

# Evaluation: Metrics and Benchmarks

Daniel Fried

11-891: Neural Code Generation

<https://cmu-codegen.github.io/s2024/>



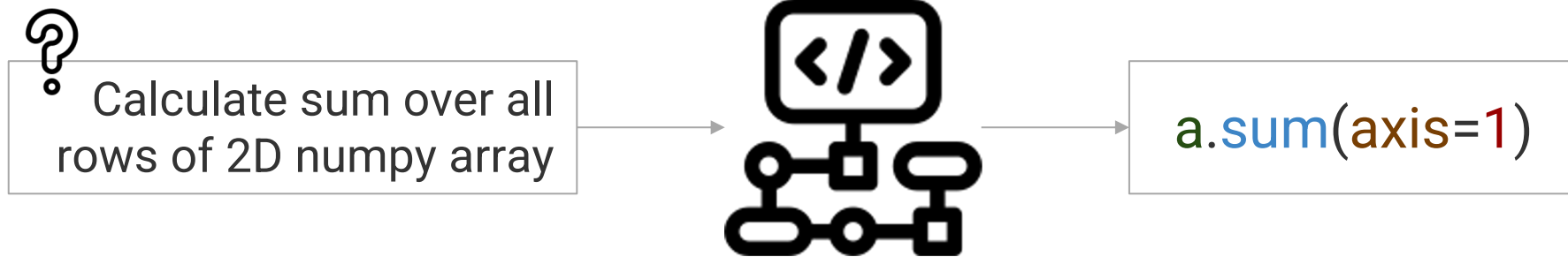
Language  
Technologies  
Institute

With slides from Zora Wang and Nikitha Rao

# The NL2Code Task

---

- ▶ Given a natural language instruction Q, generate code implementation C



# The Landscape for NL2Code Generation

## ▶ Transition of Evaluation Metrics:

- ▷ Lexical
- ▷ Neural based metrics
- ▷ Test case execution

## ▶ Domain Coverage

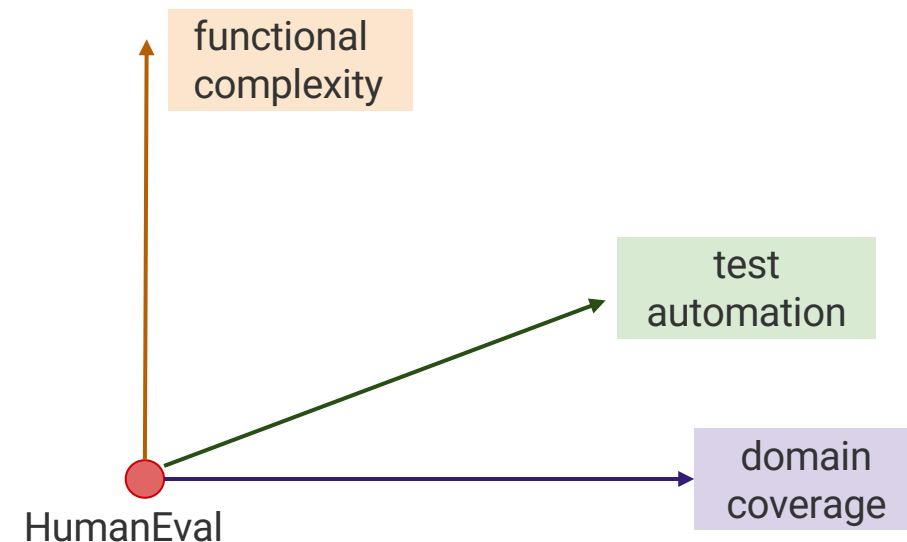
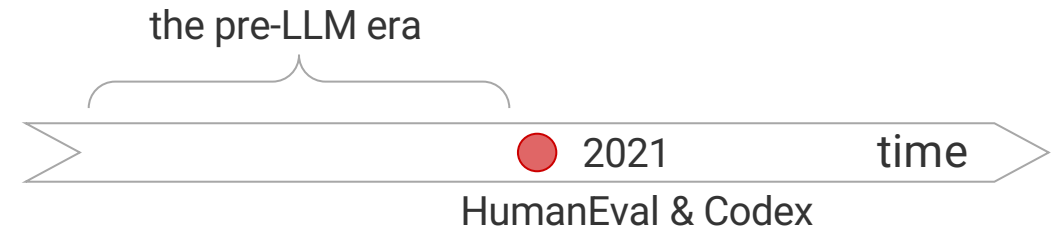
- ▷ Built-in grammar: `sum([1, 2, 4])`
- ▷ Domain-specific: data science
- ▷ Open domain: diverse Python libraries

## ▶ Functional Complexity

- ▷ Simple (toy) functions: e.g., LeetCode
- ▷ Class level
- ▷ Repository level

## ▶ Test Automation

- ▷ Human-written tests
- ▷ Fuzzing methods
- ▷ Integrating LLMs



# Pre-2020

- ▶ Most code snippets were short, and evaluated using BLEU or exact match.
- ▶ Datasets were fairly large, with dedicated training sets.

| Natural Language  | Bash Command(s)   |
|---|---|
| <i>find .java files in the current directory tree that contain the pattern 'TODO' and print their names</i> | <pre>grep -l "TODO" *.java find . -name "*.java" -exec grep -il "TODO" {} \; find . -name "*.java"   xargs -I {} grep -l "TODO" {}</pre>                                      |
| <i>display the 5 largest files in the current directory and its sub-directories</i>                         | <pre>find . -type f   sort -nk 5,5   tail -5 du -a .   sort -rh   head -n5 find . -type f -printf '%s %p\n'   sort -rn   head -n5</pre>                                       |
| <i>search for all jpg images on the system and archive them to tar ball "images.tar"</i>                    | <pre>tar -cvf images.tar \$(find / -type f -name *.jpg) tar -rvf images.tar \$(find / -type f -name *.jpg) find / -type f -name "*.jpg" -exec tar -cvf images.tar {} \;</pre> |

|              | Train | Dev | Test |
|--------------|-------|-----|------|
| # pairs      | 8,090 | 609 | 606  |
| # unique nls | 7,340 | 549 | 547  |

# Pre-2020

- ▶ Most code snippets were short, and evaluated using BLEU or exact match.
- ▶ Datasets were fairly large, with dedicated training sets.

| Dataset     | PL     | # pairs | # words | # tokens | Avg. # w. in nl | Avg. # t. in code | NL collection             | Code collection         | Semantic alignment | Introduced by                |
|-------------|--------|---------|---------|----------|-----------------|-------------------|---------------------------|-------------------------|--------------------|------------------------------|
| IFTTT       | DSL    | 86,960  | –       | –        | 7.0             | 21.8              | scraped                   | scraped                 | Noisy              | (Quirk et al., 2015)         |
| C#2NL*      | C#     | 66,015  | 24,857  | 91,156   | 12              | 38                |                           |                         |                    | (Iyer et al., 2016)          |
| SQL2NL*     | SQL    | 32,337  | 10,086  | 1,287    | 9               | 46                |                           |                         |                    | (Zhong et al., 2018)         |
| RegexLib    | Regex  | 3,619   | 13,491  | 179*     | 36.4            | 58.8*             |                           |                         | Good*              | (Ling et al., 2016)          |
| HeartStone  | Python | 665     | –       | –        | 7               | 352*              | game card                 | game card               |                    |                              |
| MTG         | Java   | 13,297  | –       | –        | 21              | 1,080*            | description               | source code             |                    |                              |
| StaQC       | Python | 147,546 | 17,635  | 137,123  | 9               | 86                | extracted using ML        | extracted using ML      | Noisy              | (Yao et al., 2018)           |
|             | SQL    | 119,519 | 9,920   | 21,413   | 9               | 60                |                           |                         |                    |                              |
| NL2RX       | Regex  | 10,000  | 560     | 45*†     | 10.6            | 26*               | synthesized & paraphrased | synthesized             | Very Good          | (Locascio et al., 2016)      |
| WikiSQL     | SQL    | 80,654  | –       | –        | –               | –                 |                           |                         |                    | (Zhong et al., 2017)         |
| NLMAPS      | DSL    | 2,380   | 1,014   | –        | 10.9            | 16.0              | synthesized given code    | expert written          | Very Good          | (Haas and Riezler, 2016)     |
| Jobs640*    | DSL    | 640     | 391     | 58†      | 9.8             | 22.9              | user written              | expert written given NL |                    | (Tang and Mooney, 2001)      |
| GEO880      | DSL    | 880     | 284     | 60†      | 7.6             | 19.1              |                           |                         |                    | (Zelle and Mooney, 1996)     |
| Freebase917 | DSL    | 917     | –       | –        | –               | –                 |                           |                         |                    | (Cai and Yates, 2013)        |
| ATIS*       | DSL    | 5,410   | 936     | 176†     | 11.1            | 28.1              |                           |                         |                    | (Dahl et al., 1994)          |
| WebQSP      | DSL    | 4,737   | –       | –        | –               | –                 |                           |                         |                    | (Yih et al., 2016)           |
| NL2RX-KB13  | Regex  | 824     | 715     | 85*†     | 7.1             | 19.0*             | turker written            |                         |                    | (Kushman and Barzilay, 2013) |
| Django*     | Python | 18,805  | –       | –        | 14.3            | –                 | expert written            | scraped                 |                    | (Oda et al., 2015)           |
| NL2Bash     | Bash   | 9,305   | 7,790   | 6,234    | 11.7            | 7.7               | given code                |                         |                    | Ours                         |

# Evaluation Metrics

# Reference Matching: BLEU

---

- ▶ Developed for machine translation (Papineni et al. 2002)
- ▶ Compares n-gram overlap between predicted and reference
- ▶ Typically, uses n-grams up to 4 (BLEU-4)

Reference: Taro visited Hanako

System: the Taro visited the Hanako

1-gram: 3/5

2-gram: 1/4

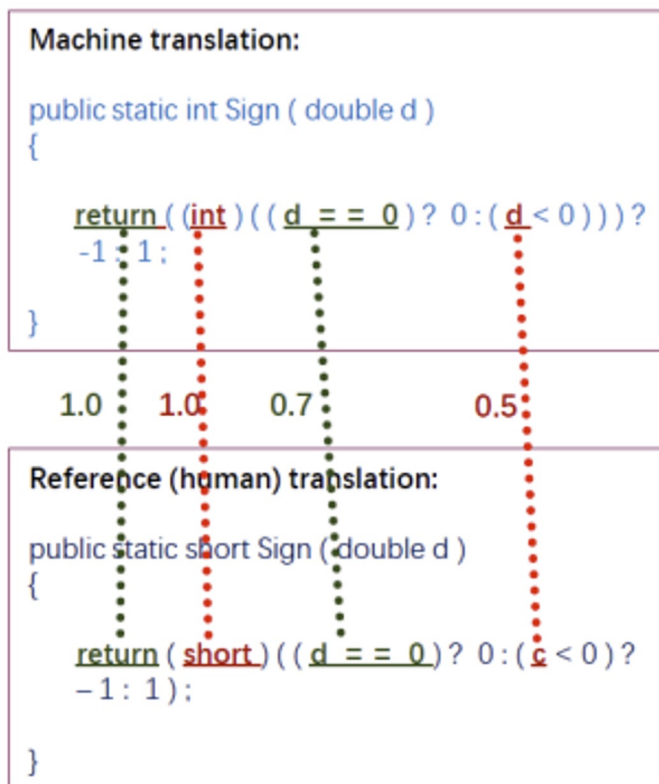
Brevity:  $\min(1, |\text{System}|/|\text{Reference}|) = \min(1, 5/3)$

brevity penalty = 1.0

$$\text{BLEU-2} = (3/5 * 1/4)^{1/2} * 1.0 \\ = 0.387$$

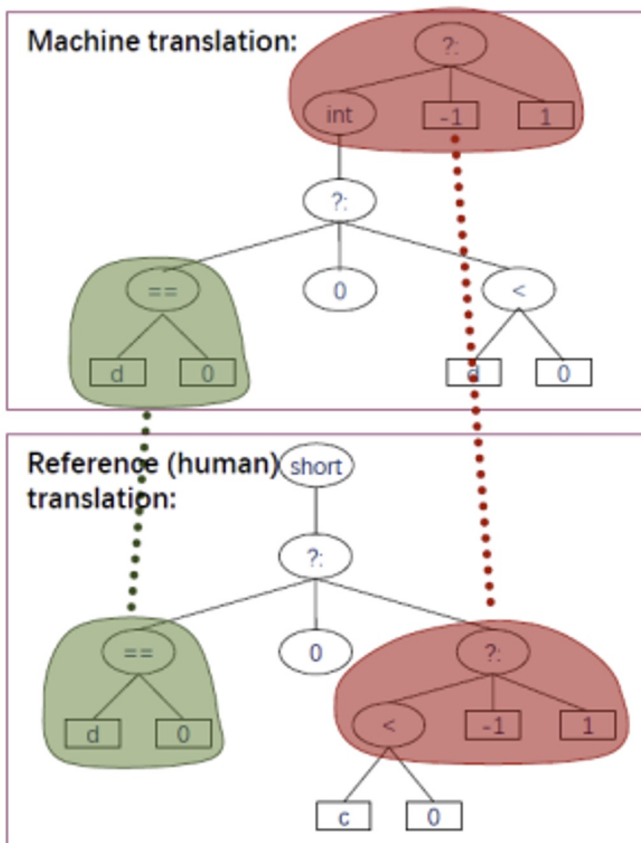
# Reference Matching: CodeBLEU

Higher weight for keywords



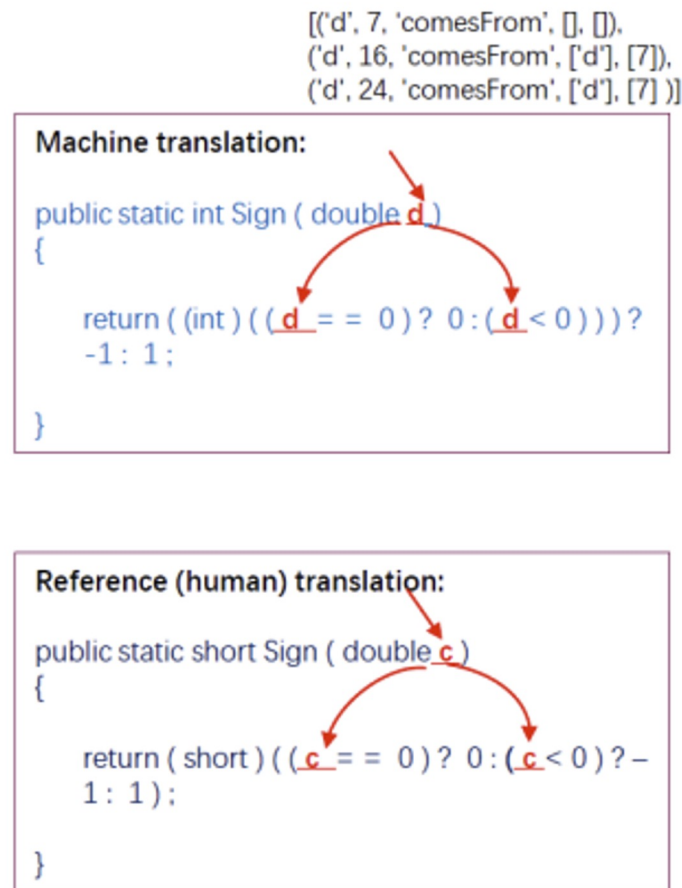
Weighted N-Gram Match

Match syntactic subtrees



Syntactic AST Match

Match data dependency graphs



Semantic Data-flow Match

$$\text{CodeBLEU} = \alpha \cdot \text{N-Gram Match (BLEU)} + \beta \cdot \text{Weighted N-Gram Match} + \gamma \cdot \text{Syntactic AST Match} + \delta \cdot \text{Semantic Data-flow Match}$$



# Reference Matching: CodeBLEU

- ▶ When evaluating evaluation metrics, check correlation with human judgements.
- ▶ In CodeBLEU: rate code outputs on a Likert scale of general quality (1=very bad; 5=very good)

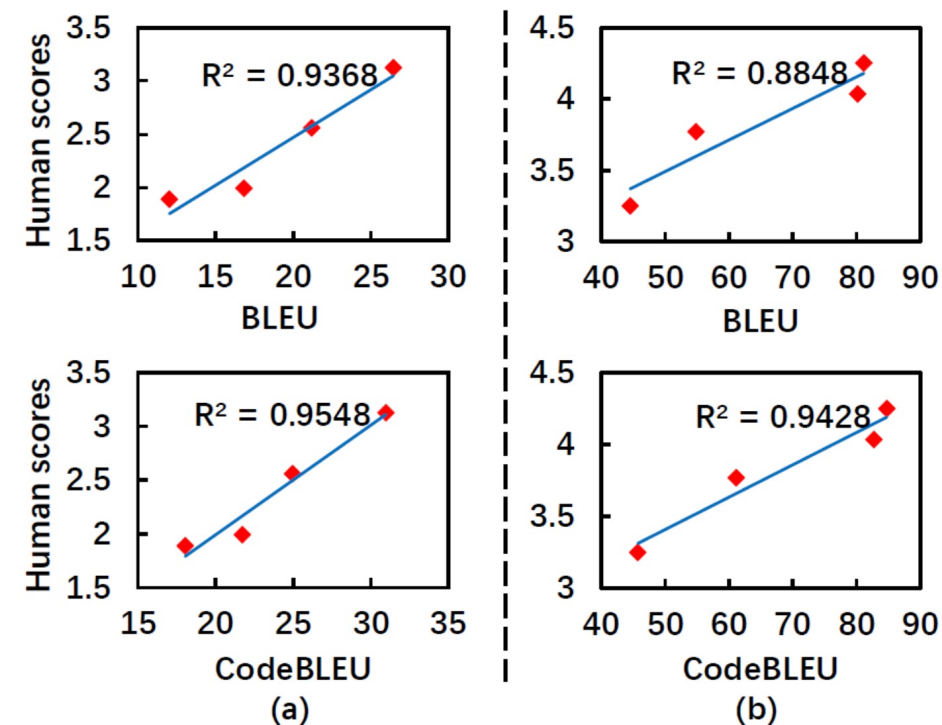
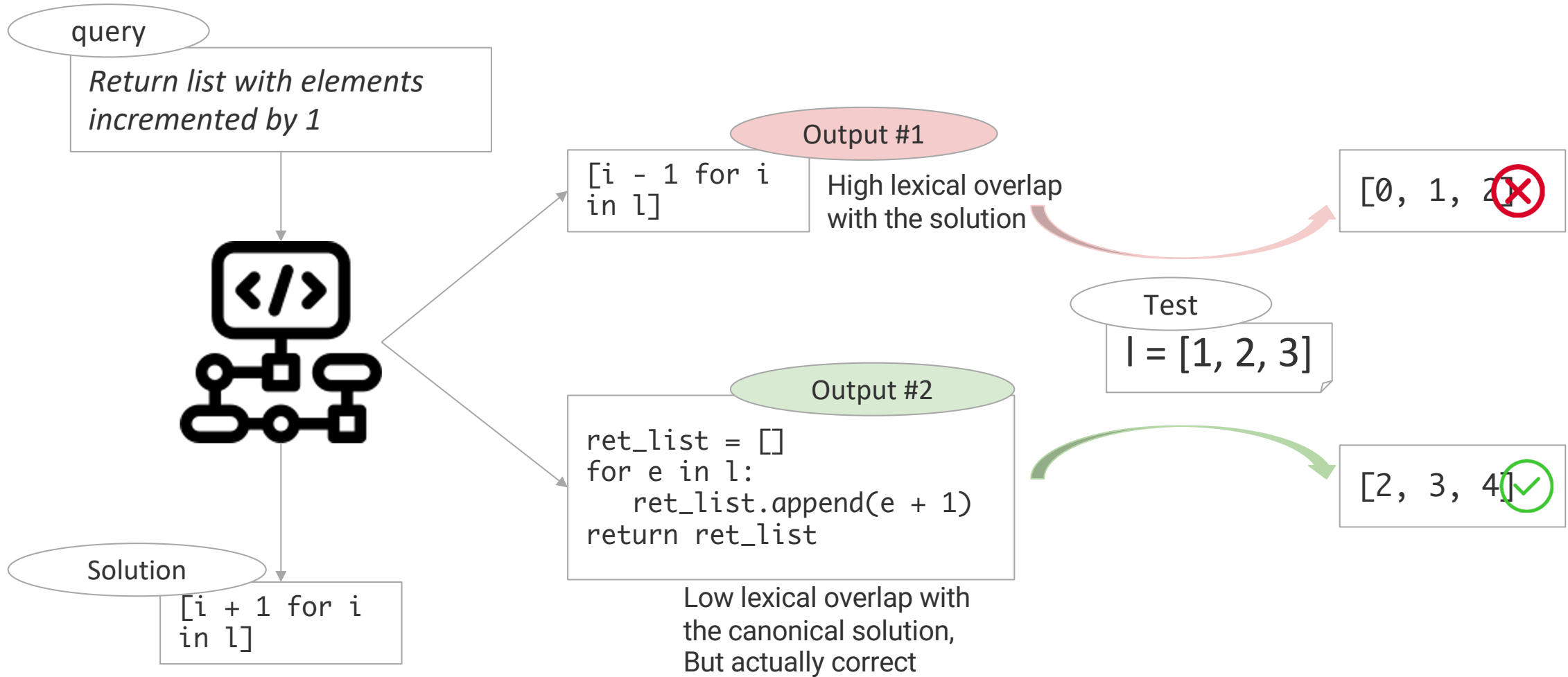


Figure 5: BLEU and CodeBLEU predict human evaluation scores. (a) Text-to-code; (b) Code translation.

# Issues: Evaluations Are Not Rigorous



# HumanEval Benchmark

---

- ▶ Evaluation: test case execution
- ▶ 164 hand-written examples, by authors of the paper
- ▶ Why human-written?
  - ▷ “It is important for these tasks to be hand-written, since our models are trained on a large fraction of GitHub, which already contains solutions to problems from a variety of sources.”

```
def solution(lst):  
    """Given a non-empty list of integers, return the sum of all of the odd elements  
    that are in even positions.  
  
    Examples  
    solution([5, 8, 7, 1]) ==>12  
    solution([3, 3, 3, 3, 3]) ==>9  
    solution([30, 13, 24, 321]) ==>0  
    """  
    return sum(lst[i] for i in range(0, len(lst)) if i % 2 == 0 and lst[i] % 2 == 1)
```

# MBPP: Mostly Basic Python Programs

---

- ▶ Similar to HumanEval, but a bit easier
- ▶ 974 short Python problems, written by crowdworkers
  - ▷ 58% mathematical, 43% list processing, 19% string processing

# MBPP: Mostly Basic Python Programs

- ▶ Model performance is sensitive to sampling temperature and number of candidates (similar findings in HumanEval/Codex paper)

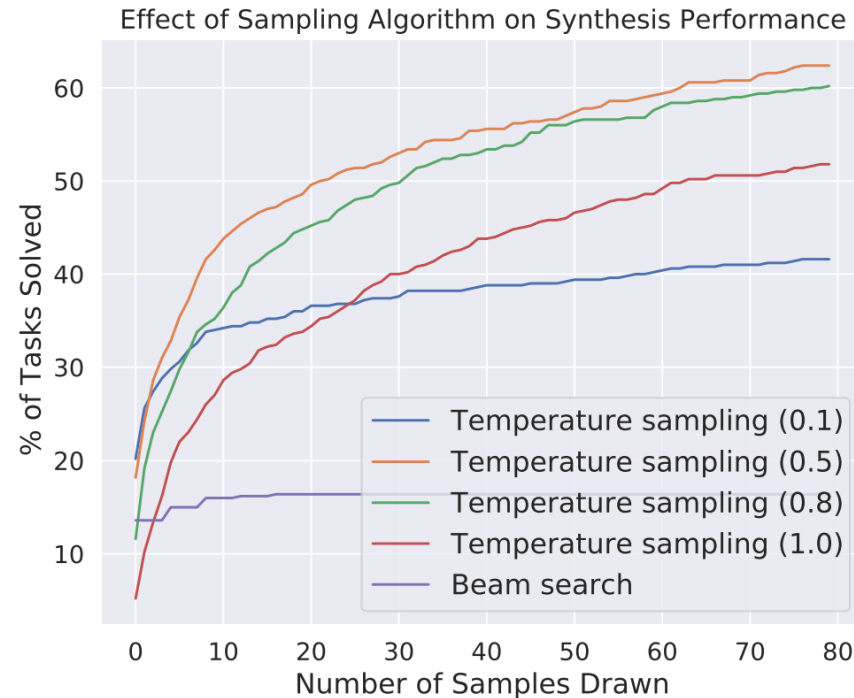


Figure 9: Higher temperatures achieve better scaling with more samples, but perform worse with a smaller budget.

# MBPP: Mostly Basic Python Programs

- ▶ BLEU against a reference solution is uncorrelated with whether samples pass execution tests (similar findings in Codex paper).

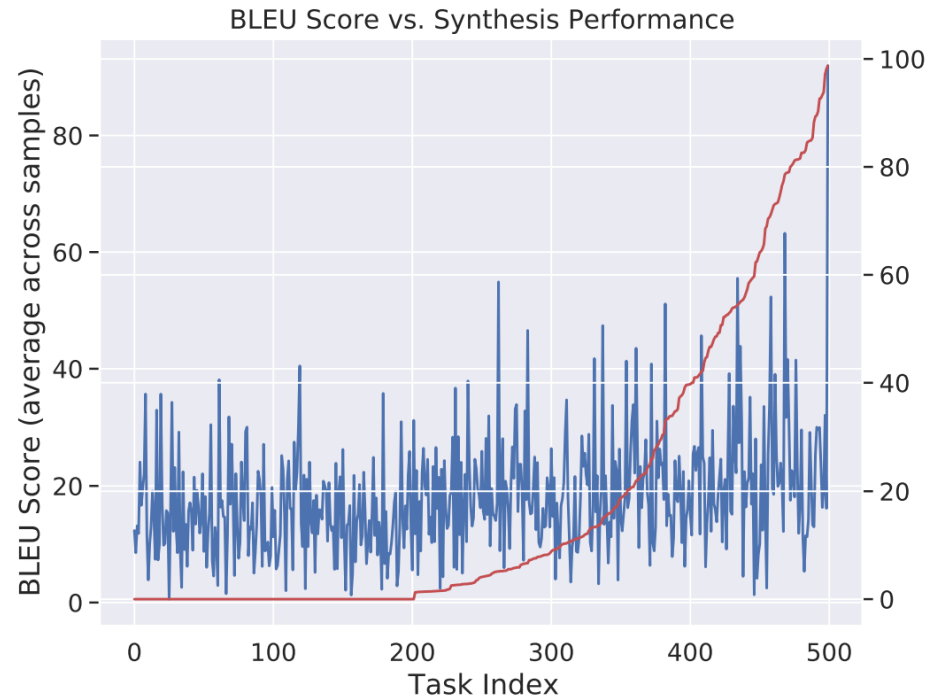


Figure 10: Comparison of BLEU score and synthesis performance for the 137B parameter model. No strong correlation is observed.

# MBPP: Mostly Basic Python Programs

- ▶ Model evaluated is a large Google LLM, LaMDA, trained mostly on natural language, which has some interaction ability.

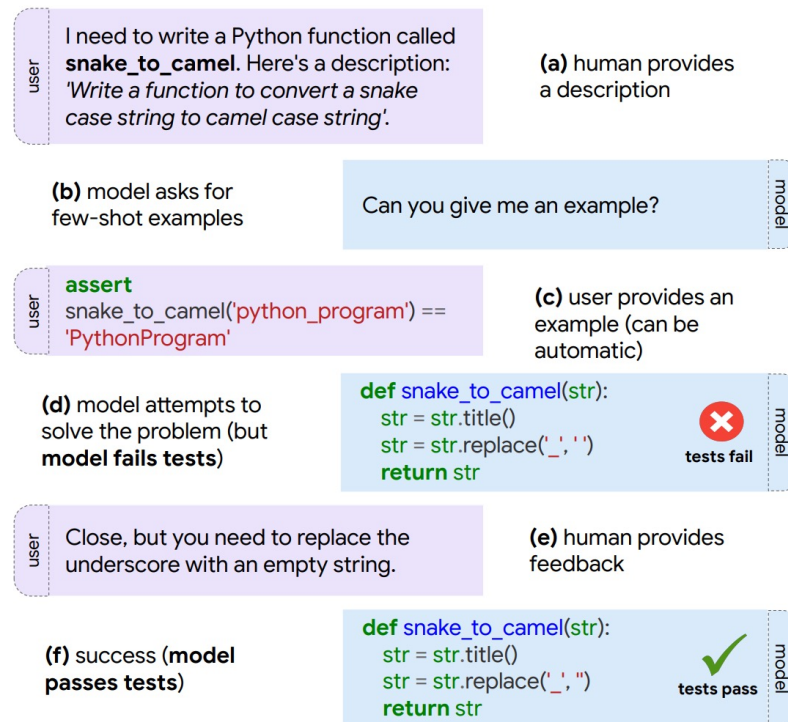


Figure 12: An overview of the “flow” of the human-model collaboration experiments. The human gives a description of the desired program and then guides the model toward the correct solution via dialog.

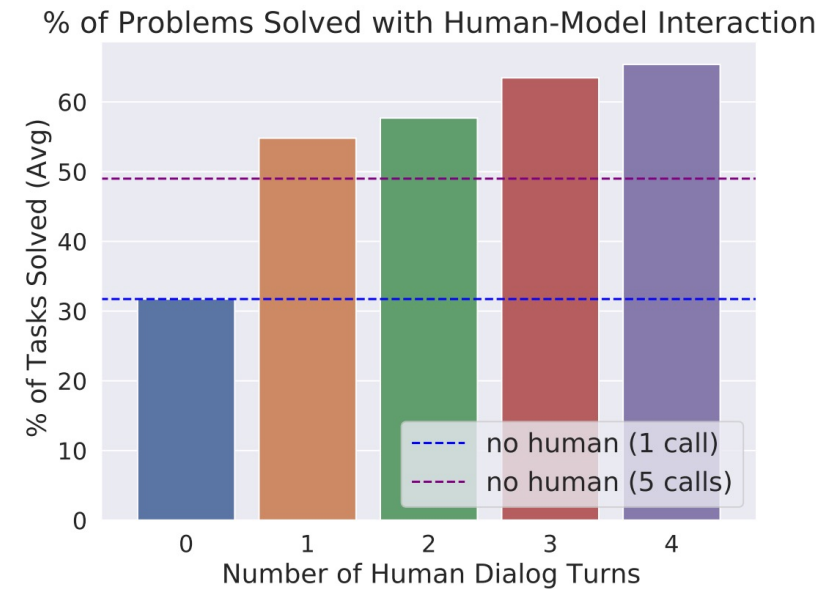
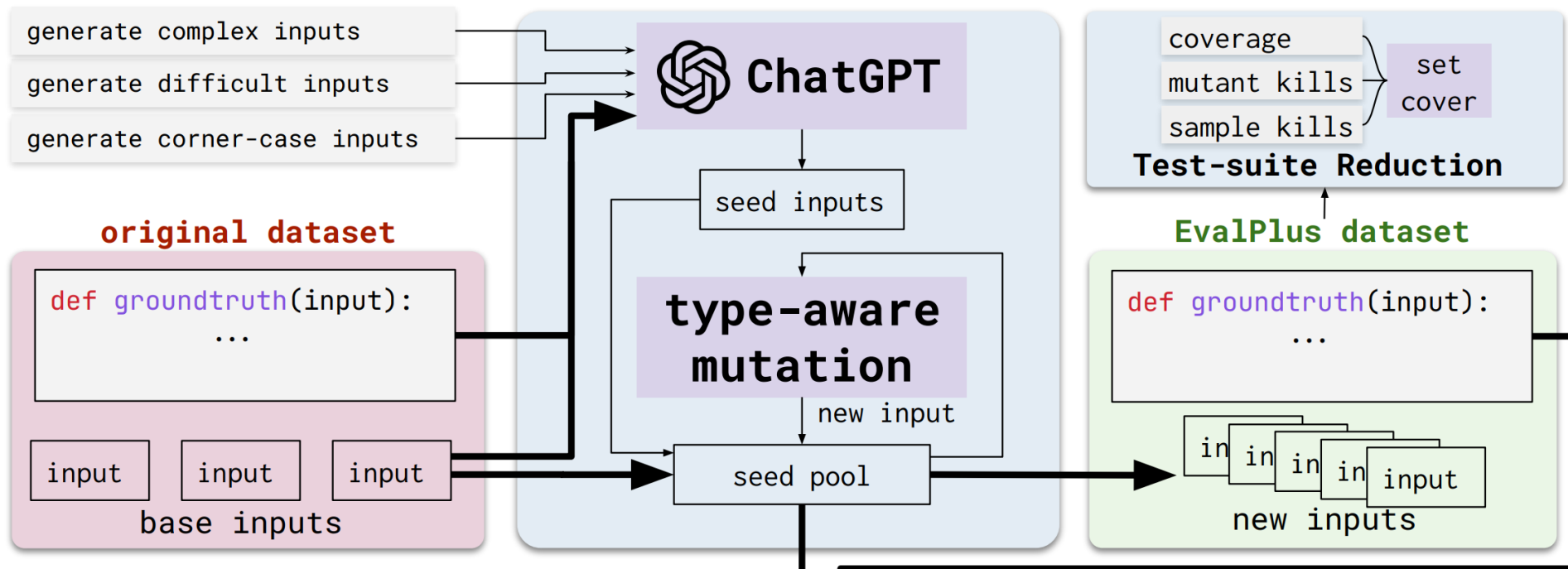


Figure 13: Percent of problems solved as the number of human dialog interventions increases. With 4 interventions, the solve rate increases from 30% to over 65%. Except for the purple horizontal baseline (which corresponds to 5 samples from the model), all pass-rates in this figure were computed using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

# Automated & Improved Testing

- ▶ EvalPlus: use LLMs and *fuzzing* (type-aware mutation) to create test cases
- ▶ Prompt ChatGPT with the GT solution, some inputs, and instructions to generate more





# Automated & Improved Testing

- ▶ EvalPlus: use LLMs and *fuzzing* (type-aware mutation) to create test cases
- ▶ Fuzzing: mutate inputs to the functions, apply the groundtruth function, and use the input-output pair to make a new test case.

Table 1: List of basic type-aware mutations over input  $x$ .

| Type                                  | Mutation  | Type               | Mutation  |
|---------------------------------------|---|--------------------|---|
| <code>int</code>   <code>float</code> | Returns $x \pm 1$   | <code>List</code>  | $\left\{ \begin{array}{l} \text{Remove/repeat a random item } x[i] \\ \text{Insert/replace } x[i] \text{ with } \text{Mutate}(x[i]) \end{array} \right.$  |
| <code>bool</code>                     | Returns a random boolean  | <code>Tuple</code> | Returns <code>Tuple(Mutate(List(x)))</code>   |
| <code>NoneType</code>                 | Returns None  | <code>Set</code>   | Returns <code>Set(Mutate(List(x)))</code>   |
| <code>str</code>                      | $\left\{ \begin{array}{l} \text{Remove a sub-string } s \\ \text{Repeat a sub-string } s \\ \text{Replace } s \text{ with } \text{Mutate}(s) \end{array} \right.$ | <code>Dict</code>  | $\left\{ \begin{array}{l} \text{Remove a key-value pair } k \rightarrow v \\ \text{Update } k \rightarrow v \text{ to } k \rightarrow \text{Mutate}(v) \\ \text{Insert } \text{Mutate}(k) \rightarrow \text{Mutate}(v) \end{array} \right.$ |

# Automated & Improved Testing

- ▶ EvalPlus: use LLMs and *fuzzing* (type-aware mutation) to create test cases
- ▶ Optionally, minify the test sets while preserving code coverage and edge case detection.

Table 2: Overview of EvalPlus-improved benchmarks.

|                              | #Tests |        |      |                  | #Tasks |
|------------------------------|--------|--------|------|------------------|--------|
|                              | Avg.   | Medium | Min. | Max.             |        |
| HUMANEVAL                    | 9.6    | 7.0    | 1    | 105 <sup>2</sup> |        |
| HUMANEVAL <sup>+</sup>       | 764.1  | 982.5  | 12   | 1,100            | 164    |
| HUMANEVAL <sup>+</sup> -MINI | 16.1   | 13.0   | 5    | 110              |        |

|                            | Size | pass@k | k=1* | k=1  | k=10 | k=100 |
|----------------------------|------|--------|------|------|------|-------|
| GPT-4 [49]                 | N/A  | base   | 88.4 |      |      |       |
|                            |      | +extra | 76.2 |      |      |       |
| Phind-CodeLlama [52]       | 34B  | base   | 71.3 | 71.6 | 90.5 | 96.2  |
|                            |      | +extra | 67.1 | 67.0 | 85.0 | 92.5  |
| WizardCoder-CodeLlama [38] | 34B  | base   | 73.2 | 61.6 | 85.2 | 94.5  |
|                            |      | +extra | 64.6 | 54.5 | 78.6 | 88.9  |
| ChatGPT [48]               | N/A  | base   | 73.2 | 69.4 | 88.6 | 94.0  |
|                            |      | +extra | 63.4 | 62.5 | 82.1 | 91.1  |

# Automated & Improved Testing

- ▶ EvalPlus: use LLMs and *fuzzing* (type-aware mutation) to create test cases
- ▶ These extra tests substantially reduce the pass@k of many models!

|                            | Size | pass@k | $k=1^*$ | $k=1$ | $k=10$ | $k=100$ |
|----------------------------|------|--------|---------|-------|--------|---------|
| GPT-4 [49]                 | N/A  | base   | 88.4    |       |        |         |
|                            |      | +extra | 76.2    |       |        |         |
| Phind-CodeLlama [52]       | 34B  | base   | 71.3    | 71.6  | 90.5   | 96.2    |
|                            |      | +extra | 67.1    | 67.0  | 85.0   | 92.5    |
| WizardCoder-CodeLlama [38] | 34B  | base   | 73.2    | 61.6  | 85.2   | 94.5    |
|                            |      | +extra | 64.6    | 54.5  | 78.6   | 88.9    |
| ChatGPT [48]               | N/A  | base   | 73.2    | 69.4  | 88.6   | 94.0    |
|                            |      | +extra | 63.4    | 62.5  | 82.1   | 91.1    |

# MultiPL-E

---

- ▶ Key idea: it's relatively easy to translate test cases on simple types (e.g. no matrices or functions) from Python to other languages.
- ▶ This allows porting HumanEval & MBPP to 18 other languages.

(a) Original Python assertion.

```
assert lsi([0]) == (None, None)
```

---

(b) Equivalent R.

```
if(!identical(lsi(c(0)), c(NULL, NULL))){  
  quit('no', 1)}
```

---

(c) Equivalent JavaScript.

```
assert.deepEqual(lsi([0]), [void 0, void 0]);
```

Figure 4: Example of a translated assertion.

(a) Original Python docstring from HumanEval #95.

Given a `dictionary`, return `True` if all keys are strings in lower case or all keys are strings in upper case, else return `False`. The function should return `False` if the given `dictionary` is empty.

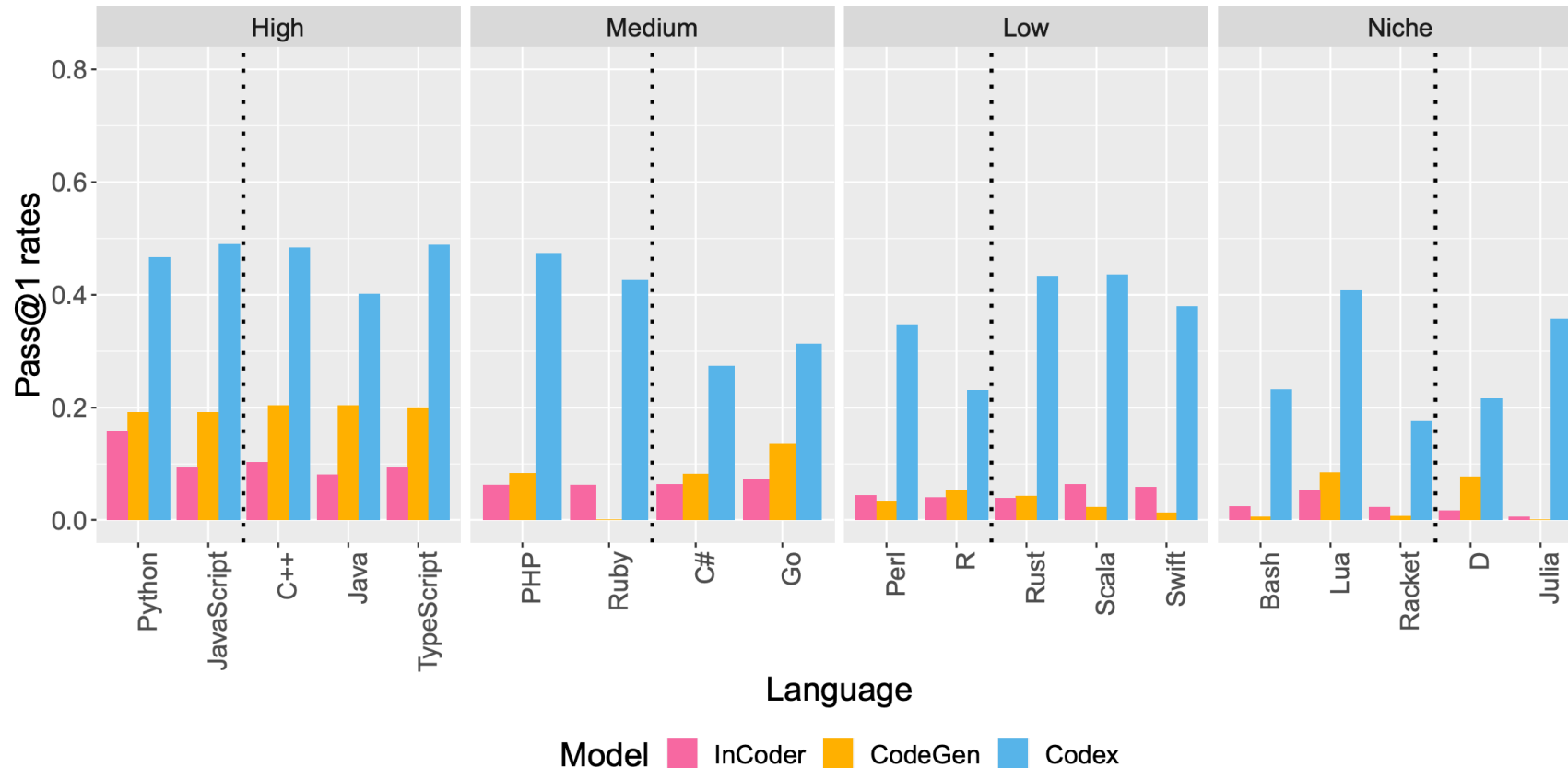
(b) Terminology translated to Perl.

Given a `hash`, return `1` if all keys are strings in lower case or all keys are strings in upper case, else return `""`. The function should return `""` if the given `hash` is empty.

Figure 5: A Python docstring and its Perl translation. Errors (e.g., "is" for "if") are from the original benchmark.

# MultiPL-E

- ▶ Models are generally better on “high-resource” languages with more code on GitHub.
- ▶ More analysis of this in the Data lecture, with Starcoder.



# Incorrect Code Can Be Valuable Too!

- ▶ Code might be easily editable to achieve a good solution.

Levenshtein distance: number of character edits required to transform.

$$\text{EDIT-SIM} = 1 - \frac{\text{lev}(\text{gen}, \text{ref})}{\max(\text{len}(\text{gen}), \text{len}(\text{ref}))}$$

Reference Code Snippet

```
def even_odd_count(num):
    even_count = 0
    odd_count = 0
    for i in str(abs(num)):
        if int(i)%2==0:
            even_count +=1
        else:
            odd_count +=1
    return (even_count, odd_count)
```

Functional Metric

pass = 0

Generated Code Snippet

```
def even_odd_count(num):
    even_count = 0
    odd_count = 0
    for i in str(num):
        if int(i) % 2 == 0:
            even_count += 1
        else:
            odd_count += 1
    return even_count, odd_count
```

Similarity Metric

edit similarity = 0.93

Human preference

preference = 0.9

# Incorrect Code Can Be Valuable Too!

---

- ▶ Dibia et al. compare metrics to evaluate 5 model outputs on HumanEval.
  - ▷ EditDistance, BLEU, Pass@1
- ▶ Professional programmers with Python experience rate on:
  - ▷ **Accuracy**: judge if the snippets are functionally equivalent (judging is easier than writing!)
  - ▷ **Value**: How useful is the snippet as a starting point?
  - ▷ **Effort**: how much effort to modify the solution into a correct one?

# Incorrect Code Can Be Valuable Too!

- ▶ Value is nearly perfectly correlated with effort (accuracy less so).
- ▶ Of all metrics, Pass@1 is most correlated with Value
- ▶ But, Edit sim > BLEU and a combination is best (as dissimilar, incorrect code is bad).

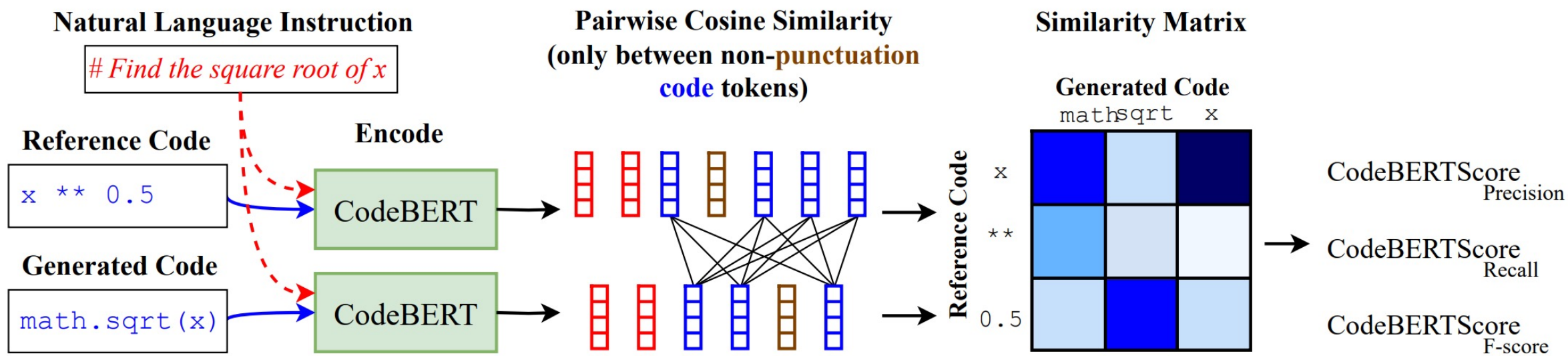
$$\mathbf{COMBINED} = \min(1.0, \text{PASS} + \text{EDIT-SIM})$$

|                 | Human Judgement |          |        | Offline Metrics |          |      |                 |
|-----------------|-----------------|----------|--------|-----------------|----------|------|-----------------|
|                 | Value           | Accuracy | Effort | Pass            | Edit Sim | bleu | <b>Combined</b> |
| Value           | 1.00            |          |        |                 |          |      |                 |
| Accuracy        | 0.87            | 1.00     |        |                 |          |      |                 |
| Effort          | 0.94            | 0.86     | 1.00   |                 |          |      |                 |
| Pass            | <b>0.61</b>     | 0.66     | 0.62   | 1.00            |          |      |                 |
| Edit Sim        | <b>0.48</b>     | 0.46     | 0.51   | 0.33            | 1.00     |      |                 |
| bleu            | <b>0.36</b>     | 0.34     | 0.39   | 0.19            | 0.68     | 1.00 |                 |
| <b>Combined</b> | <b>0.70</b>     | 0.71     | 0.72   | 0.89            | 0.61     | 0.38 | 1.00            |



# CodeBERTScore: Model-based Evaluation

- ▶ Captures some intuitions about incorrect code being useful
- ▶ BLEU and edit distance only give points for exactly matching code
- ▶ Takes NL code descriptions into account
- ▶ Use vector similarity from CodeBERT representations
- ▶ Recall: every reference vector has  $\geq 1$  candidate vector with high similarity
- ▶ Precision: every candidate vector has  $\geq 1$  reference vector with high similarity



# Domains of Code

# HumanEval Looks Like Toy Examples?

## ▶ HumanEval Examples

```
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]
```

```
def solution(lst):  
    """Given a non-empty list of integers, return the sum of all of the odd elements  
    that are in even positions.  
  
    Examples  
    solution([5, 8, 7, 1]) ==>12  
    solution([3, 3, 3, 3, 3]) ==>9  
    solution([30, 13, 24, 321]) ==>0  
    """  
    return sum(lst[i] for i in range(0, len(lst)) if i % 2 == 0 and lst[i] % 2 == 1)
```

## Real-World Development

### Asking the user for input until they give a valid response

Asked 9 years, 6 months ago Modified 1 year, 5 months ago Viewed 1.0m times

▲ I am writing a program that accepts user input.

750

```
#note: Python 2.7 users should use `raw_input`, the equivalent of 3.X's `input`  
age = int(input("Please enter your age: "))  
if age >= 18:  
    print("You are able to vote in the United States!")  
else:  
    print("You are not able to vote in the United States.")
```

transformers Public

Watch 1.1k Fork 22.9k Starred 114k

main 262 branches 139 tags

Go to file Add file Code

About

Transformers: State-of-the-art Machine Learning for Pytorch, TensorFlow, and JAX.

huggingface.co/transformers

python nlp machine-learning natural-language-processing deep-learning tensorflow pytorch transformer speech-recognition seq2seq flax pretrained-models language-models nlp-library language-model hactoberfest bert jax pytorch-transformers model-hub

|             |  |         |            |                |
|-------------|--|---------|------------|----------------|
| hi-sushanta | Removed the redundant SILUActivation class. (#27136) | 4991216 | 1 hour ago | 14,388 commits |
| .circleci   | Limit to inferior fsspec version (#27010)            |         |            | last week      |
| .github     | Dev version  |         |            | 2 hours ago    |
| docker      | [core / Quantization] AWQ integration (#27045)       |         |            | 2 days ago     |
| docs        | translate peft.md to chinese (#27215)                |         |            | 2 hours ago    |
| examples    | Dev version  |         |            | 2 hours ago    |
| model_cards | Update URL for Hub PR docs (#17532)                  |         |            | last year      |
| notebooks   | Update README.md (#25941)                            |         |            | 2 months ago   |
| scripts     | Fix stale bot for locked issues (#26711)             |         |            | 3 weeks ago    |

# Broadening Domains

- ▶ Leetcode Style: HumanEval, APPS, MBPP
  - ▷ Manually written or collected from code contest websites
  - ▷ Only uses Python built-in grammar
- ▶ Limited Domains: e.g., Data Science
  - ▷ DS-1000: StackOverflow questions
  - ▷ ARCADE: Interactive Jupyter Notebooks
  - ▷ ... ..
- ▶ Open Domain: ODEX
  - ▷ 79 Python libraries
  - ▷ Four natural languages

NumPy   SciPy   Pandas   TensorFlow   PyTorch   Scikit-learn   Matplotlib

**Problem:**  
I could not find a built-in function in Python to generate a log uniform distribution given a min and max value (the R equivalent is here), something like: `loguni[n, min, max, base]` that returns `n` log uniformly distributed in the range `min` and `max`.  
The closest I found though was `numpy.random.uniform`.  
That is, given range of `x`, I want to get samples of given size (`n`) that suit log-uniform distribution.  
Any help would be appreciated!

**A:**  
<code>  
`import numpy as np`  
`min = 1`  
`max = np.e`  
`n = 10000`  
</code>  
BEGIN SOLUTION  
<code>  
[insert]  
</code>  
END SOLUTION  
<code>  
`print(result)`  
</code>

**Reference Solution**  
`import scipy.stats`  
`result = scipy.stats.loguniform.rvs(a = min, b = max, size = n)`

**Automatic Evaluation**  
**Test case 1**  
`min = 1`  
`max = np.e`  
`n = 10000`  
`ans = ... # generated by Reference solution`  
**Test code**  
`np.testing.assert_array_equal(result.shape, ans.shape)`  
`from scipy.stats import ks_2samp`  
`# Kolmogorov-Smirnov Test judges whether the two sampled`  
`# from similar distribution`  
`assert ks_2samp(result, ans)[0] <= 0.1`  
**Surface form constraints**  
`for` and `while` should not appear in Syntax Tree

NumPy example problem involving randomness, requiring the use of a specialist knowledge test.

# APPS

- ▶ 10,000 problems taken from contest and exercise sites
- ▶ 3 difficulty levels
- ▶ Has test cases

## Problem

You are given two integers  $n$  and  $m$ . Calculate the number of pairs of arrays  $(a, b)$  such that: the length of both arrays is equal to  $m$ ; each element of each array is an integer between 1 and  $n$  (inclusive);  $a_i \leq b_i$  for any index  $i$  from 1 to  $m$ ; array  $a$  is sorted in non-descending order; array  $b$  is sorted in non-ascending order. As the result can be very large, you should print it modulo  $10^9 + 7$ . Input: The only line contains two integers  $n$  and  $m$  ( $1 \leq n \leq 1000$ ,  $1 \leq m \leq 10$ ). Output: Print one integer – the number of arrays  $a$  and  $b$  satisfying the conditions described above modulo  $10^9 + 7$ .

## Model Output

```
n, m = map(int, input().split())
mod = 10**9+7

fact = [1, 1]
factinv = [1, 1]
inv = [0, 1]

for i in range(2, n+1):
    fact.append((fact[-1] * i) % mod)
    inv.append((-inv[mod % i] * (mod//i)) % mod)
    factinv.append((factinv[-1] * inv[-1]) % mod)

def cmb(n, r, mod):
    if (r < 0) or (n < r):
        return 0
    r = min(r, n - r)
    return fact[n] * factinv[r] * factinv[n-r] % mod

print((cmb(n, m, mod) * cmb(m-1, n-1, mod)) % mod)
```

# CodeContests (AlphaCode)

- ▶ Similar to APPS but focuses on improving test coverage, via input mutation (like EvalPlus).

| Split | Problems | Tests per problem |        |           | Solutions per problem (% correct) |             |             |
|-------|----------|-------------------|--------|-----------|-----------------------------------|-------------|-------------|
|       |          | Example           | Hidden | Generated | C++                               | Python      | Java        |
| Train | 13328    | 2.0               | 14.8   | 79.1      | 493.4 (27%)                       | 281.1 (47%) | 147.9 (46%) |
| Valid | 117      | 1.5               | 12.9   | 190.0     | 231.6 (47%)                       | 137.2 (55%) | 131.1 (54%) |
| Test  | 165      | 1.7               | 9.4    | 192.7     | 196.0 (45%)                       | 97.3 (54%)  | 105.2 (51%) |

- ▶ Manual inspection shows high false-positive rate of model-produced solutions.

| Dataset          | Tests / problem | False Positive (FP) Rate | FP or Slow Rate |
|------------------|-----------------|--------------------------|-----------------|
| APPS             | 20.99           | 60%                      | 70%             |
| HumanEval        | 7.77            | 30%                      | N/A             |
| CodeContests raw | 12.4            | 62%                      | 88%             |
| CodeContests     | <b>203.7</b>    | <b>4%</b>                | <b>46%</b>      |

# DS-1000

- ▶ 1,000 data science problems, based on StackOverflow questions
- ▶ Domain-specific test cases, e.g. matplotlib plots have their elements programmatically extracted

## ① Manually Selecting and Modifying StackOverflow Problems

Here is a sample dataframe:

```
6 df = pd.DataFrame({"A": [1, 2, 3], "B": [4, 5, 6]})
```

I'd like to add inverses of each existing column to the dataframe and ... [omitted for brevity]

try:

```
5 inv_df = df.join(df.apply(lambda x: 1/x).add_prefix("inv_"))
```

- ✓ High vote
- ✓ Testable
- ✓ Useful
- ✓ Representative

## ② Adding Code Context

```
import pandas as pd
df = pd.DataFrame({"A": [1, 2, 3],
                  "B": [4, 5, 6]})
### BEGIN SOLUTION
[insert]
### END SOLUTION
print(result)
```

## ③ Implementing Automatic Tests

### Test cases

...[omit for brevity]

```
pd.testing.assert_frame_equal(result,
                               ans)
```

### Surface-form constraints

for and while should not appear in Syntax Tree

## ④ Perturbing Original Problem

... I'd like to apply the **exponential** function to each existing column ... The resulting dataframe should look like so:

```
result = pd.DataFrame({"A": [1, 2, 3],
                      "B": [4, 5, 6],
                      "exp_A": [e^1, e^2, e^3],
                      "exp_B": [e^4, e^5, e^6]})
```

... [omitted for brevity]

## ⑤ Red Teaming

```
df = pd.DataFrame({"A": [1, 2, 3],
                  "B": [4, 5, 6]})
### BEGIN SOLUTION
# A known WRONG SOLUTION
result = df.join(df.apply(lambda x: math.e**x).add_prefix('exp_'))
### END SOLUTION
print(result)
```

# DS-1000

- ▶ Perturb the problems to reduce chances of memorization, since models may have been trained on StackOverflow
- ▶ "Surface" perturbations: don't change solution. "Semantic": do, but try to keep difficulty the same (e.g. max -> min)

|                             | Pandas                    | NumPy                     | Scikit-learn               | SciPy                     | TensorFlow                 | PyTorch                    | Overall                   |
|-----------------------------|---------------------------|---------------------------|----------------------------|---------------------------|----------------------------|----------------------------|---------------------------|
| Origin <sub>surface</sub>   | 37.3                      | 61.2                      | 52.6                       | 33.0                      | 64.9                       | 64.8                       | 53.2                      |
| Surface                     | 31.9 <small>-5.4</small>  | 58.4 <small>-2.8</small>  | 55.7 <small>+3.1</small>   | 32.1 <small>-0.9</small>  | 58.0 <small>-8.9</small>   | 50.0 <small>-14.8</small>  | 49.8 <small>-3.4</small>  |
| Origin <sub>semantic</sub>  | 36.8                      | 56.7                      | 60.6*                      | 40.3                      | 71.3                       | 65.1                       | 47.2                      |
| Semantic                    | 33.2 <small>-3.6</small>  | 49.0 <small>-7.7</small>  | 38.9* <small>-21.7</small> | 34.3 <small>-6.0</small>  | 42.5 <small>-25.8</small>  | 30.5 <small>-34.6</small>  | 38.2 <small>-9.0</small>  |
| Origin <sub>difficult</sub> | 39.9                      | 52.7                      | 5.0*                       | 58.1                      | 73.0*                      | 53.8*                      | 46.8                      |
| Difficult Rewrite           | 17.7 <small>-22.2</small> | 27.1 <small>-25.6</small> | 0.0* <small>-5.0</small>   | 13.8 <small>-44.3</small> | 38.0* <small>-35.0</small> | 28.8* <small>-25.0</small> | 21.0 <small>-25.8</small> |



# ARCADE

## ▶ Executable problems from Jupyter notebooks

```

Prompt Prefix (Exemplars)
Which countries host at least two Olympic games?
# Solution: Let's solve this problem step-by-step. preamble
# Step 1: Get the counts each country hosted Olympics
count_df = df['Country'].value_counts()
# Step 2: Get the rows whose average score is above 90
filtered_df = count_df[count_df >= 2]
# Step 3: Get the country names as a list explanation
filtered_df.index.tolist()

[1] import pandas as pd
df = pd.read_csv('stores.csv')

[2] # Schema of Dataframes:
# Columns in df with example values:
# Stu_Name (Mike),Engineering (90),English (89),Math (92)

[3] Get the students with an average score above 90
for science subjects

[3a] ▶ Vanilla Prediction (no exemplars):
df['Science_Avg'] = (df['Engineering']+df['Math'])/2
df[df['Science_Avg'] > 90][['Stu_Name', 'Science_Avg']]
  
```

| Models  | pass@30 | # API | Lines of Code (LoC) | Comment Lines | Tokens /Line | API /Line |
|---|---------|-------|---------------------|---------------|--------------|-----------|
| Baseline (Tab. 2)                                   | 47.7    | 4.9   | 2.3                 | 0.1           | 21.1         | 3.2       |
| + More Context                                      | 49.3    | 4.9   | 2.3                 | 0             | 21.1         | 3.1       |
| <i>Prompting with Additional Few-shot Exemplars</i> |         |       |                     |               |              |           |
| Vanilla Code  | 49.9    | 5.3   | 2.4                 | 0.1           | 20.8         | 3.1       |
| Step-by-Step Code                                   | 51.9    | 5.6   | 3.2                 | 0.1           | 17.8         | 2.7       |
| + Preamble  | 51.9    | 5.9   | 3.5                 | 0.2           | 16.9         | 2.5       |
| + Pre. + Explanation                                | 52.5    | 6.8   | 4.2                 | 3.3           | 14.9         | 2.2       |

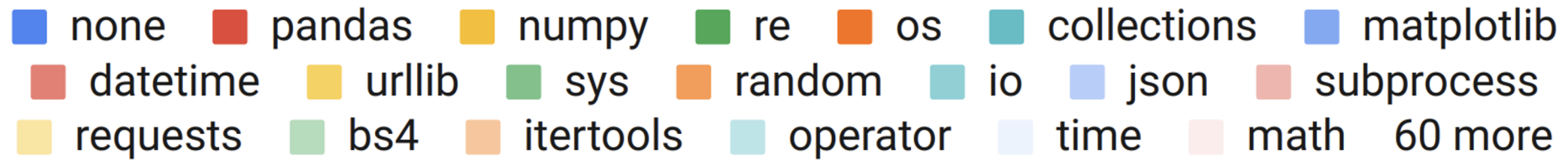
```

[3b] ▶ Step-by-Step Prompting (with exemplars):
# Solution: Let's solve this problem step-by-step. preamble
# Step 1: Create a new column with the average score of
# engineering and math explanation
df['Science_Avg'] = (df['Engineering'] + df['Math']) / 2
# Step 2: Get the rows whose average score is above 90
df_score_above_90 = df[df['Science_Avg'] > 90]
# Step 3: Return the student name and average scores
result = df_score_above_90[['Stu_Name', 'Science_Avg']]
  
```

# ODEX: Open-Domain, with Evaluation

---

## ▶ Larger Domain Coverage



## ▶ Test execution on real-world coding queries

▷ Collected from StackOverflow questions

## ▶ Support four natural languages as input

▷ English, Spanish, Japanese, Russian

```
import requests

def function(files, url, data):
    """multipartのリクエストで複数のデータ`files`, `data`を`url`にPOSTする
    (POST multiple data `files`, `data` to `url` with multipart request)
    """
```

```
return requests.post(url, files=files, data=data)

# test case
r = requests.Response()
r.status_code = 200
requests.post = Mock(return_value = r)
file_path = 'a.txt'
```

| Dataset                       | Samples | Domain | Executable? | Avg. Test Cases | Data Source                   | NL             |
|-------------------------------|---------|--------|-------------|-----------------|-------------------------------|----------------|
| JuICe (Agashe et al., 2019)   | 1,981   | open   | ✗           | -               | GitHub Notebooks              | en             |
| HumanEval (Chen et al., 2021) | 164     | 4      | ✓           | 7.7             | Hand-written                  | en             |
| MBPP (Austin et al., 2021)    | 974     | 8      | ✓           | 3.0             | Hand-written                  | en             |
| APPS (Hendrycks et al., 2021) | 10,000  | 0      | ✓           | 13.2            | Competitions                  | en             |
| DSP (Chandel et al., 2022)    | 1,119   | 16     | ✓           | 2.1             | Github Notebooks              | en             |
| MTPB (Nijkamp et al., 2022)   | 115     | 8      | ✓           | 5.0             | Hand-written                  | en             |
| Exe-DS (Huang et al., 2022)   | 534     | 28     | ✓           | -               | GitHub Notebooks              | en             |
| DS-1000 (Lai et al., 2022)    | 1,000   | 7      | ✓           | 1.6             | StackOverflow                 | en             |
| CoNaLa (Yin et al., 2018)     | 2,879   | open   | ✗           | -               | StackOverflow                 | en             |
| MCoNaLa (Wang et al., 2022)   | 896     | open   | ✗           | -               | StackOverflow                 | es, ja, ru     |
| ODEX                          | 945     | 79     | ✓           | 1.8             | StackOverflow<br>Hand-Written | en, es, ja, ru |

```
(calculate the improper integral given by the function f
from the number `n` to infinity)
"""
return
```

```
return sympy.integrate(f, (sympy.symbols('x'), n, sympy.oo))

# test case
x = sympy.symbols('x')
f = (x * x)
n = 1
assert str(function(f, n)) == 'oo'
```

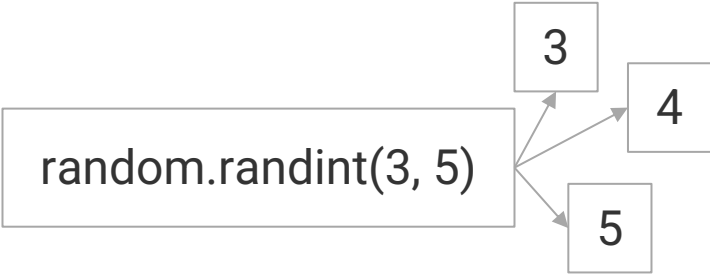
# ODEX: Unique Challenges for Execution

Closed-domain code: easy to execute and verify

```
assert func([1, 2, 10]) == [2, 3, 11]
```

Open-domain code:

▶ Random outputs



▶ Specialized verification

▶ (Potentially) not reproducible queries

▶ HTTP requests, e.g., `requests.post("https://def.xyz", data={'key': 'value'})`

```
In [1]: import numpy as np
In [2]: a = np.array([1, 2, 3])
In [3]: b = np.array([1, 2, 3])
```

```
In [4]: assert (a == b)
-----
ValueError                                Traceback (most recent call last)
Cell In[4], line 1
----> 1 assert (a == b)

ValueError: The truth value of an array with more than one element is ambiguous.

In [5]: np.array_equal(a, b)
Out[5]: True
```

# Significant Performance Gaps: Open vs. Closed

- ▶ Although Codex performs better overall
- ▶ CodeGen has smaller domain gaps

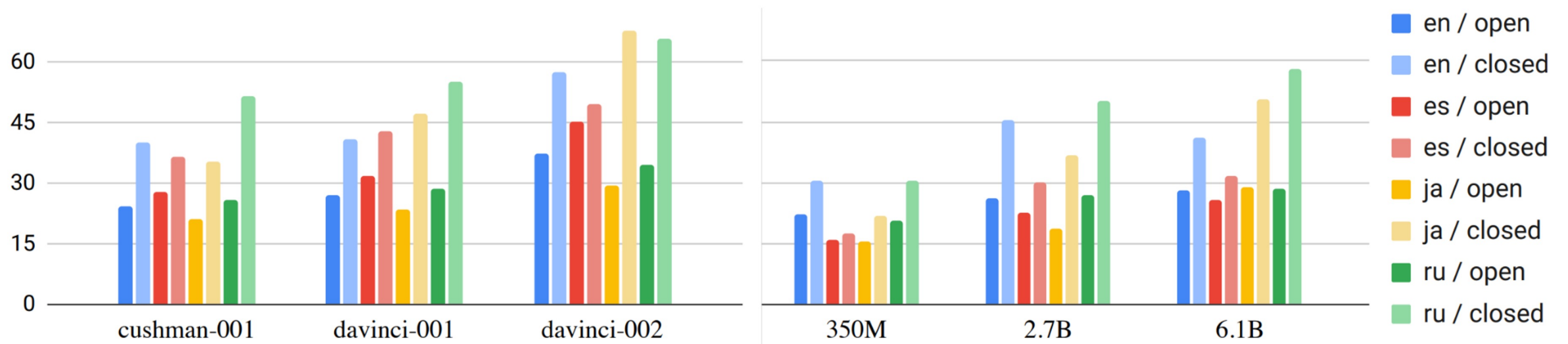


Figure 7: CODEX (left) and CODEGEN (right) *pass@1* on open- and closed-domain problems in each language.

# Code Complexity

# Functional Co

- ▶ Function Level: HumanEval, MBPP
- ▶ Class Level: ClassEval

```

from datetime import datetime
class VendingMachine:
    """This is a class to simulate a vending machine, including adding products, inserting
    coins, purchasing products, viewing balance, replenishing product inventory, and
    displaying product information. """
    def __init__(self):
        """
        Initializes the vending machine's inventory and balance.
        """
        self.inventory= []
        self.balance= {}
    def purchase_item(self, item_name):
        """ Purchases a product from the vending machine and returns the balance after the
        purchase.
    """

```

**Import Statements**

**Class Name**

**Class Description**

**Class Constructor**

**Method Signature**

**Functional Description**

**HumanEval Function Test**

```

METADATA = {\n 'author': 'jt',\n 'dataset': 'test'\n}
def check(candidate):
    assert candidate([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.3) == True
    assert candidate([1.0, 2.0, 3.9, 4.0, 5.0, 2.2], 0.05) == False
    ...

```

**MBPP Function Test**

```

[
    "assert get_ludic(10) == [1, 2, 3, 5, 7]",
    "assert get_ludic(25) == [1, 2, 3, 5, 7, 11, 13, 17, 23, 25]", ...
]

```

**ClassEval Method Test**

```

class VendingMachineTestPurchaseItem(unittest.TestCase):
    def test_purchase_item(self):
        vm = VendingMachine()
        vm.inventory = {'Coke': {'price': 1.25, 'quantity': 10}}
        vm.balance = 1.25
        self.assertEqual(vm.purchase_item('Coke'), 0.0)
        self.assertEqual(vm.inventory, {'Coke': {'price': 1.25, 'quantity': 9}})
    def test_purchase_item_2(self):
        vm = VendingMachine()
        vm.inventory = {'Coke': {'price': 1.25, 'quantity': 10}}
        vm.balance = 1.25
        self.assertEqual(vm.purchase_item('Pizza'), False)
        self.assertEqual(vm.inventory, {'Coke': {'price': 1.25, 'quantity': 10}})
    ...

```

**ClassEval Class Test**

```

class VendingMachineTestMain (unittest.TestCase):
    def setUp(self) -> None:
        self.vm = VendingMachine()
        self.vm.inventory = {'Coke': {'price': 1.25, 'quantity': 10}}
        self.vm.balance = 0
    def test_all(self):
        self.assertEqual(vm.insert_coin(1.25), 1.25)
        self.assertEqual(vm.purchase_item('Coke'), 0.0)
        self.assertEqual(vm.inventory, {'Coke': {'price': 1.25, 'quantity': 9}})
        self.assertEqual(vm.restock_item('Coke', 10), True)
        self.assertEqual(vm.inventory, {'Coke': {'price': 1.25, 'quantity': 19}})
        self.assertEqual(vm.display_items(), 'Coke - $1.25 [19]')
    ...

```

**Test Fixtures:** setUp

Figure 4: Test Cases in Existing Benchmarks and ClassEval

"Write a python function to find the first repeated character in a given string."

Figure 1: Examples in Existing Benchmarks

```

returns False.
>>> vendingMachine.inventory = {'Coke': {'price': 1.25, 'quantity': 10}}
>>> vendingMachine.restock_item('Coke', 10)
True
>>> vendingMachine.inventory
{'Coke': {'price': 1.25, 'quantity': 20}}
"""
...

```

**Parameter/Return Description**

**Example Input/Output**

Figure 2: An Example of Class Skeleton in ClassEval

# Functional Complexity

- ▶ Function Level: HumanEval, MBPP
- ▶ Class Level: ClassEval
- ▶ Repository Level:
  - ▷ RepoCoder
    - ▶ Retrieval-augmented generation
    - ▶ Multiple iterations
  - ▷ RepoEval
    - ▶ Collected 14 Github Repositories
    - ▶ Metrics:
      - ▷ exact match
      - ▷ exact similarity
      - ▷ execution

```
# Below are some referential code fragments from other files: Retrieved Code
# -----
# the below code fragment can be found in:
# tests/test_pipelines_common.py
# -----
# @unittest.skipIf(torch_device != "cuda")
# def test_to_device(self):
#     components = self.get_dummy_components()
#     pipe = self.pipeline_class(**components)
#     pipe.progress_bar(disable=None)
#     pipe.to("cpu")
# -----
"""Based on above, complete the following code:"""

@unittest.skipIf(torch_device != "cuda") Unfinished Code
def test_float16_inference(self):
    components = self.get_dummy_components()

    pipe = self.pipeline_class(**components) Model Prediction
    pipe.to(torch_device)
```

Figure 3: A visual example demonstrating the format of the RepoCoder prompt, which combines the retrieved code snippets from the repository with the unfinished code present in the target file.