

# CS598JBR MP Progress Report: Team-25

Chih-Fu Lai, Heyang Li, Jinsong Yuan, Yihang Jiao  
University of Illinois Urbana-Champaign  
[chihful2, heyangl4, jinsong3, yihangj3]@illinois.edu

## 1 Team Information

Link to the GitHub Repository:

<https://github.com/lycoris106/CS598JBR-Team-25>

Link to the Google Colab Workspace:

<https://colab.research.google.com/drive/1t-d6AL-alJ5accocCVUS9A4YbfaPTn0S?usp=sharing> listings

## 2 MP1

### 2.1

- $n$  = total number of samples generated for a problem
- $c$  = number of correct samples (passes all tests)
- $k$  = number of samples considered

$$pass@k = 1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}, n \geq k$$

This calculates the probability that at least one of the top- $k$  generated code samples passes all test cases.

### 2.2

- base-processed pass@1: 0.4
- instruct-processed pass@1: 0.6

### 2.3

See the attached table

### 2.4

Issues and resolutions:

- Remove duplicate definition: The models often output multiple versions of the function (in-place version, etc.). For example, following is part of the original base model's response of HumanEval/85:

```
1 def add(lst):
2     """Given a non-empty list of integers lst. add
3     the even elements that are at odd indices..
4
5     Examples:
6         add([4, 2, 6, 7]) ==> 2
7     """
8     return sum(lst[1::2])
9
10
11 def add(lst):
12     """Given a non-empty list of integers lst. add
13     the even elements that are at odd indices..
14
15     Examples:
16         add([4, 2, 6, 7]) ==> 2
17     """
18     return sum(lst[1::2])
19
20
```

```
21 def add(lst):
22     """Given a non-empty list of integers lst. add
23     the even elements that are at odd indices..
24
25     Examples:
26         add([4, 2, 6, 7]) ==> 2
27     """
28     return sum(lst[1::2])
29
30
31 def add(lst):
32     """Given a non-empty list of integers lst. add
33     the even elements that are at odd indices..
34
35     Examples:
36         add([4, 2, 6, 7]) ==> 2
37     """
38     return sum(lst[1::2])
39
40
```

The response is composed of multiple "add" functions while the last one is incomplete. This is resolved by splitting all the functions and extracting only the first function definition which matches the function header in the prompt when post-processing:

```
1 def split_multiple_functions(code):
2     triple_string_pattern = r'("""[\s\S
3     ]*?""")|\'\'\'[\s\S]*?\'\'\'')
4     placeholders = {}
5
6     def replacer(match):
7         key = f"__PLACEHOLDER_{len(placeholders)}"
8         placeholders[key] = match.group(0)
9         return key
10
11     protected_code = re.sub(triple_string_pattern,
12                             replacer, code)
13     raw_blocks = protected_code.split("\n\n\n")
14
15     restored_blocks = [
16         re.sub("|".join(map(re.escape,
17                             placeholders.keys()))
18                 lambda m: placeholders[m.group(0)],
19                 block)
20         for block in raw_blocks
21     ]
22
23     return [b.strip() for b in restored_blocks if
24             b.strip()]
25
26 def filter_response(prompt, response):
27     """
28     Filters out any extraneous text before the
29     first function definition.
30     """
```

**Table 1: Pairwise comparison of pass/fail results for each of the 20 programs under two different settings, before and after post-processing.**

Problem ID	Base Results	Base Results Processed	Instruct Results	Instruct Results Processed
HumanEval/96	✗	✗	✓	✓
HumanEval/117	✗	✗	✗	✓
HumanEval/85	✗	✗	✗	✓
HumanEval/148	✗	✓	✓	✓
HumanEval/64	✗	✗	✗	✗
HumanEval/149	✓	✓	✗	✗
HumanEval/13	✗	✗	✗	✓
HumanEval/95	✗	✗	✗	✗
HumanEval/126	✗	✗	✗	✗
HumanEval/123	✗	✗	✗	✗
HumanEval/162	✗	✗	✗	✗
HumanEval/41	✗	✗	✗	✓
HumanEval/21	✗	✓	✗	✓
HumanEval/61	✗	✗	✗	✓
HumanEval/97	✗	✓	✓	✓
HumanEval/27	✓	✓	✗	✓
HumanEval/1	✗	✗	✓	✓
HumanEval/68	✓	✓	✗	✗
HumanEval/150	✗	✓	✗	✓
HumanEval/112	✗	✓	✗	✗

```

26 # find the first function in prompt using
    regex
27 import re
28 match = re.search(r'def\s+\w+\s*\(.*\?\)\s*(->\s*.*\s*)?\s*:', prompt)
29 if not match:
30     return None
31
32 clean_prompt = match.group(0)
33 functions = split_multiple_functions(response)
34 for func in functions:
35     # if the func string contains the
    clean_prompt, return it
36     if clean_prompt in func:
37         return func
38 return None
39

```

- Remove redundant statements: The models often generate redundant, sometimes syntactically-wrong testing statements.

For example, the original instruct model's response of HumanEval/117 contains some test statement at the end, including a call to an unknown "select" function:

```

1 def select_words(s, n):
2     """Given a string s and a natural number n,
    you have been tasked to implement
3     a function that returns a list of all words
    from string s that contain exactly
4     n consonants, in order these words appear in
    the string s.
5     If the string s is empty then the function
    should return an empty list.
6     Note: you may assume the input string contains
    only letters and spaces.
7     Examples:

```

```

8     select_words("Mary had a little lamb", 4) ==>
    ["little"]
9     select_words("Mary had a little lamb", 3) ==>
    ["Mary", "lamb"]
10    select_words("simple white space", 2) ==> []
11    select_words("Hello world", 4) ==> ["world"]
12    select_words("Uncle sam", 3) ==> ["Uncle"]
13    """
14    vowels = "aeiou"
15    words = s.split()
16    result = []
17    for word in words:
18        consonants = 0
19        for letter in word:
20            if letter.lower() not in vowels:
21                consonants += 1
22        if consonants == n:
23            result.append(word)
24    return result
25
26
27 print(select_words("Mary had a little lamb", 4))
    # ["little"]
28 print(select_words("Mary had a little lamb", 3))
    # ["Mary", "lamb"]
29 print(select("simple white space", 2)) # []
30 print(select_words("Hello world", 4)) # ["world"]
31 print(select_words("Uncle sam", 3)) # ["Uncle"]
32

```

We naturally remove these statements with our function extraction shown previously.

- Completion-only: The models always return several full functions instead of completions which are required by HumanEval (like the example HumanEval/117 above). Our

post-process converts the full function output to the completion by deleting function definition and comments.

```

1 def get_completion(prompt, response):
2     response = filter_response(prompt, response)
3     if response:
4         # remove the definition line from the
         response
5         lines = response.split("\n")
6         if len(lines) > 1:
7             lines = lines[1:]
8
9         # remove the comments/docstring if any (
         including the content)
10        while lines and (lines[0].strip().
         startswith('"""') or lines[0].strip().
         startswith("'''")):
11            lines.pop(0)
12            while lines and not (lines[0].strip().
         endswith('"""') or lines[0].strip().endswith
         ("''')):
13                lines.pop(0)
14            if lines:
15                lines.pop(0) # remove the ending
         docstring line
16            while lines and lines[0].strip().
         startswith("#"):
17                lines.pop(0)
18
19        return "\n".join(lines)
20    return None
21

```

The base model's output for HumanEval/85 after post-processing:

```

1     return sum(lst[1::2])

```

The instruct model's output for HumanEval/117 after post-processing:

```

1     vowels = "aeiou"
2     words = s.split()
3     result = []
4     for word in words:
5         consonants = 0
6         for letter in word:
7             if letter.lower() not in vowels:
8                 consonants += 1
9         if consonants == n:
10            result.append(word)
11    return result
12

```

- Wrong interpretation: The models sometimes fail to understand more complex function descriptions and generate wrong solutions. This is difficult to resolve with post-processing without human inspection.
- Edge cases: The models sometimes do not consider all edge cases of inputs (even when explicitly stated in prompts). This is difficult to resolve with post-processing without human inspection.

### 3 MP2

TBA

### 4 MP3

TBA