

# Regression Language Models for Code



Yash Akhauri<sup>\*1</sup>, Xingyou Song<sup>\*2</sup>, Arissa Wongpanich<sup>2</sup>, Bryan Lewandowski<sup>2</sup>, Mohamed S. Abdelfattah<sup>1</sup>  
<sup>1</sup>Cornell University, <sup>2</sup>Google

## Goal

We study code-to-metric regression:

*Predicting a numeric outcome of a code execution.*

Why challenging?

- Programming languages very open-ended.
- Strictly non-tabular data, very understudied!

## Example

Consider the following examples from APPS (Python LeetCode):

Can you predict memory usage precisely?

Problem (Distance Value)

Given two integer arrays `arr1` and `arr2`, and the integer `d`, return the distance value between the two arrays. The distance value is defined as the number of elements `arr1[i]` s.t. there is not any element `arr2[j]` where  $|arr1[i] - arr2[j]| \leq d$ .

Memory-efficient (O(1) extra space)

```

From typing import List

class Solution:
    def findTheDistanceValue(self, arr1: List[int], arr2: List[int], d: int) -> int:
        count = 0
        for a in arr1:
            for b in arr2:
                if abs(a - b) <= d:
                    # found a conflict
                    break
            else:
                count += 1
        return count

```

Less memory-efficient (hash structures)

```

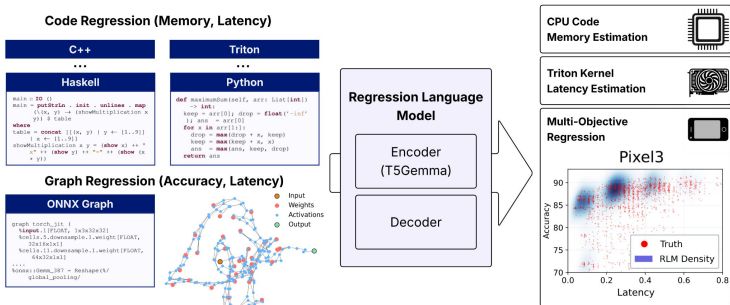
From typing import List
from collections import Counter

class Solution:
    def findTheDistanceValue(self, arr1: List[int], arr2: List[int], d: int) -> int:
        # overheard! build dict
        arr2_counts = Counter(arr2)
        # overheard! build hash set
        arr2_set = set(arr2)

        total = 0
        for a in arr1:
            target = a - d, a + d + 1
            # overheard! new set
            arr2_candidates = set(arr2_set.intersection(target))
            # overheard!
            total += arr2_candidates
        return total

```

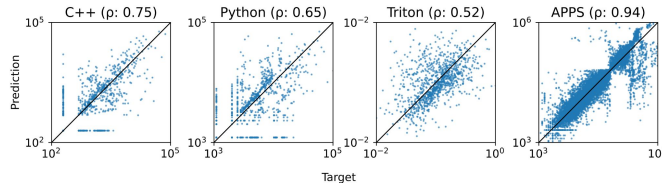
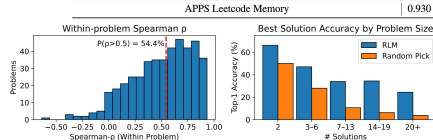
## Overview



## Programming Languages

A single model can predict decently across 24+ programming languages!

Language	$\rho$	Language	$\rho$	Language	$\rho$	Language	$\rho$
C++	0.748	Go	0.670	Python	0.647	Kotlin	0.634
C	0.741	D	0.656	OCaml	0.643	Swift	0.630
Lisp	0.625	Lua	0.618	Haskell	0.611	Rust	0.611
Perl	0.592	C#	0.583	Java	0.560	Scala	0.537
Fortran	0.527	TypeScript	0.463	Pascal	0.461	Ruby	0.460
Bash	0.455	FW	0.439	JavaScript	0.395	PHP	0.347



## Computational Graphs

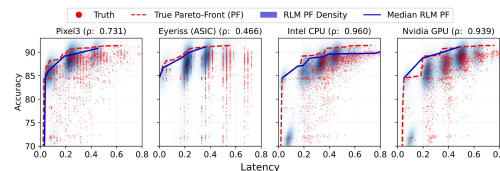
Consider using ONNX Graphs for Neural Architecture Search (NAS):

Matches SoTA Graph Neural Networks (GNNs)!

Method	NASNet	Amoeba	PNAS	ENAS	DARTS	Average
MLP (Adjacency Enc.)	0.002	0.032	0.082	0.021	0.124	0.052
Arch2Vec (Graph Enc.)	0.209	0.107	0.184	0.224	0.333	0.212
CATE (Transformer Enc.)	0.150	0.160	0.217	0.236	0.425	0.238
GNN	0.364	0.376	<b>0.444</b>	0.438	0.523	0.429
FLAN (Previous SoTA)	0.344	0.470	0.430	<b>0.484</b>	<b>0.567</b>	0.459
<b>RLM (Ours)</b>	<b>0.382</b>	<b>0.488</b>	0.427	0.481	0.528	<b>0.461</b>

## Multi-Objective Predictions

Decoding twice allows multi-objective predictions across different hardware:



## Why Prompting Doesn't Work

Consider in-context learning (ICL) with Gemini.

Each example graph takes 32K tokens; 1M context = 50-shots only.

- Can't absorb data
- Way more expensive
- Not trained for this task

Weight-based updates vs. In-context learning

