

# **Mercury:** An Efficiency Benchmark for LLM Code Synthesis

*Du* **Mingzhe**

# What is Code Efficiency?

```
# Sort an array of integers in ascending order.
```

```
# Quick Sort - 121 ms
```

```
def sortArray(self, nums):  
    def quicksort(nums, l, r):  
        if r - l ≤ 1: return  
        pivot = partition(nums, l, r)  
        quicksort(nums, l, pivot)  
        quicksort(nums, pivot + 1, r)  
    quicksort(nums, 0, len(nums))  
    return nums
```

```
# Bubble Sort - 5714 ms
```

```
def sortArray(self, nums):  
    i = 0  
    while i < len(nums)-1:  
        j = i + 1  
        while j < len(nums):  
            if nums[i] > nums[j]:  
                nums[i],nums[j] = nums[j],nums[i]  
            j += 1  
        i += 1  
    return nums
```

🕒 Runtime

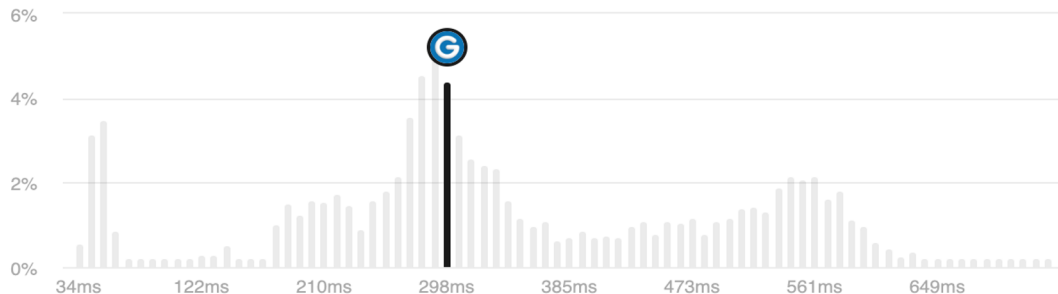
**305 ms**

👏 Beats 54.74% of users with Python3

💾 Memory

**73.67 MB**

Beats 42.86% of users with Python3



$$\text{Beyond}@K = \frac{\sum_{N,K}^{n=0,k=0} \left[ 1 - \frac{\max(R_A^n) - R_k^n}{\max(R_A^n) - \min(R_A^n)} \right]}{N \cdot K}$$

# Why Efficiency is Important?

One Example - GTA Online 

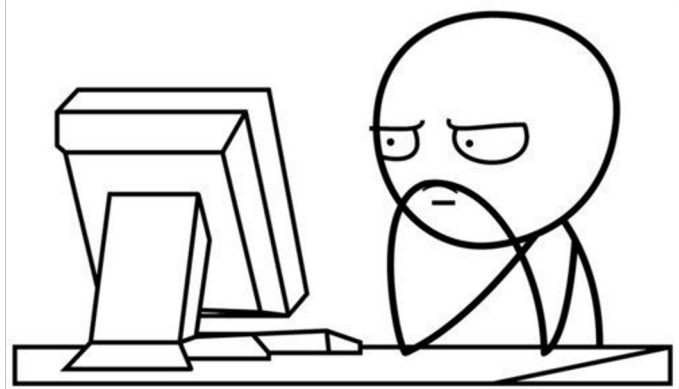
Slow Game Loading.

Lasted 7 Years.

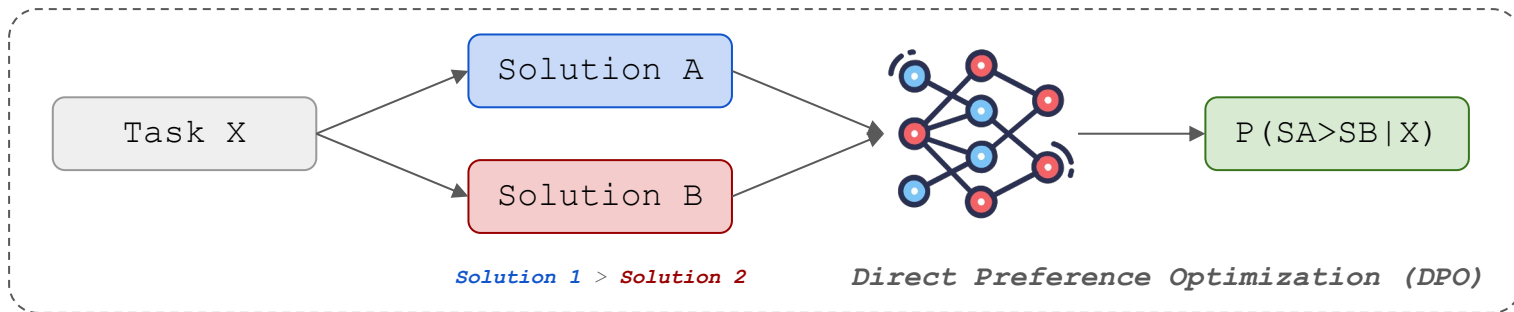
6 Mins to 1 Mins.

**3,046,557** Active Monthly Users.

$7 * 5 * 3,046,557 = 106,629,495$  mins are wasted!



# How to Improve Efficiency?



# Results

Model Name	Mercury-easy			Mercury-medium			Mercury-hard			Mercury-Average		
	B@1	B@3	B@5	B@1	B@3	B@5	B@1	B@3	B@5	B@1	B@3	B@5
deepseek-coder-6.7b	63.55	72.55	<u>75.19</u>	<u>58.80</u>	65.70	67.28	<u>45.72</u>	48.52	50.58	<u>55.99</u>	62.22	64.33
+ SFT	57.70	72.57	75.18	54.87	65.24	66.77	41.87	47.54	49.62	51.42	61.74	63.83
+ DPO	<u>58.29</u>	<u>73.03</u>	75.65	55.86	<u>66.15</u>	<u>67.77</u>	43.70	<u>49.52</u>	<u>51.54</u>	52.56	<u>62.87</u>	<u>64.96</u>
CodeLlama-7b	<u>46.49</u>	55.33	56.37	36.08	46.02	49.05	6.83	9.46	9.48	29.72	36.79	38.12
+ SFT	42.57	52.48	54.54	37.48	40.54	45.48	5.11	3.97	4.86	28.23	32.22	34.79
+ DPO	45.90	<u>55.87</u>	<u>56.83</u>	<u>40.58</u>	<u>46.54</u>	<u>49.57</u>	<u>8.32</u>	<u>10.41</u>	<u>10.47</u>	<u>31.45</u>	<u>37.47</u>	<u>38.78</u>
CodeLlama-13b	<u>42.30</u>	48.47	52.79	39.95	53.32	55.60	7.94	14.61	15.75	29.88	38.50	41.09
+ SFT	42.21	54.90	56.36	35.54	44.47	47.71	3.10	8.78	9.68	26.81	35.92	37.76
+ DPO	40.08	<u>54.50</u>	<u>59.67</u>	<u>40.99</u>	<u>59.19</u>	<u>61.18</u>	<u>13.70</u>	<u>19.16</u>	<u>20.54</u>	<u>31.40</u>	<u>43.97</u>	<u>46.85</u>
CodeLlama-34b	62.50	74.35	77.58	54.02	65.84	67.17	23.17	24.02	28.51	46.45	54.55	57.61
+ SFT	64.49	78.51	83.99	53.88	69.47	71.00	15.17	22.09	27.52	44.37	56.48	60.69
+ DPO	<u>71.62</u>	<u>84.24</u>	<u>87.89</u>	<u>63.15</u>	<u>74.81</u>	<u>76.14</u>	<u>27.30</u>	<u>32.19</u>	<u>36.41</u>	<u>53.88</u>	<u>63.57</u>	<u>66.68</u>
deepseek-coder-33b	67.92	76.61	78.78	67.18	73.34	76.27	52.36	58.79	59.43	62.40	69.52	71.41
+ SFT	68.58	79.85	83.04	67.16	76.10	79.38	45.09	55.55	57.87	60.15	70.41	73.33
+ DPO	<b>77.83</b>	<b>86.06</b>	<b>88.83</b>	<b>75.15</b>	<b>83.08</b>	<b>85.20</b>	<u>60.52</u>	<b>66.70</b>	<b>67.64</b>	<b>71.10</b>	<b>78.53</b>	<b>80.48</b>
gpt-4-preview-1106	70.54	-	-	62.53	-	-	<b>70.06</b>	-	-	67.84	-	-

Table 4: Code Efficiency Evaluation Results Categorized by Difficulty Levels. The notation  $B@K$  represents the *Beyond@K* score with  $K$  generated solutions. **The bolded value** indicates the top performance for each metric, while the underlined values denote the most effective approaches among the original model and the baselines.

**Thank you** 😊