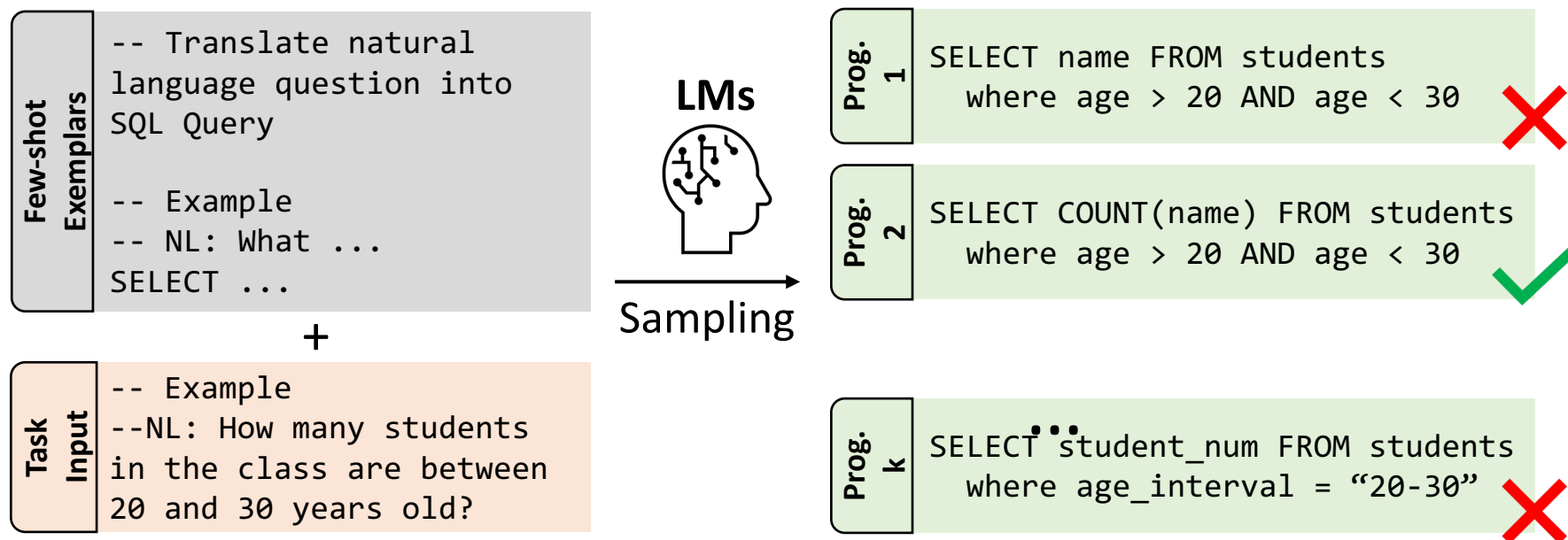
**LMs**

Greedy Decoding →

**Prog.**
```
SELECT name FROM students
    where age > 20 AND age < 30
```
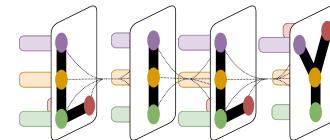❌

- Task: natural language to code generation using large language models without parameter updates (i.e., finetuning)

- Observation:
  - While CodeLMs struggles with precision in the few-shot setting, it often produces the correct output when enough samples are drawn.

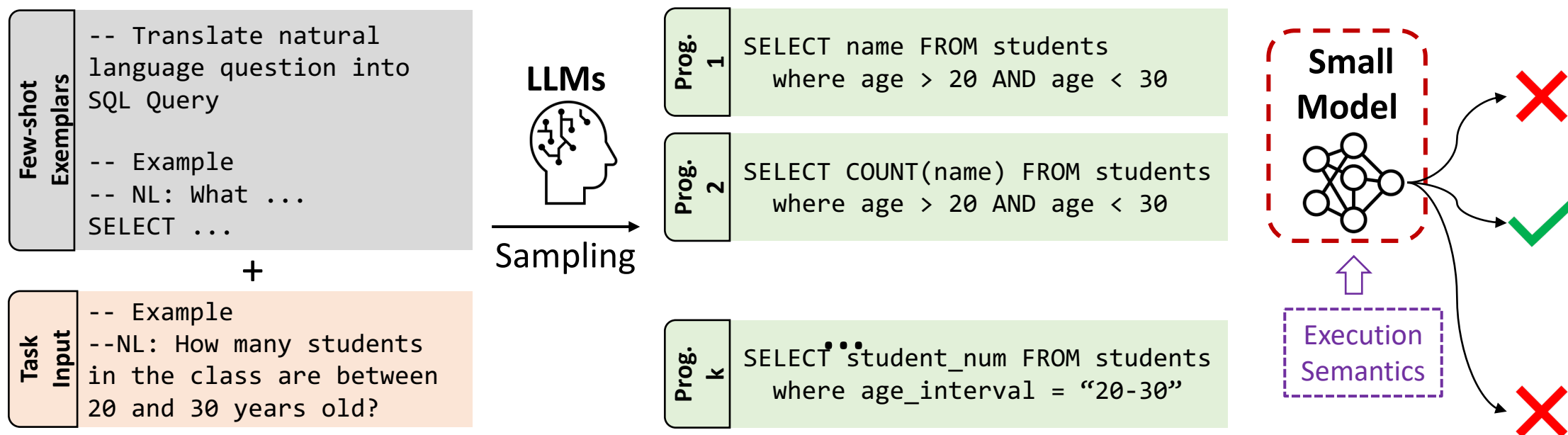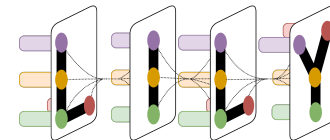- If we can not directly finetune the LLMs...
  - Can we train *a separate (relatively) small model* as an "add-on", to rerank the programs samples from LLMs?
  - Can we incorporate execution semantics in this model instead?

- We propose LEVER, which <u>le</u>arns to <u>veri</u>fy language-to-code generation by LLMs trained code (CodeLMs), with the help of execution.

- LEVER has three main steps: **1) Generation**

- We propose LEVER, which <u>le</u>arns to <u>ver</u>ify language-to-code generation by LLMs trained code (CodeLMs), with the help of execution.

- LEVER has three main steps: 1) Generation; **2) Execution**
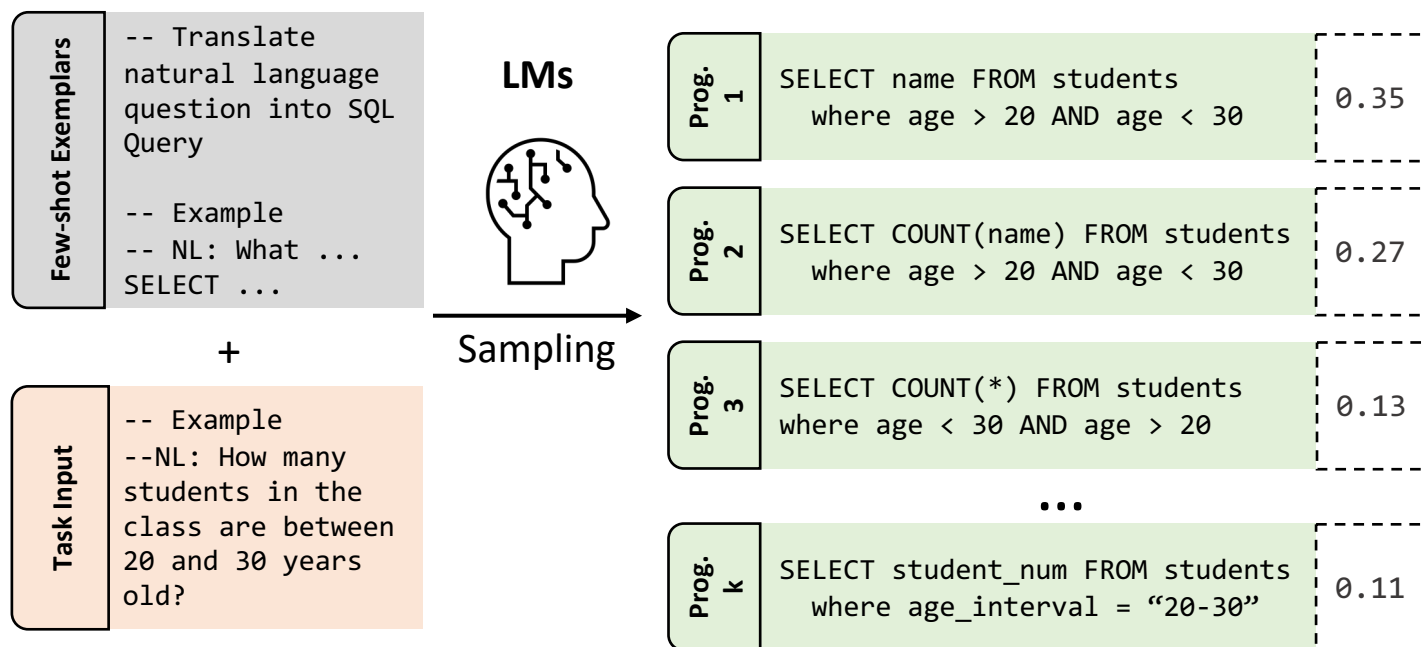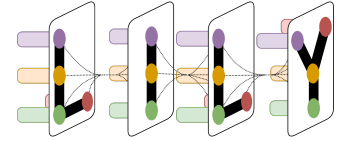
- We propose LEVER, which <u>le</u>arns to <u>ver</u>ify language-to-code generation by LLMs trained code (CodeLMs), with the help of execution.

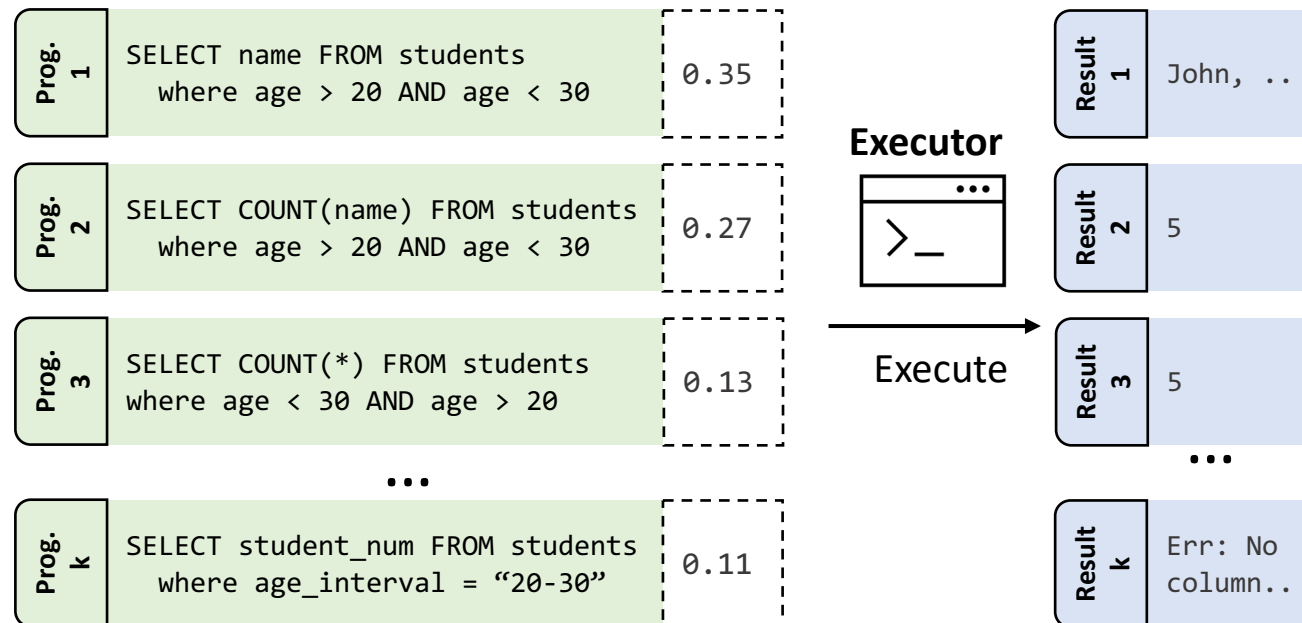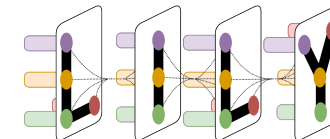- LEVER has three main steps: 1) Generation; 2) Execution; **3) Verification**
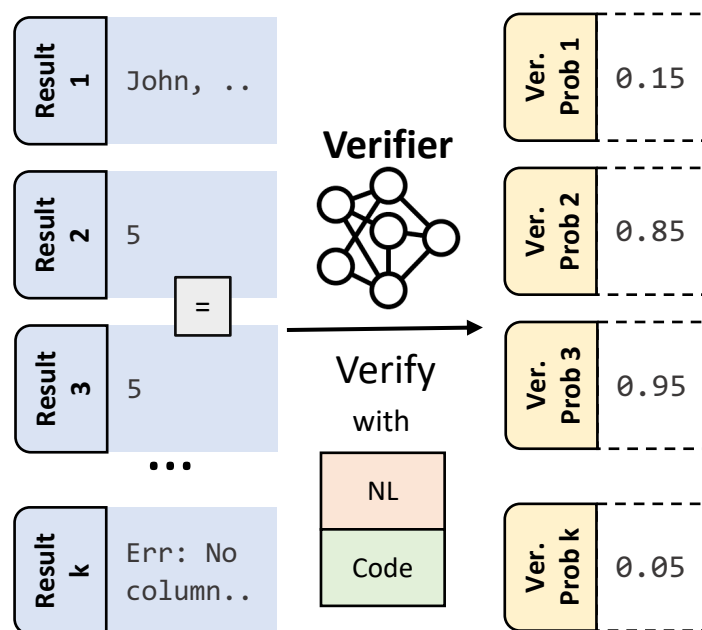
- Detailed formulation:
  - We parameterize the verifier as **a binary classifier**, with the input as:
    - natural language $x$; program sample $\hat{y}$; and its execution result $\mathcal{E}(\hat{y})$

$$P_\theta(v_{=1}|x, \hat{y}, \mathcal{E}(\hat{y}))$$

  - Given the input $x$ and a program sample $\hat{y} \in S$, we obtain the **reranking probability** as the joint probability of generation and passing verification:

$$P_R(\hat{y}, v_{=1}|x) = P_{\mathbf{LM}}(\hat{y}|x) \cdot P_\theta(v_{=1}|x, \hat{y}, \mathcal{E}(\hat{y}))$$

  - We further aggregate the reranking probability of the programs in the samples that executes to the same result, and obtains **the final score**

$$R(x, \hat{y}) = P_R(\mathcal{E}(\hat{y}), v_{=1}|x) = \sum_{y \in S, \mathcal{E}(y)=\mathcal{E}(\hat{y})} P_R(y, v_{=1}|x)$$

Yale

- ## Detailed formulation:
  - We further aggregate the reranking probability of the programs in the samples that executes to the same result, and obtains **the final score**

- Training data creation - for each input $x$:
  - Sample program candidates on the training set examples $\widehat{y_i} \sim P_{LM}(y|x)$
  - Execute the programs to obtain their execution results $\widehat{z_i} = \mathcal{E}(\widehat{y_i})$
  - Use gold exec. result or test cases to automatically label their correctness $v_i$
  - We created a set of training examples $\{x, \widehat{y_i}, \widehat{z_i}, v_i\}_{i=1}^{n}$ for each input $x$

- Learning objective:
  - Negative log-likelihood, normalized by the number of program candidates

$$\mathcal{L}_\theta(x, S) = -\frac{1}{|S|} \cdot \sum_{\hat{y}_i \in S} \log P_\theta(v_i | x, \hat{y}_i, \hat{z}_i)$$

Yale

- ## Datasets:

  **SQL**
  - ➤ Spider (Yu et al., 2018): text-to-SQL semantic parsing;
  - ➤ WikiTQ (Pasupat et al., 2015): table question answering

  **Python**
  - ➤ GSM8k (Cobbe et al., 2021): math word problems
  - ➤ MBPP (Austin et al., 2021): basic python programming

**Complex question**  What are the name and budget of the departments with average instructor salary greater than the overall average?

**Complex SQL**
```
SELECT T2.name, T2.budget
FROM instructor as T1 JOIN department as
T2 ON T1.department_id = T2.id
GROUP BY T1.department_id
HAVING avg(T1.salary) >
    (SELECT avg(salary) FROM instructor)
```

Spider

| Year | City | Country | Nations |
|------|------|---------|---------|
| 1896 | Athens | Greece | 14 |
| 1900 | Paris | France | 24 |
| 1904 | St. Louis | USA | 12 |
| ... | ... | ... | ... |
| 2004 | Athens | Greece | 201 |
| 2008 | Beijing | China | 204 |
| 2012 | London | UK | 204 |

$x$ = Greece held its last Summer Olympics in which year?

$y$ = 2004

WikiTQ

**Problem:** Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?
**Solution:** Beth bakes 4 2 dozen batches of cookies for a total of 4*2 = <<4*2=8>>8 dozen cookies
There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of 12*8 = <<12*8=96>>96 cookies
She splits the 96 cookies equally amongst 16 people so they each eat 96/16 = <<96/16=6>>6 cookies
**Final Answer:** 6

**Problem:** Mrs. Lim milks her cows twice a day. Yesterday morning, she got 68 gallons of milk and in the evening, she got 82 gallons. This morning, she got 18 gallons fewer than she had yesterday morning. After selling some gallons of milk in the afternoon, Mrs. Lim has only 24 gallons left. How much was her revenue for the milk if each gallon costs $3.50?
Mrs. Lim got 68 gallons - 18 gallons = <<68-18=50>>50 gallons this morning.
So she was able to get a total of 68 gallons + 82 gallons + 50 gallons = <<68+82+50=200>>200 gallons.
She was able to sell 200 gallons - 24 gallons = <<200-24=176>>176 gallons.
Thus, her total revenue for the milk is $3.50/gallon x 176 gallons = $<<3.50*176=616>>616.
**Final Answer:** 616

GSM8k

**prompt**  Write a python function to check if a given number is one less than twice its reverse. Your code should satisfy these tests:

**assert** check(70) == False
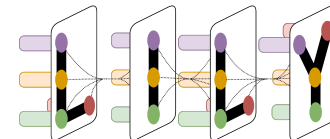**assert** check(23) == False
**assert** check(73) == True

**model**
```
def check(n) :
    if n == 2*int(str(n)[::-1])-1 :
        return True
    else :
        return False
```
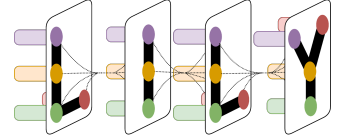
MBPP

- Datasets: 4 language-to-code benchmarks from different domains
- CodeLMs:
  - <u>Codex-davinci-002</u>: best CodeLM available, accessible through API
  - <u>InCoder-6B</u>: open-source
  - <u>CodeGen-16B</u>: open-source
- Evaluation metric
  - Execution Accuracy (i.e., pass@1)
- Baselines:
  - <u>Greedy</u>: choose most likely token per decoding step
  - Maximum Likelihood (<u>ML</u>): choose the program with highest generation prob.
  - Error Pruning + ML (<u>EP+ML</u>): first prune out the programs with execution error
  - <u>EP + Voting</u>: majority vote of the error-free execution results

Yale

- ## Codex + LEVER achieves new SoTA results on all 4 benchmarks
  - ### It outperforms all previous <u>finetuning</u> and <u>few-shot learning</u> results

| Methods | Dev |
|---|---|
| *Previous Work without Finetuning* | |
| Rajkumar et al. (2022) | 67.0 |
| MBR-Exec (Shi et al., 2022) | 75.2 |
| Coder-Reviewer (Zhang et al., 2022) | 74.5 |
| *Previous Work with Finetuning* | |
| T5-3B (Xie et al., 2022) | 71.8 |
| PICARD (Scholak et al., 2021) | 75.5 |
| RASAT (Qi et al., 2022) | 80.5 |
| *This Work with code-davinci-002* | |
| Greedy | 75.3 |
| EP + ML | 77.3 |
| LEVER🎏 | $\mathbf{81.9}_{\pm 0.1}$ |

Spider

| Methods | Dev | Test |
|---|---|---|
| *Previous Work without Finetuning* | | |
| Codex QA* (Cheng et al., 2022) | 50.5 | 48.7 |
| Codex SQL (Cheng et al., 2022) | 60.2 | 61.1 |
| Codex Binder (Cheng et al., 2022) | **65.0** | 64.6 |
| *Previous Work with Finetuning* | | |
| TaPEX* (Liu et al., 2021) | 60.4 | 59.1 |
| TaCube (Zhou et al., 2022) | 61.1 | 61.3 |
| OmniTab* (Jiang et al., 2022) | - | 63.3 |
| *This Work with code-davinci-002* | | |
| Greedy | 49.6 | 53.0 |
| EP + ML | 52.7 | 54.9 |
| LEVER🎏 | $64.6_{\pm 0.2}$ | $\mathbf{65.8}_{\pm 0.2}$ |

WikiTQ

| Methods | Dev | Test |
|---|---|---|
| *Previous Work without Finetuning* | | |
| PAL (Gao et al., 2022) | - | 72.0 |
| Codex + SC[†] (Wang et al., 2022) | - | 78.0 |
| PoT-SC (Chen et al., 2022b) | - | 80.0 |
| *Previous Work with Finetuning* | | |
| Neo-2.7B + SS (Ni et al., 2022) | 20.7 | 19.5 |
| Neo-1.3B + SC (Welleck et al., 2022) | - | 24.2 |
| DiVeRSe*[†] (Li et al., 2022b) | - | 83.2 |
| *This Work with codex-davinci-002* | | |
| Greedy | 68.1 | 67.2 |
| EP + ML | 72.1 | 72.6 |
| LEVER🎏 | $\mathbf{84.1}_{\pm 0.2}$ | $\mathbf{84.5}_{\pm 0.3}$ |

GSM8k

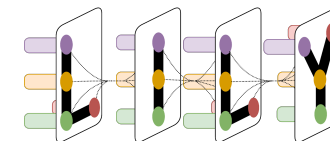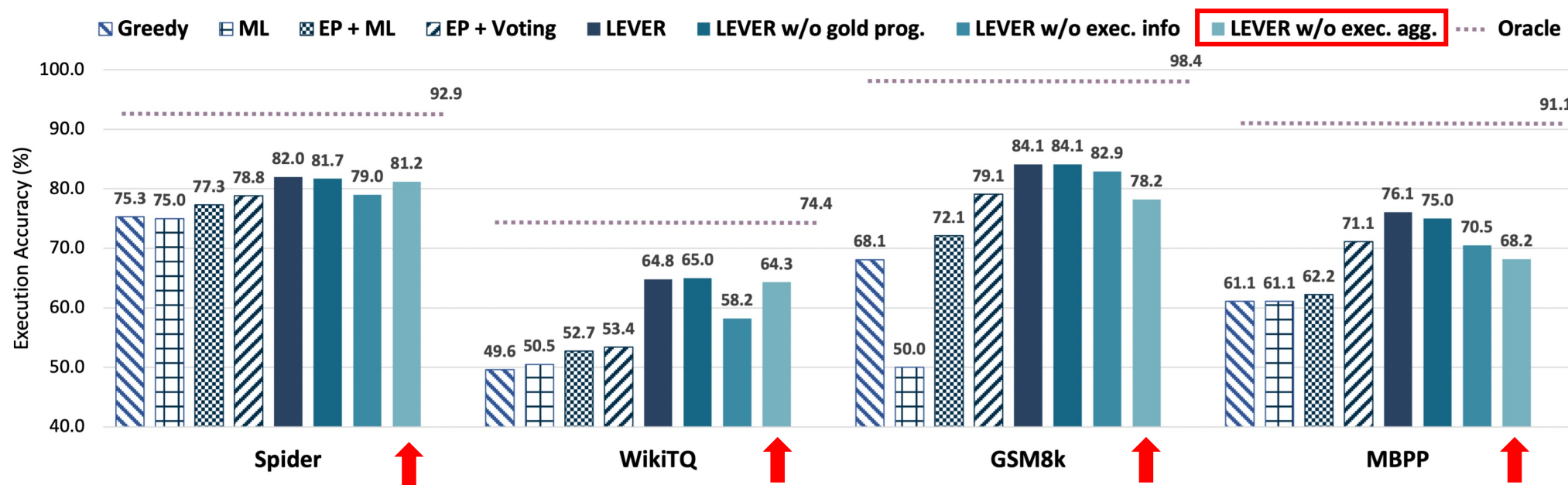| Methods | Dev | Test |
|---|---|---|
| *Previous Work without Finetuning* | | |
| MBR-Exec (Shi et al., 2022) | - | 63.0 |
| Reviewer (Zhang et al., 2022) | - | 66.9 |
| *This Work with codex-davinci-002* | | |
| Greedy | 61.1 | 62.0 |
| EP + ML | 62.2 | 60.2 |
| LEVER🎏 | $\mathbf{75.4}_{\pm 0.7}$ | $\mathbf{68.9}_{\pm 0.4}$ |

MBPP

Yale

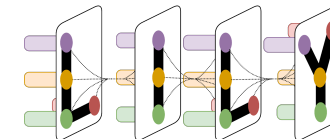- Execution information are crucial to the performance improvement

- Execution information are crucial to the performance improvement
- Exec. agg. works well for Python but not SQL generation datasets

- Execution information are crucial to the performance improvement
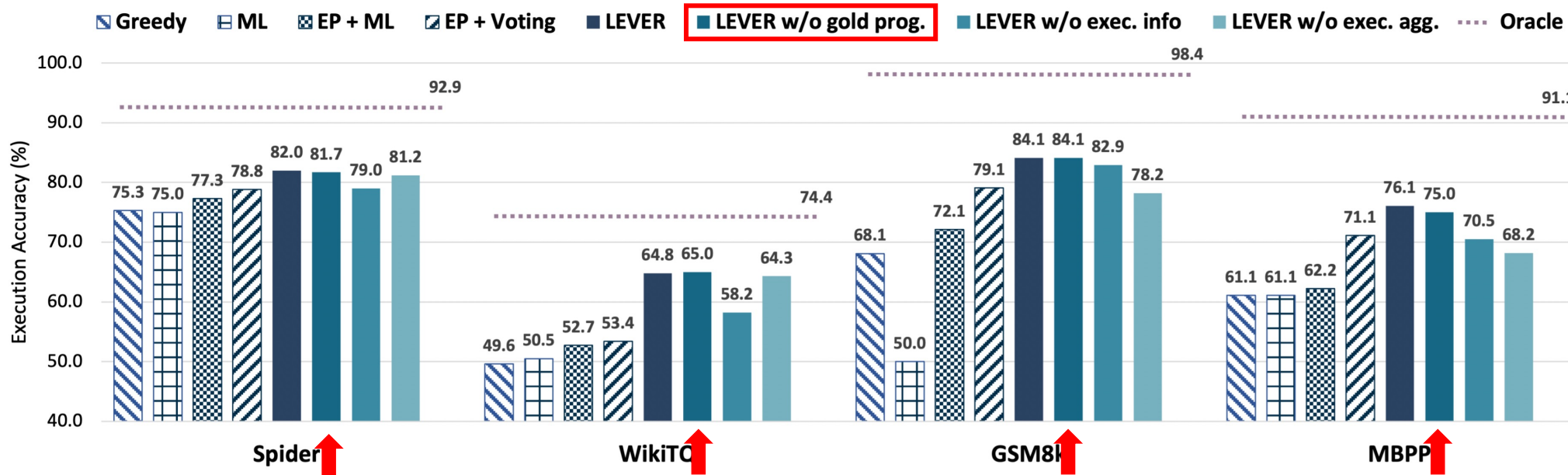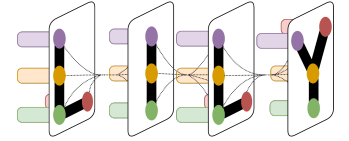- Exec. agg. works well for Python but not SQL generation datasets
- LEVER works well with weakly-supervised setting, where gold programs are not provided for learning
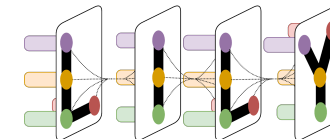


44

- Even larger improvements (e.g., up to 30.0%) are observed for InCoder and CodeGen models;

- Similar findings for ablation study
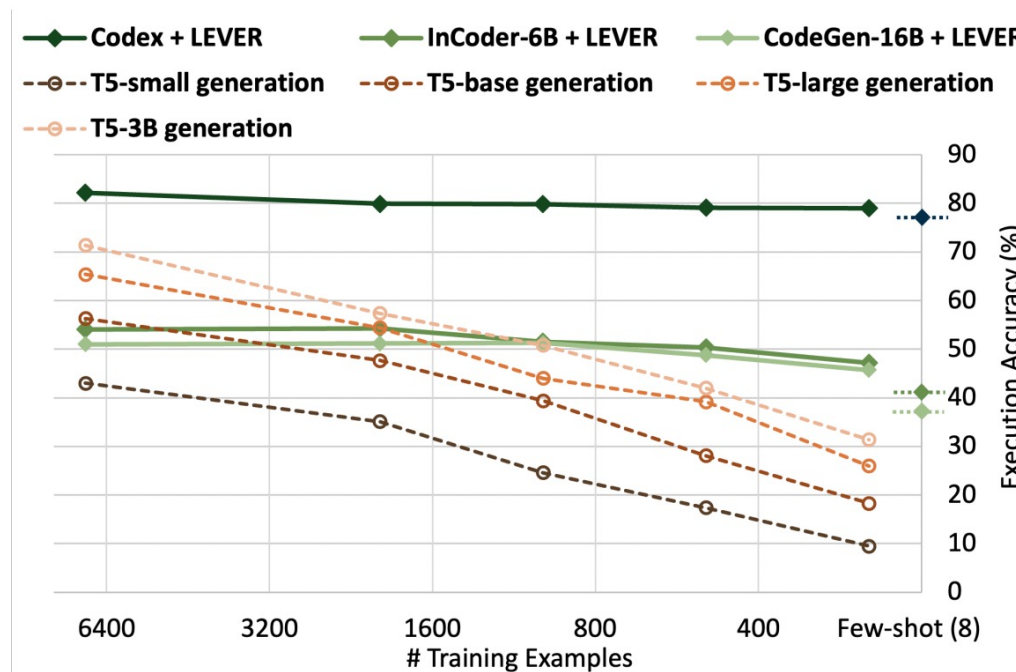  - With the exception that voting & exec. agg. methods decreases the performance

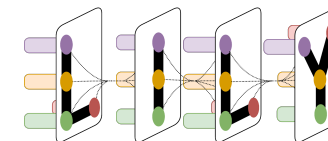| Methods | InCoder-6B | | CodeGen-16B | |
|---|---|---|---|---|
| | Spider | GSM8k | Spider | GSM8k |
| *Previous work:* | | | | |
| MBR-EXEC | 38.2 | - | 30.6 | - |
| Reviewer | 41.5 | - | 31.7 | - |
| *Baselines*: | | | | |
| Greedy | 24.1 | 3.1 | 24.6 | 7.1 |
| ML | 33.7 | 3.8 | 31.2 | 9.6 |
| EP + ML | 41.2 | 4.4 | 37.7 | 11.4 |
| EP + Voting | 37.4 | 5.9 | 37.1 | 14.2 |
| LEVER🎣 | 54.1 | **11.9** | 51.0 | **22.1** |
| — gold prog. | 53.4 | - | **52.3** | - |
| — exec. info | 48.5 | 5.6 | 43.0 | 13.4 |
| — exec. agg. | **54.7** | 10.6 | 51.6 | 18.3 |
| Oracle | 71.6 | 48.0 | 68.6 | 61.4 |

- LEVER works better under few-resource settings than generative finetuning
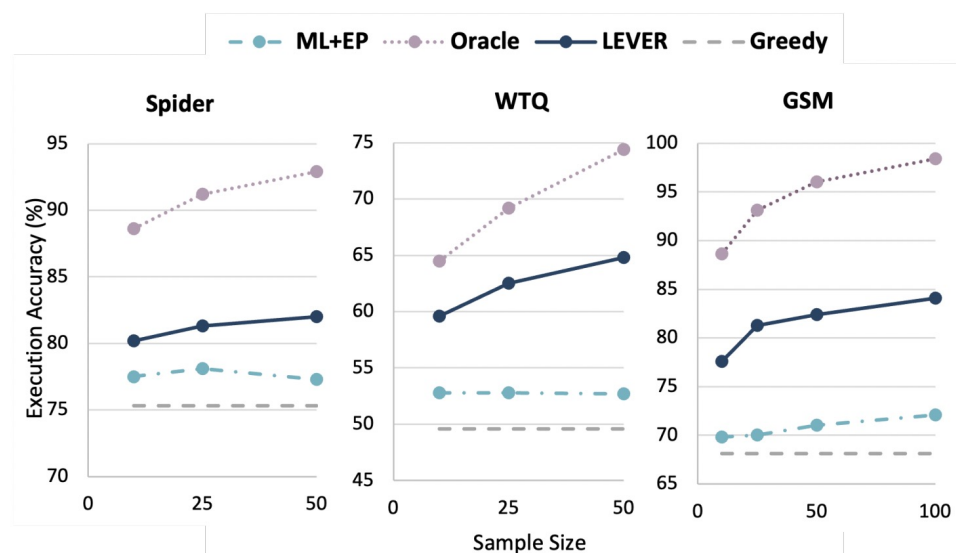


Verification vs. Generation w/ decreasing training data
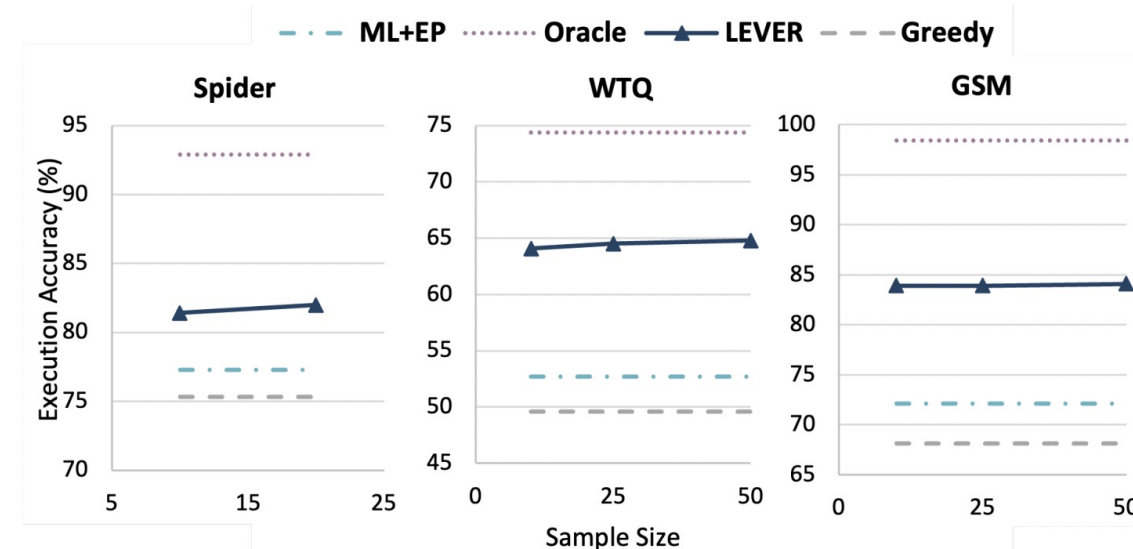
- LEVER is sensitive to the sample size at inference time but not training time thus a higher sampling budget should be applied during inference
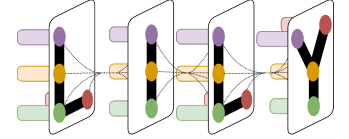


Ablation on sample size at inference time



Ablation on sample size at training time

- LEVER still non-trivially improves the baseline performance in most cases;
- Transfer typically works better when the percentage of positive labels are closer

| Target CodeLM & ML+EP Baseline | | Source CodeLM (% Positive Labels) | | |
|---|---|---|---|---|
| | | Codex (64.0%) | InCoder (9.2%) | CodeGen (8.6%) |
| Codex | 77.3 | 82.0 (+4.7) | 81.7 (+4.4) | 80.8 (+3.5) |
| InCoder | 41.2 | 46.4 (+5.2) | 54.1 (+12.9) | 47.6 (+6.4) |
| CodeGen | 37.7 | 44.7 (+7.0) | 48.9 (+11.2) | 51.0 (+13.3) |

Between-LM Transfer results for Spider

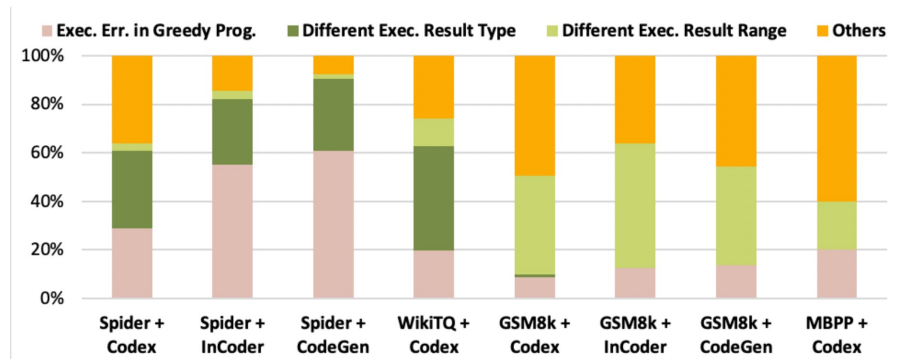| Target CodeLM & ML+EP Baseline | | Source CodeLM (% Positive Labels) | | |
|---|---|---|---|---|
| | | Codex (53.4%) | InCoder (2.3%) | CodeGen (5.0%) |
| Codex | 72.1 | 83.7 (+11.6) | 70.0 (-2.1) | 71.9 (-0.2) |
| InCoder | 4.3 | 8.3 (+4.0) | 11.9 (+7.6) | 12.3 (+8.0) |
| CodeGen | 9.6 | 18.4 (+8.8) | 20.7 (+11.1) | 22.1 (+12.5) |

Between-LM Transfer results for GSM8k
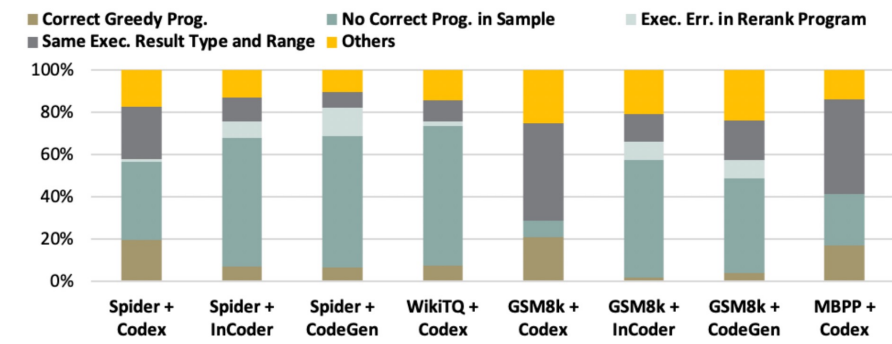
- **When LEVER succeeds:**
  - It is often because the execution results provide crucial information such as execution errors, variable type and range

- **When LEVER fails:**
  - The most common reason is that no correct program can be found in the samples (i.e., upper-bound is reached), which is especially the case for weaker CodeLMs
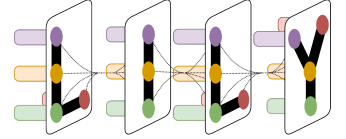


When LEVER reranks a correct program at the top but the greedy decoding fails.



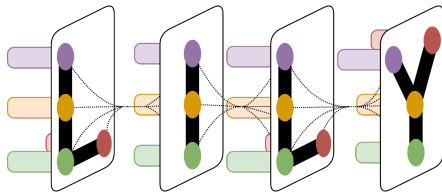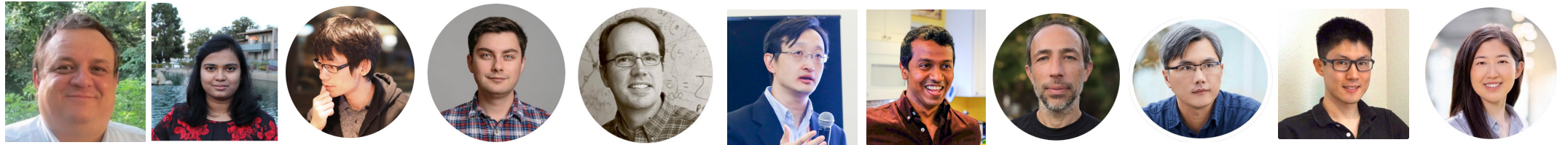When LEVER fails to rank a correct program at the top.

- How can we use a smaller model to help improve LLMs?

  - One way is to train a separate model that operates on the output of LLMs, such as verification, reranking, etc.

- We can also try to incorporate additional information (e.g., execution) in the separate model with the blackbox LLMs

- Neural-symbolic-neural approach is possible!

# Thanks!

**Also thanks to my wonderful collaborators:**

Yale LILY Lab

Meta AI

Microsoft Research

Yale