
Measuring the Impact of Programming Language Distribution

Gabriel Orlanski^{1 2*} Kefan Xiao² Xavier Garcia³ Jeffrey Hui² Joshua Howland² Jonathan Malmaud²
Jacob Austin³ Rishabh Singh^{2*} Michele Catasta^{2*}

Abstract

Current benchmarks for evaluating neural code models focus on only a small subset of programming languages, excluding many popular languages such as Go or Rust. To ameliorate this issue, we present the BabelCode framework for execution-based evaluation of any benchmark in any language. BabelCode enables new investigations into the qualitative performance of models’ memory, runtime, and individual test case results. Additionally, we present a new code translation dataset called Translating Python Programming Puzzles (TP3) from the Python Programming Puzzles (Schuster et al., 2021) benchmark that involves translating expert-level python functions to any language. With both BabelCode and the TP3 benchmark, we investigate if balancing the distributions of 14 languages in a training dataset improves a large language model’s performance on low-resource languages. Training a model on a balanced corpus results in, on average, 12.34% higher *pass@k* across all tasks and languages compared to the baseline. We find that this strategy achieves 66.48% better *pass@k* on low-resource languages at the cost of only a 12.94% decrease to high-resource languages. In our three translation tasks, this strategy yields, on average, 30.77% better low-resource *pass@k* while having 19.58% worse high-resource *pass@k*.¹

1. Introduction

In the 2022 StackOverflow Developer Survey, Rust was the 14th most popular programming language despite not rank-

*Work Done While At Google ¹Department of Computer Science, New York University, New York, New York ²Google Labs ³Google Brain. Correspondence to: Gabriel Orlanski <go533@nyu.edu>, Kefan Xiao <kfxiao@google.com>, Xavier Garcia <xgarcia@google.com>.

Proceedings of the 40th International Conference on Machine Learning, Honolulu, Hawaii, USA. PMLR 202, 2023. Copyright 2023 by the author(s).

¹<https://github.com/google-research/babelcode>

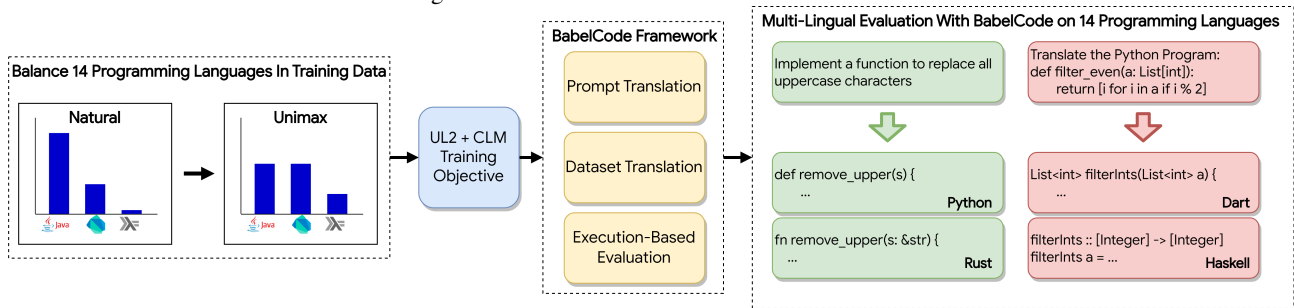
ing in the survey taken five years prior. However, the 13th most popular language, Go, has nearly doubled Rust’s number of StackOverflow questions in this time frame. Further, despite their similar popularity, Go has nearly 350% more source code available (Kocetkov et al., 2022). These disparities highlight the problem that many popular programming languages are starkly low-resource, especially compared to the most popular languages.

Despite their impressive generative capabilities, especially in code, Large Language Models (LLM) are adversely impacted by this language resource imbalance. Thus, developers will likely find minimal utility from LLMs if they are not using the extremely popular languages. It is therefore imperative to investigate how to mitigate the discrepancy between a language’s popularity and the amount of data available for it. Prior works focusing on code generation (Ahmad et al., 2021) and multilingual natural language processing (Arivazhagan et al., 2019; Conneau et al., 2019) use temperature-based strategies to balance the training languages. Such a strategy duplicates extremely low-resource languages thousands of times, which has been shown to significantly reduce performance (Allamanis, 2019).

Beyond the the language balancing strategy, evaluating code LLMs in a multi-lingual setting presents significant challenges. Existing datasets are either mono-lingual (Chen et al., 2021; Austin et al., 2021; Lai et al., 2022) or limited to only a subset of popular programming languages (Roziere et al., 2020). Each problem in these datasets, which we henceforth refer to as a *benchmark*, contains an input, and a canonical solution along with the test-cases for checking correctness. Creating a new benchmark for each language of interest would require insurmountable engineering and monetary costs. To address both of these problems, we present the BabelCode framework for execution-based evaluation of *any benchmark in any language* and use it to investigate the impact of programming language distribution on code generation and translation.

BabelCode is open-sourced, has an extensive test suite, and supports evaluating four benchmarks in 14 languages. It is designed specifically to enable future research directions such as the evaluation of custom data-structures. BabelCode allows investigation of novel research directions through

Figure 1. Overview of this work’s contributions.



the measurement of memory and runtime usage for a given prediction, as well as the outcomes of individual test cases. Furthermore, we can use BabelCode to build multi-lingual execution based benchmarks from existing mono-lingual datasets. We demonstrate this functionality by creating a new dataset called Translating Python Programming Puzzles (TP3) from the Python Programming Puzzles (Schuster et al., 2021) benchmark, where the objective is to translate expert-level python programs to other languages. The source programs for TP3 are the hand-crafted verification functions for each problem in P3. As the authors hand-wrote each function, they are significantly more complex than the current state-of-the-art code translation benchmarks, such as Transcoder (Roziere et al., 2020), for which code LLMs are already achieving highly impressive results.

Our presented framework is closely related to the concurrent work of MBXP (Athiwaratkun et al., 2023) and Multi-PLE (Cassano et al., 2022). While MBXP is quite similar to BabelCode, it is not open-sourced and requires that the input benchmarks be in Python. Multi-PLE is open-sourced, but only supports generation tasks and contains significant errors in multiple languages. BabelCode addresses these issues through an extensive test suite that ensures that the code generated is correct, and that crucial functionality, such as data structure equivalence, works when executed.

With the BabelCode framework, we investigate remedies to the problems of programming language imbalance. We utilize the Unimax algorithm (Chung et al., 2023) to limit the maximum number of times to duplicate a language’s data to a constant N . We then train 1B, 2B, and 4B parameter decoder-only models on both the natural and Unimax N distributions. We utilize the UL2 (Tay et al., 2022) and causal language modeling training objective. We find that models trained on the balanced dataset significantly outperform the baseline models on low-resource languages across all tasks. Further, we find that the resulting performance drop on high-resource languages is mitigated by increasing the model size.

This paper makes the following key contributions:

- We propose and release BabelCode, a new execution-based evaluation framework that allows for multilingual evaluation of code generation and translation capabilities of code language models. It also supports the easy addition of new benchmark tasks and execution-based metrics.
- We show that the code language models trained on the natural distributions of GitHub source code have poor performance on low-resource languages in both generation and translation tasks.
- We propose a new data balancing strategy for programming languages to improve performance on low-resource languages. We demonstrate that the resulting models outperform the baseline models across all tasks by an average of 12.34% $pass@k$ for all languages, with a further improvement of 39.70% $pass@k$ to low-resource languages.
- We find that the average improvements on low-resource languages from training on balanced data do not scale with model size. But scaling model sizes significantly helps the average $pass@k$ loss compared to the baselines on high-resource languages going from a loss of 39.70% with the 1B model to a loss of 2.47% with the 4B model.

2. The BabelCode Framework

BabelCode enables the evaluation of a collection of problems, each consisting of a prompt and a set of test cases, in any language through four stages: 1) represent each test case in our domain specific language (DSL) defined in Figure 2, 2) use this generic form to generate the test cases in the target language from the input and output values, 3) use a Jinja² template to generate a testing script in the target language, and 4) execute the target script through the command line. This is done autonomously, requiring minimal human intervention. We provide an overview of how an example

²<https://jinja.palletsprojects.com/en/3.1.x/>

Table 1. Differences between BabelCode and prior works. NL2C is natural language to code, while C2C is code to code datasets. BabelCode has an extensive test-suite that automatically tests each language’s implementation and correctness when executed.

Name	Open Sourced	# Lang.	NL2C Support	C2C Support	Mem. & Time Metrics	Test Suite	Indiv. Test Case Results	Lang. Agnostic Datasets
MultiPL-E	✓	18	✓	✗	✗	✗	✗	✗
MBXP	✗	10	✓	✓	✗	✗	✓	✗
BabelCode	✓	14	✓	✓	✓	✓	✓	✓

Figure 2. BabelCode’s domain specific language for representing the input and output types of a question. Prior works require that the source dataset be written in Python, while our DSL removes this restriction and allows users to create datasets in *any* language. This enables seamless additions of new languages while simplifying future expansions to features such as custom data structures.

```

S → LeafTypes | ListType | MapType
ListType → list < S > | set < S >
MapType → map < CoreTypes; S >
LeafTypes → CoreTypes | boolean | double | float | long
CoreTypes → character | integer | string

```

problem is translated in Figure 8. Overall the key novel elements of BabelCode are: I) the use of a DSL to translate programming questions, II) type-specific equivalence, III) the ability to measure the performance of a given program at a low level (i.e., memory used, runtime, which tests passed), and IV) a large scale test-suite for ensuring correctness of generated code.

2.1. Framework Design

BabelCode shares many design similarities to the concurrent work from Athiwaratkun et al. (2023). Specifically, we follow the same approach to inferring argument and return types. We follow the respective documentation and tutorials for each language to determine which native types to use. We also use these docs to determine the docstring formatting and naming convention. These mappings are used to generate unit and integration tests for each language automatically. They ensure that each language’s implementation is syntactically correct and that, when executed, the equality comparison is correct. We describe the framework design and similarities to Athiwaratkun et al. (2023) in Appendix A.

DSL Representations: Using a DSL in the first phase, we do not force the inputs to be Python, thus enabling more flexibility to represent more generic tasks. For example, given the inputs from two test cases: {"a": [[1], [], [80]]} and {"a": []}, we only represent the *types* in our generic DSL. Thus, the resulting type string for this input is

map<string;list<integer>>. We do not represent the actual values in the generic form as we can easily translate literals across languages. This allows users to create a dataset from any language by requiring that they only represent the types of the inputs and outputs in this generic form. The language agnostic nature of the DSL enables future extensions of BabelCode to incorporate complex inputs and outputs such as custom data-structures. For example, the representation of a node class in a BST could be BSTNode<integer;integer>.

Equality Checking: We support floating point equivalence to a precision of $\epsilon = 1e-6$ for floats and $\epsilon = 1e-9$ for doubles. To determine if a given value is a float or a double, we count the number of digits after the decimal place. We apply this same logic to int and long by counting the total number of digits. Languages such as C# do not, by default, support deep equivalence of data structures. In such cases, we serialize the objects to JSON and check that the resulting strings are equal. Otherwise, we use the language built-in deep equality functionality.

Test Statement Execution: We opt to print the result of each test case (i.e. TEST-0...PASSED) to the standard output in a parseable format across all languages. Along with try-catch blocks, this allows the evaluation of *every* test case for a given problem. This allows finer analysis of individual programs when compared to using assert statements as it identifies if specific corner cases fail.

Prompt Translation: As Wang et al. (2022a) showed, LLMs are sensitive to the input prompts for code generation. Therefore BabelCode supports prompt translation and construction for multiple different problem formulations. We replace the names of languages, such as Python, with the target language. We use the language-specific naming convention to properly format the signature in the best practice style. If an argument uses a reserved keyword, we append arg to its name so that it retains the same meaning but will no longer conflict. We replace Python-specific terms with their equivalent names in the target language. For tasks formulated as code-completion, we support formatting the problem description as a native docstring. We do *not* translate the import statements in the header. Instead, we exclude the headers from all languages to provide

a language-agnostic format. Translating prompts to a target language is not novel by itself, as both Athiwaratkun et al. (2023) and Cassano et al. (2022) proposed methods to accomplish this. BabelCode’s builds on those works by translating reserved characters. For example, in Julia, the “\$” in docstrings will raise errors if not properly escaped. Thus, we implement methods to automatically handle such cases and ensure correctness.

2.2. Differences To Prior Works

We summarize the high-level differences between BabelCode and prior works in Table 1. The **MBXP** framework from Athiwaratkun et al. (2023) is the most similar to our work as discussed in subsection 2.1. Similar to BabelCode, MBXP does have individual test-case results; however, it uses `assert` statements and thus can only determine the first test-case that fails. MBXP does use language experts to review the generated code’s quality and discuss the validation it supports to ensure that generated code parses and/or compiles for its respective language. BabelCode also has this functionality but, additionally, it ensures correctness through a test suite that covers the execution of generated code. We provide scripts to allow validating that source solutions to a dataset pass the generated code. For languages that do not have a solution in the respective dataset, we generate “mock” predictions that return the expected output type. This allows us to ensure that generated code is correct in *all* supported languages even if no solution exists.

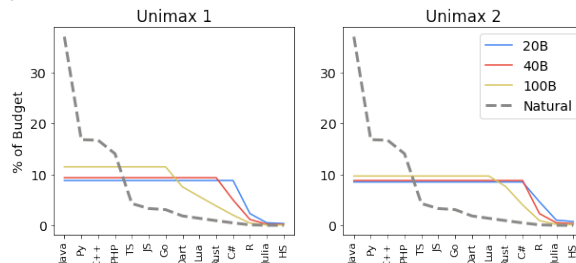
The **MultiPL-E** framework from Cassano et al. (2022) supports 18 languages compared to BabelCode’s 16. However, we support four datasets, while MultiPL-E only currently has support for two datasets. In addition, BabelCode also supports fine-grained evaluation metrics for memory, running time, and individual test cases. Our extensive test suite and validation scripts have also exposed many language-specific idiosyncrasies that naive methods of translation fail to handle. For example, in Julia, any “\$” will be treated as string interpolation, even if it is in a docstring. Thus, in the majority of cases, these must be escaped. We automatically rename variables that use reserved keywords. In languages such as C#, the `==` operator checks equivalence by *reference* instead of *value*. Besides corner cases, our DSL and templates allow us to effectively implement proper floating point equivalence for problems that return a float. Finally, in many languages, MultiPL-E uses types that are *not* considered best practice, such as in Scala, where it relies on the Java types `ArrayList` instead of the native `List`.

3. Low-Resource Code Language Models

Because the data availability can vary greatly by programming language, we can consider the goal of building a multilingual code model as a data-imbalanced multi-task learning

problem. Previous work in the multilingual natural language community (Conneau et al., 2019; Arivazhagan et al., 2019) and in the program synthesis space (Ahmad et al., 2021) have used sampling strategies relying on temperature-scaling. In this work, we use the Unimax (Chung et al., 2023) strategy to address this imbalance. The Unimax algorithm assumes that we are given a budget of how many examples we plan to consume during training and a maximum number of times, N , any single example can be duplicated in the training corpus. Then, we separate the data into buckets by programming language and add N epochs of each of the lowest-resource languages until we can safely distribute the remaining budget across all the remaining languages without exceeding N epochs over any one of these remaining languages. This will allow us to control the number of epochs N we perform over the low-resource languages to minimize overfitting while allowing fair distribution of the compute budget to the remaining high-resource languages. We will ablate the choice of N in our experiments.

Figure 3. Different distributions for Unimax with different budgets.



4. Experimental Setup

4.1. Models

To understand the impact of training decoder-only models on the different programming language distributions, we train models in 3 sizes: 1B, 2B, and 4B. For each of these sizes, we train 5 different models on each distribution: Natural and Unimax N , where $N \in \{1, 2, 3, 4\}$. The parameters and training differences are listed in Table 2. We follow Chowdhery et al. (2022) for all other architecture choices. Every model has a context window of 2048 and is trained identically with the same vocabulary described in subsection 4.3. We use a base learning rate of 0.01 and a constant warmup with a step inverse decay. The number of warmup steps is kept to 10% of the total training steps per model. The total number of training steps is 38000, 77000, 190000 for the 1B, 2B, and 4B models, respectively. We use the Adafactor optimizer (Shazeer & Stern, 2018) and a batch size of 256. We prepend `[code]` to the beginning and add the tag `[eod]` to the end of each file from our training data. Finally, we use the T5X and SeqIO (Roberts et al., 2022) frameworks. We use the UL2 (Tay et al., 2022) objective

Table 2. Hyperparameters for models trained (BC) compared with those used to train PaLM-Coder(PC). For PaLM-Coder, we report the number of code tokens trained on. Each BC model is trained on each of the naturally occurring distributions of the GitHub data and each of the distributions is detailed in section 3 where $N \in \{1, 2, 3, 4\}$

Model	# of Layers	Heads	d_{model}	Train Tokens(B)
BC 1B	16	8	8192	20.2
BC 2B	24	16	10240	40.4
BC 4B	26	16	14336	100
PC 8B	32	16	4096	46.8
PC 62B	64	32	8192	46.8

with an additional causal language modeling objective as described in Appendix D.

4.2. Training Data

Our curated source code corpus was obtained by collecting publicly available code data on the web using a custom code data collection system. We apply a similar license filter as Kocetkov et al. (2022) to remove any files with non-permissible licenses, use simple heuristics to filter out low-quality code and apply near-deduplication to obtain our corpus of high quality, permissive source code. After preprocessing, we select 14 programming languages by their file extensions according to the mapping used by GitHub’s Linguist library³ to segment the dataset by language. To calculate the number of examples per language, we use SeqIO’s caching feature and take the number of examples after post-processing (Roberts et al., 2022). We list the percentages of all examples and file extensions used per language in Appendix C. With these numbers, we consider the top 7 languages to be **high-resource**(HR): Java, Python, C++, PHP, TypeScript, JavaScript, and Go. We further consider the bottom 7 languages to be **low-resource**(LR): Dart, Lua, Rust, C#, R, Julia, and Haskell.

4.3. Vocabulary

The original PaLM (Chowdhery et al., 2022) vocabulary focuses on multilingual natural language. In contrast, we trained our SentencePiece (Kudo & Richardson, 2018) vocabulary with 64k tokens from the training data directly. Each programming language is uniformly sampled to build the vocabulary. In previous works, such as Chen et al. (2021), a list of tokens that consists of a different number of whitespace is manually added to represent code more efficiently. In our work, we rely on the SentencePiece model to learn the whitespace tokens by allowing extra whitespace tokens and whitespace-only tokens. In the end, the model can

represent up to 12 whitespaces into one token. In addition, numbers are split into individual tokens.

4.4. Benchmarks

BabelCode currently supports 4 datasets. To allow the translation of any dataset to any language, we modify each benchmark as well as remove problems that were incompatible. These changes are described in Appendix B. For HumanEval (Chen et al., 2021), MBPP (Austin et al., 2021), and Transcoder (Roziere et al., 2020), we add the prefix **BabelCode-** (**BC**) to indicate that we are using the BabelCode specific version. Further, for Transcoder, we use the same version as in Chowdhery et al. (2022). **BC-HumanEval** (**BC-HE**) has 161 out of the original 164 HumanEval questions. **BC-MBPP** has 855 of the original 999 questions. **BC-Transcoder** (**BC-TC**) has 524 of the original 560 questions.

We additionally introduce a new dataset called **Translating Python Programming Puzzles (TP3)**. We take the verification functions from the questions in the original Python Programming Puzzles dataset (Schuster et al., 2021) to create this dataset. These functions are hand-crafted by the authors and are used to check if an answer satisfies the constraints of the puzzle. These puzzles range in difficulty from basic character checking to competitive programming problems. Thus, each verification function is written by an expert python programmer and requires a significant understanding of programming to translate. In total, there are 370 python functions to translate. Examples from TP3 can be found in subsection B.4.

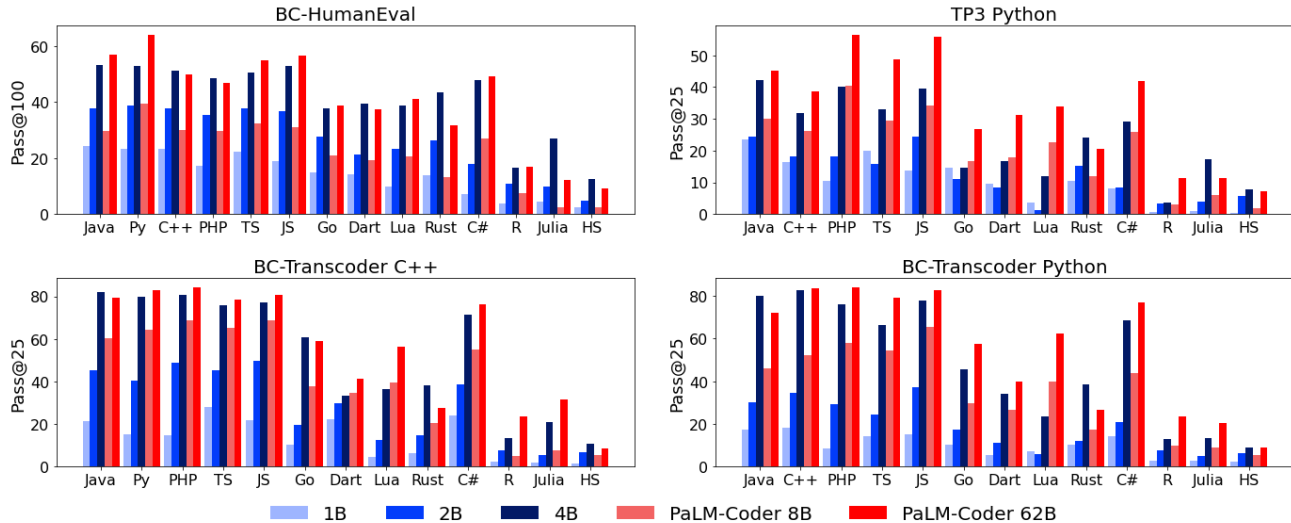
4.5. Evaluation

For BC-HumanEval, we follow Chen et al. (2021) and generate 200 programs per problem. Further, we use a zero-shot prompt described in subsection E.1. We use the built-in docstring translation of BabelCode. We generate 50 programs per problem on our three translation tasks and use the prompts described in subsection E.2. We consider these prompts zero-shot as we do not provide any additional examples. However, we provide the translated signature without the docstring in the prompt. We do not consider this to be data leakage as it is trivial to translate signatures with libraries such as Treesitter⁴.

For every dataset, we use $T = 0.8$, $top_p = 0.95$, and do not use top_k . We use the $pass@k$ estimator (Chen et al., 2021) to measure the performance. We use $k = 100$ and $k = 25$ for generation and translation, respectively.

³<https://github.com/github/linguist/>

Figure 4. Comparison of the models trained with PaLM-Coder models. For each dataset, we use Chen et al. (2021) $pass@k$ estimator with $n = 2 * k$. We then generate n samples per problem with $T = 0.8$. Full results can be found in Appendix F. Languages in the X-Axis are sorted from high to low resource. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.



5. Results

5.1. Baseline Models

We report the baseline results for our trained models and PaLM-Coder in Figure 4. On BC-HumanEval, we find that the 2B model has a better $pass@100$ than that of PaLM-Coder 8B on all but C# and Python. On average, the BC-2B model trained on the natural distribution of GitHub data has average improvements of 48.17% compared to PaLM-Coder 8B despite having a quarter of the number of parameters and training on 6.4B fewer code tokens. Further, we find that the 4B model outperforms PaLM-Coder 62B on 6 of the 14 languages evaluated. This likely results from the 4B model seeing over 53B tokens more than what PaLM-Coder 62B did. Another likely factor in this discrepancy is that the data PaLM-Coder was fine-tuned on included all languages on GitHub in contrast to our filtered training dataset.

We also observe that performance on languages do not scale with respect to their resource level nor the model’s size. C#, Dart, Julia, and Haskell have significantly higher gains when scaling to 4B model size when compared to the other languages. While this may be due to the increased number of training tokens, it is not consistent across all LR languages as the increase in performance for R and Lua when scaling from 1B to 2B is similar to that when scaling from 2B to 4B. Instead, this result is likely due to better transfer from languages such as Java, Python, and C++.

The importance of scale for multi-lingual code models is

⁴<https://tree-sitter.github.io/tree-sitter/>

further demonstrated by the results of the translation tasks. We find that in BC-TP3, the 1B and 2B models’ performance is similar. However, the most significant gains are from scaling up to 4B where it beats PaLM-Coder 8B on all but three languages in this zero-shot translation. We do make note, though, that while we do not provide any examples for in-context learning, we do provide the signature in the target language during generation. This finding is less pronounced in BC-Transcoder as the scaling observed in all languages is more akin to that seen in BC-HumanEval.

5.2. Impact of Balancing Programming Languages

Figure 5 shows the mean $pass@k$ scores of different models trained on each of the 5 distributions for each of the 4 datasets. As expected, the natural distribution is optimal if the focus is solely HR languages as the performance losses when training on Unimax balanced data are 15.47%, 14.00%, and 9.35% for the 1B, 2B, and 4B models, respectively. However, for any LR language, Unimax is clearly better given that there is an average $pass@100$ improvement on these languages of 111.85%, 68.38%, and 19.22% for the 1B, 2B, and 4B size models, respectively. For generation tasks, we find that $N = 3$ is optimal with respect to the difference between performance gained on LR and performance lost on HR languages. On the 1B, 2B, and 4B models, the ones trained on the Unimax 3 dataset had differences of 130.17%, 87.80%, and 36.00%, respectively.

We observe similar scaling trends on TP3, as training on a Unimax distribution yielded average $pass@25$ improvements to LR languages of 124.45% for the 1B model,

Figure 5. Effects of scale on the average $pass@k$ of the high and low resource languages for each of four datasets. Full tabulated results are located in Appendix F.

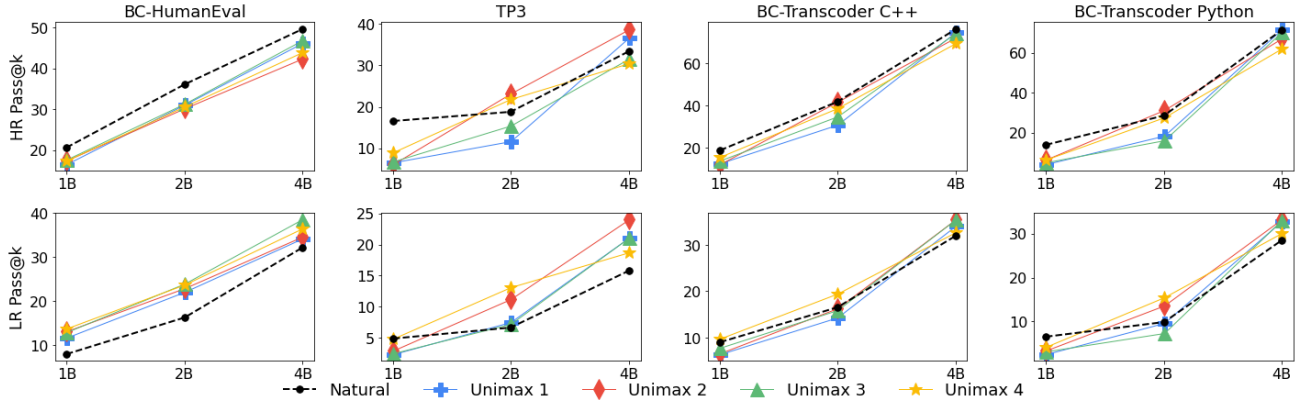
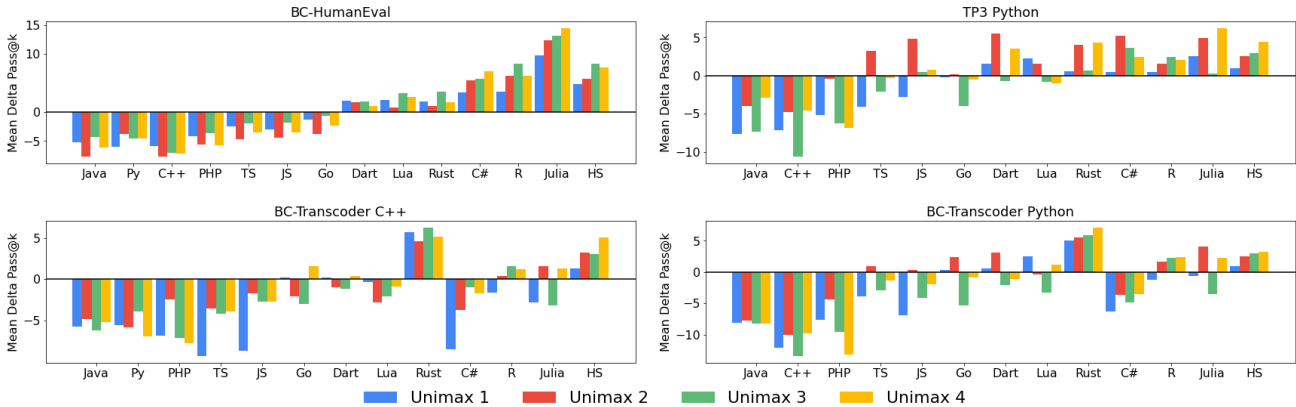


Figure 6. Mean relative difference of $pass@k$ for each of the models trained on the different Unimax distributions compared to the $pass@k$ of the same sized model trained on the Natural distribution. The X-Axis is the language sorted from high to low resource. HS is Haskell and Py is Python. The percent changes for each delta for HR languages are shown in Table 12 and Table 13 for LR languages.



64.51% for the 2B model, and 51.29% for the 4B model when compared to the same sized models trained on the natural distribution. Unlike BC-HumanEval, training the 4B on Unimax Distributions yielded *better* average HR performance with an increase of 6.80%. As shown in Figure 6, training a 4B model on the Unimax 2 distribution had a mean $pass@25$ improvement of 71.59% in LR languages and an improvement of 20.31% on HR languages when compared to the natural distribution. Training on other Unimax distributions does not see as large of improvements. For the 4B model, we find mean LR improvements of 42.39%, 52.91%, and 38.26% when trained on the Unimax 1, 3, and 4 distributions, respectively. This indicates that for TP3, at least, balancing the training data for each language improves translation capabilities. However, less Python data adversely affects understanding the source code necessary to translate it properly.

When evaluated on BC-Transcoder, we find that LR performance *increased* with size. When the source language is C++, training on the Unimax distributions yielded an average $pass@25$ improvements of 7.57%, 6.76%, and 11.80% for the 1B, 2B, and 4B models, respectively. Translating Python to other languages followed this trend with an average change of -26.04%, 15.1%, and 22.47% for the 1B, 2B, and 4B models, respectively. On BC-Transcoder, we find similar benefits when translating from Python to other languages, although the performance on higher resource languages is significantly worse. When translating from C++ to other languages, we find that training both a 1B and 2B model on the UM 4 distribution improves performance on 5 of the 7 LR languages. For 4B sized models, the UM 2 distribution is optimal as LR performance increased by an average of 20.47% when compared to training on the natural distribution. As the source code of BC-Transcoder

focuses on language-agnostic algorithm implementations, this scaling trend is most likely due to the importance of a surface-level understanding of the target language. Further, the fact that this trend does not appear for BC-HumanEval or TP3 indicates that neither model size nor duplication of language data enables the model to have a deep understanding of these low-resource languages.

5.3. Effects Of Language Balance on Predictions

We find that, as is expected, decreasing the number of tokens for a language negatively impacts its performance on that language. To compare the overall effects of language balancing at each size, we focus on the Unimax 1 and Unimax 2 distributions as they represent the largest change in proportions of HR languages when compared to the Natural distribution. Figure 7 shows that on BC-HumanEval, training on either UM 1 or UM 2 will cause the model to generate fewer correct solutions than when the model is trained on the Natural distribution with respect to HR languages. However, this is *not* due to those models generating more programs with either compilation or run-time errors as the raw average increase is only 0.40 and 1.15 for the models trained on the Unimax 1 and Unimax 2 respectively. Rather, we find that the largest decrease is in the mean % test cases passed per problem. Training on the Unimax 1 and Unimax 2 distributions results in 5.50% and 9.09% fewer test cases respectively when compared to the model trained on the natural distribution.

On LR languages, the Unimax 1 distribution yielded the best improvements compared to the other distributions. Specifically, the programs generated by the model trained on the Natural distribution passed, on average, 5.13% of the test cases per problem. In comparison, 9.53% and 10.48% of average test cases per problem were solved by the models trained on the Unimax 1 and Unimax 2 distributions. The less than 1% improvement when going from Unimax 1 to Unimax 2 suggests that, for generation tasks, multi-lingual models of code benefit the most from seeing unique data.

In our translation task of TP3, we observe consistent improvements in the mean number of test cases passed for both HR and LR languages. For the former, we observe an average improvement of 2.58% and 3.06% compared to the Natural distribution for the UM 1 and 2 distributions respectively. On LR languages, we find average improvements of 3.40% and 4.99% over the Natural distribution for the UM 1 and UM 2 distributions respectively. These results, along with the performance improvements discussed in subsection 5.2, indicate that translation tasks benefit highly from uniformly balanced languages. This is, likely, due to the task formulation where natural language understanding is not necessary. Higher resource languages are more likely to contain diverse natural language and code pairs due to the

language’s popularity.

Thus, performance on NL2Code tasks, such as BC-HumanEval, depends on the unique samples of code and doc-strings in the training corpus. Translation, on the other hand, does not have this constraint. Rather, it appears that uniformly balancing languages is the optimal strategy for this task.

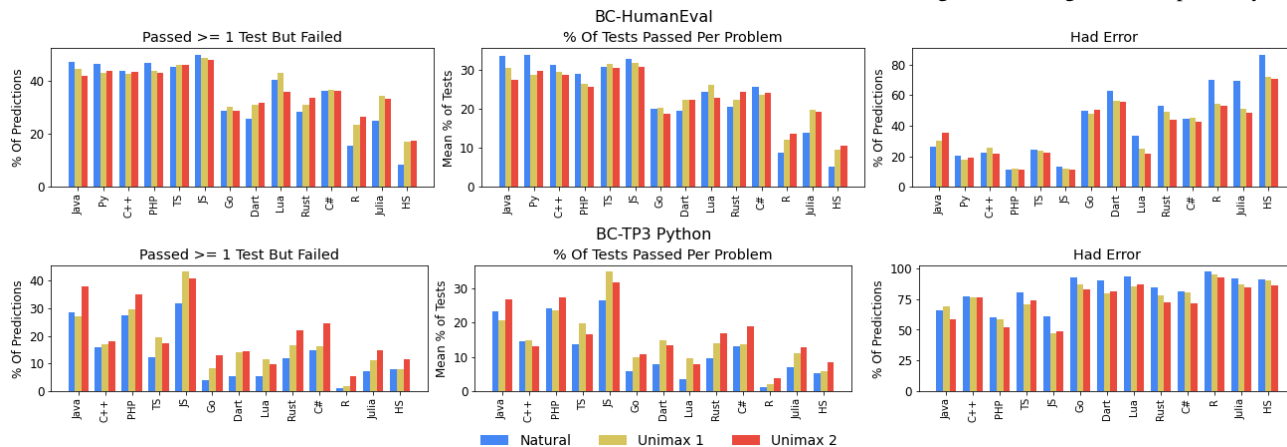
6. Related Works

Code Evaluation Existing code benchmarks have primarily focused on surface matching evaluation (Lu et al., 2021; Yin et al., 2018; Wang et al., 2022b; Husain et al., 2019). Recent works have introduced new execution-based benchmarks for both generation (Austin et al., 2021; Hendrycks et al., 2021; Chen et al., 2021; Lai et al., 2022) and repair (Yasunaga & Liang, 2021) tasks, however, these have been limited to only Python. Additional works have introduced generation (Li et al., 2022) and translation (Roziere et al., 2020) tasks in multiple-languages, but are limited to only C++, Java, and Python. We acknowledge concurrent works by Cassano et al. (2022) and Athiwaratkun et al. (2023) on translating HumanEval and MBPP into multiple programming languages. As we note in subsection 2.2, BabelCode supports deeper analysis on a wider range of tasks while including significant methods for ensuring correctness.

Code LLMs Recent years has seen significant interest in LLMs for code. CodeBERT (Feng et al., 2020) is the first work to train an encoder only model on code. CodeT5 (Wang et al., 2021), PLBART (Ahmad et al., 2021), and additional works (Clement et al., 2020; Orlanski & Gittens, 2021; Chakraborty et al., 2022) examine training encoder-decoder models on code. Similar to this work, Ahmad et al. (2021) investigate difference data balancing strategies for pre-training. Our work differs in that we focus on balancing many programming languages in pre-training data. AlphaCode (Li et al., 2022), Codex (Chen et al., 2021), PaLM (Chowdhery et al., 2022), and other works (Nijkamp et al., 2022; Fried et al., 2022; Allal et al., 2023; Christopoulou et al., 2022) have shown that decoder-only code language models achieve exceptional performance on a wide range of tasks. Additional works have investigated different training strategies (Roziere et al., 2020; Bavarian et al., 2022) and different pre-training data (Rozière et al., 2021; Orlanski et al., 2022; Austin et al., 2021).

Language Balancing Choosing a proper sampling distribution from a mixture of datasets of various size is a difficult problems. Initial attempts at studying this in the multilingual natural language processing literature relied on temperature-based approaches (Conneau et al., 2019; Arivazhagan et al., 2019). These approaches oversample the low-resource tasks and downsample the high-resource ones. Other works have

Figure 7. Results on BC-HumanEval and BC-TP3 at a prediction level. Left to right: 1) The % of predictions that passed at least one test, but not all 2) The average, per question, percent of tests passed for each prediction 3) The % of predictions that had either a compilation error, runtime error, or timed out. Full results for BC-HumanEval and BC-TP3 can be found in Figure 9 and Figure 10, respectively.



adopted more dynamic approaches, adapting the sampling rates in an online fashion during training (Wang et al., 2020).

7. Conclusion

We proposed the BabelCode framework for multi-lingual execution-based evaluation and a new strategy for balancing programming language distributions. We highlight the ease of creating new benchmarks with BabelCode by proposing the Translating Python Programming Puzzles. Our experiments demonstrate that adjusting how much we oversample low-resource languages and downsample high-resource languages greatly improves low-resource performance with minimal impact to the performance of high-resource languages in tasks involving either a single or multiple programming language. By open-sourcing BabelCode, future work can investigate improved balancing strategies along with new multi-lingual programming language questions.

Acknowledgements

We thank Michael Janner, Owen Lewis, Alex Polozov, Uros Popovic, Devjeet Roy, Tal Schuster, and Charles Sutton for their helpful discussions and feedback on the paper.

References

Ahmad, W., Chakraborty, S., Ray, B., and Chang, K.-W. Unified pre-training for program understanding and generation. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 2655–2668, Online, June 2021. Association for Computational Linguistics. URL <https://www.aclweb.org/>

[org/anthology/2021.naacl-main.211](https://www.aclweb.org/anthology/2021.naacl-main.211).

Allal, L. B., Li, R., Kocetkov, D., Mou, C., Akiki, C., Ferrandis, C. M., Muennighoff, N., Mishra, M., Gu, A., Dey, M., et al. Santacoder: don’t reach for the stars! *arXiv preprint arXiv:2301.03988*, 2023.

Allamanis, M. The adverse effects of code duplication in machine learning models of code. In *Proceedings of the 2019 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*, pp. 143–153, 2019.

Arivazhagan, N., Bapna, A., Firat, O., Lepikhin, D., Johnson, M., Krikun, M., Chen, M. X., Cao, Y., Foster, G., Cherry, C., et al. Massively multilingual neural machine translation in the wild: Findings and challenges. *arXiv preprint arXiv:1907.05019*, 2019.

Athiwaratkun, B., Gouda, S. K., Wang, Z., Li, X., Tian, Y., Tan, M., Ahmad, W. U., Wang, S., Sun, Q., Shang, M., Gougonadla, S. K., Ding, H., Kumar, V., Fulton, N., Farahani, A., Jain, S., Giaquinto, R., Qian, H., Ramanathan, M. K., Nallapati, R., Ray, B., Bhatia, P., Sengupta, S., Roth, D., and Xiang, B. Multi-lingual evaluation of code generation models. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=Bo7eeXm6An8>.

Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., et al. Program synthesis with large language models. *arXiv preprint arXiv:2108.07732*, 2021.

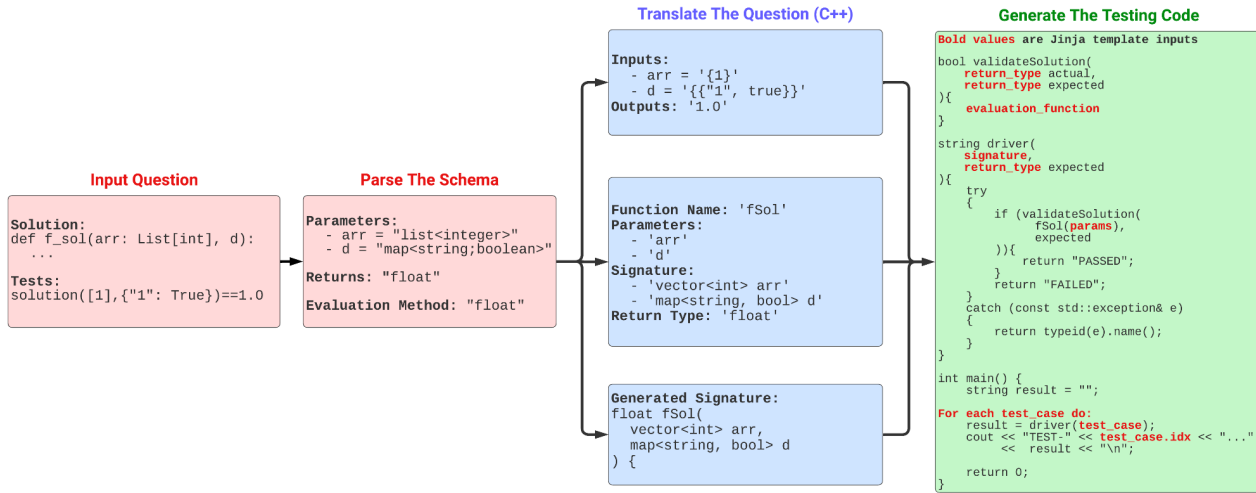
Bavarian, M., Jun, H., Tezak, N., Schulman, J., McLeavey, C., Tworek, J., and Chen, M. Efficient training of

- language models to fill in the middle. *arXiv preprint arXiv:2207.14255*, 2022.
- Cassano, F., Gouwar, J., Nguyen, D., Nguyen, S., Phipps-Costin, L., Pinckney, D., Yee, M. H., Zi, Y., Anderson, C. J., Feldman, M. Q., et al. A scalable and extensible approach to benchmarking nl2code for 18 programming languages. *arXiv preprint arXiv:2208.08227*, 2022.
- Chakraborty, S., Ahmed, T., Ding, Y., Devanbu, P., and Ray, B. Natgen: generative pre-training by “naturalizing” source code. *Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2022.
- Chen, M., Tworek, J., Jun, H., Yuan, Q., Ponde, H., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F. P., Cummings, D. W., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W. H., Nichol, A., Babuschkin, I., Balaji, S. A., Jain, S., Carr, A., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M. M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., and Zaremba, W. Evaluating large language models trained on code. *ArXiv*, abs/2107.03374, 2021.
- Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- Christopoulou, F., Lampouras, G., Gritta, M., Zhang, G., Guo, Y., Li, Z.-Y., Zhang, Q., Xiao, M., Shen, B., Li, L., Yu, H., Yu Yan, L., Zhou, P., Wang, X., Ma, Y., Iacobacci, I., Wang, Y., Liang, G., Wei, J., Jiang, X., Wang, Q., and Liu, Q. Pangu-coder: Program synthesis with function-level language modeling. *ArXiv*, abs/2207.11280, 2022.
- Chung, H. W., Garcia, X., Roberts, A., Tay, Y., Firat, O., Narang, S., and Constant, N. Unimax: Fairer and more effective language sampling for large-scale multilingual pretraining. In *The Eleventh International Conference on Learning Representations*, 2023. URL <https://openreview.net/forum?id=kXwdL1cWOAi>.
- Clement, C., Drain, D., Timcheck, J., Svyatkovskiy, A., and Sundaresan, N. PyMT5: multi-mode translation of natural language and python code with transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 9052–9065, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.728. URL <https://aclanthology.org/2020.emnlp-main.728>.
- Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., Grave, E., Ott, M., Zettlemoyer, L., and Stoyanov, V. Unsupervised cross-lingual representation learning at scale. In *Annual Meeting of the Association for Computational Linguistics*, 2019.
- Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., and Zhou, M. CodeBERT: A pre-trained model for programming and natural languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 1536–1547, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.findings-emnlp.139. URL <https://aclanthology.org/2020.findings-emnlp.139>.
- Fried, D., Aghajanyan, A., Lin, J., Wang, S. I., Wallace, E., Shi, F., Zhong, R., tau Yih, W., Zettlemoyer, L., and Lewis, M. Incoder: A generative model for code infilling and synthesis. *ArXiv*, abs/2204.05999, 2022.
- Hendrycks, D., Basart, S., Kadavath, S., Mazeika, M., Arora, A., Guo, E., Burns, C., Puranik, S., He, H., Song, D., and Steinhardt, J. Measuring coding challenge competence with APPS. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021. URL <https://openreview.net/forum?id=sD93G0zH3i5>.
- Husain, H., Wu, H., Gazit, T., Allamanis, M., and Brockschmidt, M. Codesearchnet challenge: Evaluating the state of semantic code search. *ArXiv*, abs/1909.09436, 2019.
- Kocetkov, D., Li, R., Allal, L. B., Li, J., Mou, C., Ferrandis, C. M., Jernite, Y., Mitchell, M., Hughes, S., Wolf, T., et al. The stack: 3 tb of permissively licensed source code. *arXiv preprint arXiv:2211.15533*, 2022.
- Kudo, T. and Richardson, J. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*, 2018.
- Lai, Y., Li, C., Wang, Y., Zhang, T., Zhong, R., Zettlemoyer, L., Yih, S., Fried, D., yi Wang, S., and Yu, T. Ds-1000: A natural and reliable benchmark for data science code generation. *ArXiv*, abs/2211.11501, 2022.
- Li, Y., Choi, D. H., Chung, J., Kushman, N., Schrittwieser, J., Leblond, R., Tom, Eccles, Keeling, J., Gimeno, F., Lago, A. D., Hubert, T., Choy, P., de, C., d’Autume, M., Babuschkin, I., Chen, X., Huang, P.-S., Welbl, J., Goyal, S., Alexey, Cherepanov, Molloy, J., Mankowitz, D. J., Robson, E. S., Kohli, P., de, N., Freitas, Kavukcuoglu, K., and Vinyals, O. Competition-level code generation with alphacode. *Science*, 378:1092 – 1097, 2022.

- Lu, S., Guo, D., Ren, S., Huang, J., Svyatkovskiy, A., Blanco, A., Clement, C. B., Drain, D., Jiang, D., Tang, D., Li, G., Zhou, L., Shou, L., Zhou, L., Tufano, M., Gong, M., Zhou, M., Duan, N., Sundaresan, N., Deng, S. K., Fu, S., and Liu, S. Codexglue: A machine learning benchmark dataset for code understanding and generation. *ArXiv*, abs/2102.04664, 2021.
- Nijkamp, E., Pang, B., Hayashi, H., Tu, L., Wang, H., Zhou, Y., Savarese, S., and Xiong, C. A conversational paradigm for program synthesis. *arXiv preprint arXiv:2203.13474*, 2022.
- Orlanski, G. and Gittens, A. Reading stackoverflow encourages cheating: Adding question text improves extractive code generation. *ArXiv*, abs/2106.04447, 2021.
- Orlanski, G., Yang, S., and Healy, M. Evaluating how fine-tuning on bimodal data effects code generation. *ArXiv*, abs/2211.07842, 2022.
- Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P. J., et al. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21(140):1–67, 2020.
- Roberts, A., Chung, H. W., Levskaya, A., Mishra, G., Bradbury, J., Andor, D., Narang, S., Lester, B., Gaffney, C., Mohiuddin, A., Hawthorne, C., Lewkowycz, A., Salcianu, A., van Zee, M., Austin, J., Goodman, S., Soares, L. B., Hu, H., Tsvyashchenko, S., Chowdhery, A., Bastings, J., Bulian, J., Garcia, X., Ni, J., Chen, A., Kenealy, K., Clark, J. H., Lee, S., Garrette, D., Lee-Thorp, J., Raffel, C., Shazeer, N., Ritter, M., Bosma, M., Passos, A., Maitin-Shepard, J., Fiedel, N., Omernick, M., Saeta, B., Sepassi, R., Spiridonov, A., Newlan, J., and Gsmundo, A. Scaling up models and data with t5x and `seqio`. *arXiv preprint arXiv:2203.17189*, 2022. URL <https://arxiv.org/abs/2203.17189>.
- Roziere, B., Lachaux, M.-A., Chatussot, L., and Lample, G. Unsupervised translation of programming languages. *Advances in Neural Information Processing Systems*, 33, 2020.
- Rozière, B., Lachaux, M.-A., Szafraniec, M., and Lample, G. Dobf: A deobfuscation pre-training objective for programming languages. In *Neural Information Processing Systems*, 2021.
- Schuster, T., Kalyan, A., Polozov, A., and Kalai, A. T. Programming puzzles. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2021. URL https://openreview.net/forum?id=fe_hCc4RBrq.
- Shazeer, N. and Stern, M. Adafactor: Adaptive learning rates with sublinear memory cost. In *International Conference on Machine Learning*, pp. 4596–4604. PMLR, 2018.
- Tay, Y., Dehghani, M., Tran, V. Q., Garcia, X., Bahri, D., Schuster, T., Zheng, H. S., Houshy, N., and Metzler, D. Unifying language learning paradigms. *arXiv preprint arXiv:2205.05131*, 2022.
- Wang, S., Li, Z., Qian, H., Yang, C., Wang, Z., Shang, M., Kumar, V., Tan, S., Ray, B., Bhatia, P., Nallapati, R., Ramanathan, M. K., Roth, D., and Xiang, B. Recode: Robustness evaluation of code generation models. 2022a.
- Wang, X., Tsvetkov, Y., and Neubig, G. Balancing training for multilingual neural machine translation. *arXiv preprint arXiv:2004.06748*, 2020.
- Wang, Y., Wang, W., Joty, S., and Hoi, S. C. CodeT5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp. 8696–8708, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.emnlp-main.685. URL <https://aclanthology.org/2021.emnlp-main.685>.
- Wang, Z., Cuenca, G., Zhou, S., Xu, F. F., and Neubig, G. Mconala: A benchmark for code generation from multiple natural languages. *ArXiv*, abs/2203.08388, 2022b.
- Yasunaga, M. and Liang, P. Break-it-fix-it: Unsupervised learning for program repair. In *International Conference on Machine Learning (ICML)*, 2021.
- Yin, P., Deng, B., Chen, E., Vasilescu, B., and Neubig, G. Learning to mine aligned code and natural language pairs from stack overflow. *2018 IEEE/ACM 15th International Conference on Mining Software Repositories (MSR)*, pp. 476–486, 2018.

A. BabelCode Design

Figure 8. Sample problem translated from Python to C++ using BabelCode



BabelCode’s design shares many similarities to Athiwaratkun et al. (2023) and Cassano et al. (2022). For translation, we too implement a recursive visitor pattern to translate input and output values to the corresponding code in the target language. When converting a coding dataset, we follow prior works by parsing `assert` statements using AST parsing libraries to determine the inputs and outputs for a given question. To find the function name for a problem, we once again use AST parsers to find the function definition located in the ground truth solution. The found tree is additionally used for parsing the argument names and types. If the types for either the arguments or returns do not exist, we infer them based on the types found from the literal values of the inputs and outputs. While our implementation differs, the overall process is similar to Athiwaratkun et al. (2023) and Cassano et al. (2022). Following Cassano et al. (2022), we execute the generated code through the command line using each language’s recommended commands to compile and run a given script. As Athiwaratkun et al. (2023) is not open sourced, we cannot compare the similarities of this portion.

B. Dataset Changes

B.1. Incompatible Problems

```

1 def encode_cyclic(s: str):
2     """
3     returns encoded string by cycling groups of three characters.
4     """
5     # split string to groups. Each of length 3.
6     groups = [s[(3 * i):min((3 * i + 3), len(s))] for i in range((len(s) + 2) // 3)]
7     # cycle elements in each group. Unless group has fewer elements than 3.
8     groups = [(group[1:] + group[0]) if len(group) == 3 else group for group in groups]
9     return "".join(groups)
10
11
12 def decode_cyclic(s: str):
13     return encode_cyclic(encode_cyclic(s))
14
15 from random import randint, choice
16 import string
17 letters = string.ascii_lowercase
18 for _ in range(100):
19     str = ''.join(choice(letters) for i in range(randint(10, 20)))
20     encoded_str = encode_cyclic(str)
21     assert decode_cyclic(encoded_str) == str

```

B.2. Changes To HumanEval

Original:

```

1 def reverse_delete(s,c):
2     """ Task
3     We are given two strings s and c, you have to deleted all the characters in s that are
4     equal to any character in c
5     then check if the result string is palindrome.
6     A string is called palindrome if it reads the same backward as forward.
7     You should return a tuple containing the result string and True/False for the check.
8     Example
9     For s = "abcde", c = "ae", the result should be ('bcd',False)
10    For s = "abcdef", c = "b" the result should be ('acdef',False)
11    For s = "abcdedcba", c = "ab", the result should be ('cdedc',True)
12    """
13    s = ''.join([char for char in s if char not in c])
14    return (s,s[::-1] == s)
15 assert reverse_delete('abcde', 'ae') == ('bcd', False)
16 assert reverse_delete('abcdef', 'b') == ('acdef', False)
17 assert reverse_delete('abcdedcba', 'ab') == ('cdedc', True)

```

Modified:

```

1 def reverse_delete(s,c):
2     """ Task
3     We are given two strings s and c, you have to deleted all the characters in s that are
4     equal to any character in c
5     then check if the result string is palindrome.
6     A string is called palindrome if it reads the same backward as forward.
7     You should return a two element list containing the result string and "True" if the
8     check passed, otherwise "False".
9     Example
10    For s = "abcde", c = "ae", the result should be ('bcd',False)
11    For s = "abcdef", c = "b" the result should be ('acdef',False)
12    For s = "abcdedcba", c = "ab", the result should be ('cdedc',True)
13    """
14    s = ''.join([char for char in s if char not in c])
15    return [s,str(s[::-1] == s)]
16 assert reverse_delete('abcde', 'ae') == ['bcd', 'False']
17 assert reverse_delete('abcdef', 'b') == ['acdef', 'False']
18 assert reverse_delete('abcdedcba', 'ab') == ['cdedc', 'True']

```

B.3. Changes To Transcoder

Original:

```

1 int difference_between_highest_and_least_frequencies_in_an_array ( int arr [ ], int n ) {
2     sort ( arr, arr + n );
3     int count = 0, max_count = 0, min_count = n;
4     for ( int i = 0;
5     i < ( n - 1 );
6     i ++ ) {
7         if ( arr [ i ] == arr [ i + 1 ] ) {
8             count += 1;
9             continue;
10        }
11        else {
12            max_count = max ( max_count, count );
13            min_count = min ( min_count, count );
14            count = 0;
15        }
16    }

```

```

17 return ( max_count - min_count );
18 }

```

Modified:

```

1 int difference_between_highest_and_least_frequencies_in_an_array(vector<int> arr, int n) {
2   sort(arr.begin(), arr.end());
3   int count = 0, max_count = 0, min_count = n;
4   for ( int i = 0;
5     i < ( n - 1 );
6     i ++ ) {
7     if ( arr [ i ] == arr [ i + 1 ] ) {
8       count += 1;
9       continue;
10    }
11    else {
12      max_count = max ( max_count, count );
13      min_count = min ( min_count, count );
14      count = 0;
15    }
16  }
17  return ( max_count - min_count );
18 }

```

B.4. TP3 Examples

```

1 def sat(inds: List[int], string):
2   return inds == sorted(inds) and ''.join((string[i] for i in inds)) == 'intelligent'
3
4 assert sat([-10, -5, -1, 0, 2, 2, 3, 4, 7, 8, 12], 'enlightenment') == True
5 assert sat([-11, -10, -8, -6, -4, -4, -3, -2, -1, 1, 3], 'inntGetlige') == True
6 assert sat([-10, -5, -1, 0, 2, 2, 3, 4, 7, 8, 12], ' einliJSgeteq ne CAlti') == False

```

C. Training Languages

Table 3. Languages used for training and the extensions we used to filter files. The percentages of the data are calculated after caching and postprocessing using SeqIO.

Language	Extensions	% Of Data
C#	.cs, .cake, .csx, .linq	0.49%
C++	.cpp, .c++, .cc, .cp, .cxx, .h, .h++, .hh, .hpp, .hxx, .inl, .ino, .ipp, .ixx, .re, .tcc, .tpp	16.68%
Dart	.dart	1.85%
Go	.go	3.09%
Haskell	.hs, .hs-boot, .hsc	0.02%
Java	.java, .jav, .jsh	36.95%
JavaScript	.js, .cjs, .mjs	3.31%
Julia	.jl	0.03%
Lua	.lua	1.39%
PHP	.php, .aw, .ctp, .fcgi, .inc, .php3, .php4, .php5, .phps, .phpt	14.05%
Python	.py, .py3, .pyi, .pyw, .pxi	16.80%
R	.r, .rd, .rsx	0.11%
Rust	.rs, .rs.in	0.93%
TypeScript	.ts, .cts, .mts	4.28%

D. Training Objective

This paper uses a variant of the UL2 objective (Tay et al., 2022) for training the code language models. The UL2 objective consists of a mixture of span corruption and prefix language modeling objectives, as defined in Raffel et al. (2020). In this work, we select two span corruption instances using the implementation provided in the T5 library.⁵ The only differences between these two instances consist of different values for the `noise_density` and `mean_noise_span_length` arguments. In particular, we use (3.0, 0.15) and (32, 0.5) for the (`noise_density`, `mean_noise_span_length`) arguments for each span corruption instance respectively.

The prefix language modeling objective randomly breaks text into two pieces, and the model is tasked to reconstruct the latter, given the former. Finally, we add an additional objective which consists of causal language modeling, which can be considered a special case of prefix language modeling; the first piece consists of the empty string. We assign the probabilities 10%, 10%, 20%, and 60% for each objective, respectively.

E. Prompts Used

E.1. Generation Tasks

```
1 You are an expert {{ Language }} programmer, complete the implementation.
2 Solution in {{ Language }}:
3 [BEGIN]
4
5 {{ Signature With Docstring }}
```

Each `{{...}}` represents a field that is filled in.

Example from HumanEval for generating C# code:

```
1 You are an expert C# programmer, complete the implementation.
2 Solution in C#:
3 [BEGIN]
4
5 class Solution {
6     /**
7      * Return length of given string
8      * >>> GetStringLength("")
9      * 0
10     * >>> GetStringLength("abc")
11     * 3
12     */
13     public int GetStringLength(string s) {
```

E.2. Translation Tasks

```
1 Translate the following {{ Source Language }} program to {{ Target Language }}:
2 Input:
3
4 {{ Source Code }}
5
6 {{ Target Language }} Translation:
7 [BEGIN]
8
9 {{ Target Signature }}
```

Each `{{...}}` represents a field that is filled in. The `{{fields}}` correspond to the source language we are translating from, while `{{fields}}` correspond to the target language to translate too.

Example For TP3 translation from Python to Haskell:

```
1 Translate the following Python program to Haskell:
2 Input:
```

⁵See <https://github.com/google-research/text-to-text-transfer-transformer/blob/main/t5/data/preprocessors.py#L1923>

```

3
4 def sat(i: int) -> bool:
5     return i % 123 == 4 and i > 10 ** 10
6
7 Haskell Translation:
8 [BEGIN]
9
10 sat :: Integer -> Bool
11 sat i =
    
```

Figure 9. Qualitative Comparison of the 4B model trained on the Natural, the Unimax 1, and Unimax 2 distributions when evaluated on BC-HumanEval. The results can be found in Table 16 and Table 17.

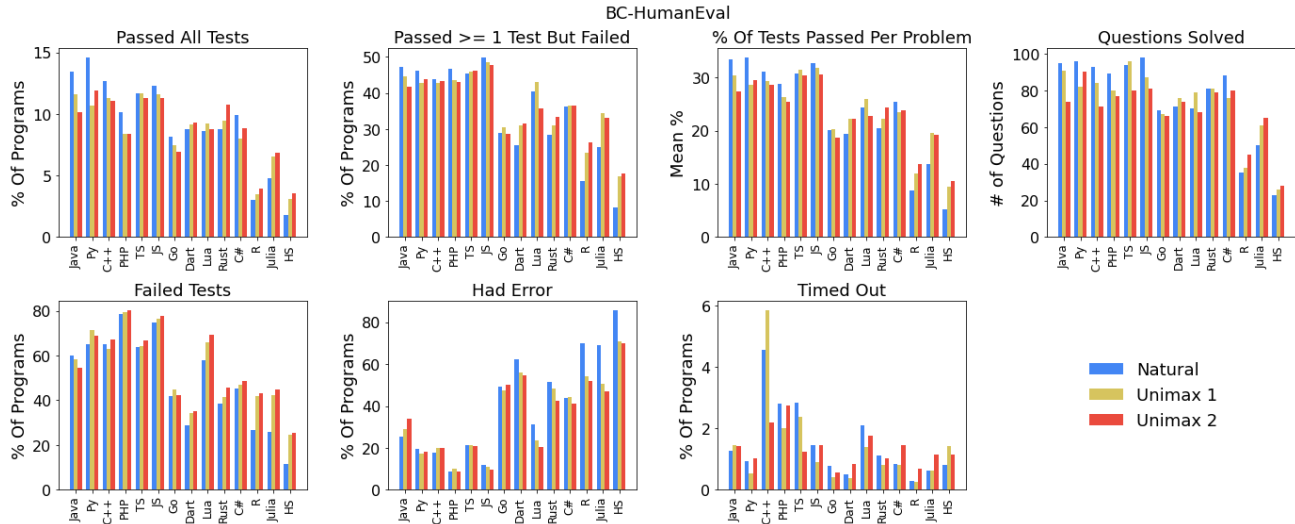
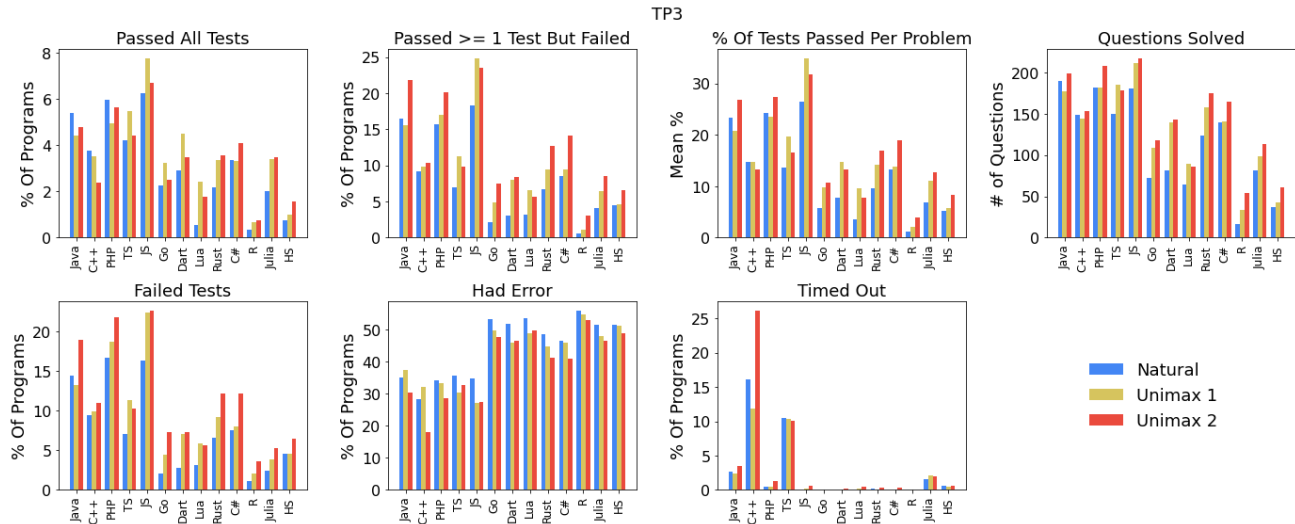


Figure 10. Qualitative Comparison of the 4B model trained on the Natural, the Unimax 1, and Unimax 2 distributions when evaluated on TP3. The results can be found in Table 18 and Table 19.



F. Full Results

Table 4. BC-HumanEval *pass*@1 values for the different models and training distributions. Used $T = 0.8$ and sampled 200 programs per problem. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	Py	R	Rust	TS
1B	Nat	1.0	3.6	2.3	2.5	0.7	3.8	3.6	0.5	1.8	2.8	4.8	0.5	1.6	4.0
	UM 1	1.7	3.0	3.0	2.6	1.3	2.8	4.0	2.1	2.2	2.5	3.9	1.2	2.8	4.4
	UM 2	2.0	3.2	3.0	2.7	1.6	2.7	3.9	2.1	2.1	2.3	4.2	1.4	3.1	4.3
	UM 3	1.6	1.5	2.6	2.6	1.4	2.8	4.0	2.5	2.2	2.2	4.0	1.8	2.6	4.1
	UM 4	1.7	2.7	3.1	2.9	1.5	2.8	3.7	2.6	2.2	2.2	3.5	2.1	2.5	4.1
2B	Nat	2.6	7.5	5.0	5.4	1.0	8.0	7.6	1.2	4.5	6.2	9.1	1.4	3.9	7.9
	UM 1	5.3	6.0	6.1	5.1	1.9	6.6	7.6	4.4	5.4	5.6	7.8	2.1	6.4	7.5
	UM 2	5.2	6.1	5.6	4.5	2.1	5.7	6.4	4.5	5.2	4.8	7.0	2.8	5.8	7.0
	UM 3	5.5	6.2	5.2	4.7	2.4	6.2	6.8	5.1	4.9	4.8	7.5	3.5	6.1	7.0
	UM 4	4.9	6.1	5.4	4.7	2.9	5.7	6.5	4.6	4.8	4.6	7.5	3.3	5.6	7.1
4B	Nat	9.9	12.7	8.7	8.2	1.8	13.5	12.3	4.7	8.6	10.1	14.6	3.0	8.7	11.7
	UM 1	8.0	11.3	9.2	7.5	3.1	11.6	11.6	6.6	9.2	8.4	10.7	3.5	9.5	11.7
	UM 2	8.9	11.1	9.3	7.0	3.6	10.2	11.3	6.8	8.7	8.4	11.9	4.0	10.7	11.3
	UM 3	9.2	9.9	9.0	7.6	4.5	10.5	12.3	8.9	9.2	9.6	11.2	4.5	10.6	11.6
	UM 4	10.4	11.2	8.9	7.7	5.0	10.5	10.6	7.9	9.2	8.0	10.0	5.1	11.0	11.0
8B	PaLM	2.2	3.3	2.5	2.1	0.1	2.5	4.1	0.1	2.2	2.6	3.6	0.2	1.0	4.2
	PaLM-C	2.6	4.4	3.2	3.3	0.3	3.9	5.8	0.1	3.7	4.9	8.1	0.4	1.5	5.6
62B	PaLM	5.9	6.5	3.9	5.3	0.3	6.9	8.5	0.7	6.8	6.2	9.1	1.5	1.8	7.9
	PaLM-C	7.6	9.6	5.7	6.6	0.8	10.4	10.7	1.4	7.5	7.2	11.0	1.9	3.5	9.7

Table 5. BC-TP3 *pass*@1 values for the different models and training distributions. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	R	Rust	TS
1B	Nat	0.5	1.1	0.6	1.0	0.0	1.7	1.2	0.1	0.2	0.6	0.0	0.6	2.1
	UM 1	0.1	0.2	0.1	0.5	0.2	0.3	0.7	0.1	0.1	0.1	0.0	0.3	0.7
	UM 2	0.3	0.1	0.1	0.3	0.1	0.4	1.3	0.3	0.0	0.1	0.0	0.4	1.0
	UM 3	0.2	0.1	0.1	0.3	0.2	0.7	1.1	0.1	0.0	0.0	0.0	0.3	0.7
	UM 4	0.2	0.3	0.4	0.3	0.8	0.6	1.1	0.7	0.1	0.2	0.0	0.7	2.1
2B	Nat	1.0	2.2	1.3	1.9	0.8	2.9	4.1	0.3	0.1	2.8	0.4	2.2	3.1
	UM 1	1.3	0.7	0.7	0.7	0.5	1.9	1.0	0.3	0.3	1.2	0.1	1.1	0.4
	UM 2	1.9	2.1	2.8	0.9	1.0	2.7	6.8	0.6	0.2	4.0	0.1	1.8	5.4
	UM 3	1.1	0.4	0.2	0.4	0.8	1.9	3.6	0.3	0.1	1.7	0.4	0.6	1.0
	UM 4	3.2	1.8	2.4	2.7	1.5	3.7	5.5	2.1	0.5	2.8	0.4	2.9	4.1
4B	Nat	5.9	6.5	5.1	3.9	1.3	9.4	10.9	3.5	0.9	10.4	0.6	3.8	7.3
	UM 1	5.8	6.1	7.8	5.7	1.7	7.7	13.5	5.9	4.2	8.6	1.2	5.8	9.6
	UM 2	7.1	4.1	6.1	4.4	2.7	8.3	11.7	6.1	3.1	9.8	1.3	6.2	7.7
	UM 3	8.7	5.8	7.1	3.6	2.6	7.8	12.1	2.9	1.3	9.5	2.1	6.9	11.1
	UM 4	5.0	4.8	5.7	4.0	1.9	6.8	9.4	2.4	1.3	4.3	2.2	6.3	7.3
8B	PaLM	1.7	4.6	4.9	4.8	0.3	2.6	7.4	0.3	2.9	6.4	0.1	2.2	6.9
	PaLM-C	3.4	5.2	4.8	4.2	0.1	4.7	8.6	0.4	3.6	7.7	0.2	2.4	7.3
62B	PaLM	7.0	7.9	6.6	6.1	1.3	7.9	11.8	1.3	6.2	12.2	1.0	3.6	12.0
	PaLM-C	8.4	8.3	7.6	6.6	1.5	9.9	14.2	1.6	8.0	14.1	2.6	4.0	12.7

Table 6. BC-Transcoder with Python source *pass*@1 values for the different models and training distributions where the source language is Python. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	R	Rust	TS
1B	Nat	1.7	2.1	0.4	1.3	0.2	2.2	2.0	0.1	0.6	0.8	0.2	1.0	2.0
	UM 1	0.1	0.1	0.0	0.4	0.2	0.3	0.5	0.0	0.0	0.1	0.1	0.8	0.8
	UM 2	0.3	0.1	0.1	0.3	0.4	0.6	1.3	0.2	0.0	0.1	0.2	0.7	1.1
	UM 3	0.4	0.3	0.0	0.1	0.3	0.5	0.9	0.1	0.0	0.1	0.2	0.7	0.7
	UM 4	0.3	0.2	0.2	0.1	1.2	0.6	1.1	0.2	0.1	0.4	0.3	0.9	1.3
2B	Nat	2.9	5.5	1.0	4.4	1.0	4.9	8.2	0.3	0.4	3.8	1.3	3.5	5.2
	UM 1	2.9	2.6	0.8	1.2	0.9	3.8	2.5	0.1	0.4	1.5	0.8	1.8	1.0
	UM 2	4.4	5.6	3.9	3.2	1.5	4.9	10.1	1.0	0.3	3.6	2.3	3.2	5.9
	UM 3	2.1	1.0	0.3	0.3	1.4	3.0	3.9	0.0	0.1	1.7	0.5	1.2	1.5
	UM 4	4.8	4.7	2.9	3.5	1.7	4.8	8.4	2.7	2.5	2.6	2.4	4.0	5.7
4B	Nat	23.7	28.4	6.8	11.7	2.3	29.5	27.9	1.7	2.4	23.4	2.8	8.3	15.3
	UM 1	16.7	23.7	9.7	18.6	2.4	18.6	35.3	3.6	8.1	20.8	2.6	12.7	22.4
	UM 2	16.0	16.1	8.4	15.0	3.3	16.6	26.2	5.1	5.3	17.4	5.0	11.3	17.4
	UM 3	21.8	30.6	12.5	14.6	3.5	23.2	37.1	0.9	3.5	20.3	6.1	17.0	28.2
	UM 4	14.5	17.6	3.6	13.0	1.4	14.9	26.6	2.0	4.5	5.0	3.6	14.5	14.7
8B	PaLM	2.9	11.8	4.7	7.3	0.9	4.3	16.3	0.1	5.1	8.8	1.7	3.2	11.6
	PaLM-C	8.5	10.8	5.3	8.6	1.1	8.9	24.2	1.0	9.4	13.7	2.0	4.0	14.3
62B	PaLM	21.4	29.1	7.3	17.8	1.9	17.7	35.6	3.4	16.9	25.6	4.3	7.3	29.3
	PaLM-C	28.7	33.0	9.6	21.4	2.2	23.6	38.4	4.2	22.1	32.4	8.1	7.3	29.6

Table 7. BC-Transcoder with C++ Source *pass*@1 values for the different models and training distributions. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	Dart	Go	HS	Java	JS	Julia	Lua	PHP	Py	R	Rust	TS
1B	Nat	3.0	3.1	1.3	0.1	2.6	2.3	0.1	0.3	1.0	1.8	0.1	0.6	3.7
	UM 1	0.4	3.1	0.9	0.2	1.1	1.1	0.0	0.1	0.8	2.3	0.2	0.8	1.9
	UM 2	1.1	1.6	0.5	0.4	1.1	2.7	0.0	0.0	1.1	1.5	0.1	0.7	1.9
	UM 3	1.9	1.6	0.8	0.4	1.6	2.9	0.0	0.0	0.7	2.3	0.1	1.1	1.5
	UM 4	1.3	3.7	1.8	1.3	1.7	3.0	0.2	0.5	2.4	1.8	0.3	1.6	3.7
2B	Nat	8.9	13.3	5.5	1.2	9.2	16.0	0.3	1.5	12.4	11.2	1.4	4.6	12.1
	UM 1	4.1	7.9	3.4	1.4	6.7	6.2	0.2	2.1	5.0	6.8	0.5	3.6	4.3
	UM 2	8.6	18.1	6.8	2.6	9.4	20.2	0.4	2.1	17.3	8.4	1.2	5.4	15.2
	UM 3	4.6	12.0	4.0	2.1	4.6	12.9	0.5	2.0	8.2	7.7	1.1	2.4	10.3
	UM 4	7.7	15.1	5.9	3.0	7.1	14.4	1.9	2.0	9.6	5.5	1.2	5.4	13.0
4B	Nat	34.5	17.3	20.6	3.2	37.6	32.9	3.3	6.9	34.0	31.7	2.5	10.3	29.2
	UM 1	27.0	18.3	23.5	3.7	27.9	41.2	1.6	9.9	34.5	31.3	2.6	14.3	33.1
	UM 2	19.3	21.1	18.7	4.4	22.0	34.1	4.3	6.4	26.5	25.2	4.0	12.2	24.2
	UM 3	31.5	20.8	16.0	4.6	32.3	42.6	1.0	7.0	39.9	33.5	5.0	16.4	40.2
	UM 4	25.0	15.5	16.4	3.1	21.1	31.9	1.3	6.1	9.7	20.4	2.6	11.6	28.7
8B	PaLM	17.5	16.0	8.3	1.3	14.9	28.1	0.7	8.5	21.2	14.8	1.1	5.0	21.8
	PaLM-C	20.4	15.3	11.2	1.4	20.9	30.8	0.6	12.1	26.5	23.2	1.1	5.3	22.0
62B	PaLM	27.3	17.9	20.6	2.6	24.0	42.4	6.5	16.3	41.3	26.7	4.3	8.6	37.3
	PaLM-C	35.7	17.4	22.1	3.0	30.3	44.3	8.5	19.7	46.6	42.4	8.9	9.4	40.7

Table 8. BC-HumanEval *pass*@100 values for the different models and training distributions. Used $T = 0.8$ and sampled 200 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	Py	R	Rust	TS
1B	Nat	7.3	23.2	14.4	14.9	2.4	24.3	19.0	4.4	9.8	17.1	23.3	4.0	13.9	22.4
	UM 1	12.3	16.2	14.0	12.0	7.5	17.3	18.2	13.0	13.1	15.0	19.9	7.9	14.5	17.8
	UM 2	14.5	16.9	13.8	11.9	8.3	19.6	19.1	15.8	13.5	14.8	21.2	10.4	16.5	19.1
	UM 3	13.7	13.5	13.4	15.4	10.0	21.4	18.4	14.2	12.8	14.6	21.1	10.4	16.0	18.6
	UM 4	15.8	16.6	13.8	12.3	9.7	19.7	18.1	16.6	14.3	15.3	20.6	10.6	15.9	19.6
2B	Nat	17.9	37.8	21.3	27.8	4.9	37.8	36.8	9.7	23.3	35.3	38.8	10.9	26.5	37.9
	UM 1	28.5	31.8	24.6	26.2	12.2	32.0	33.8	23.8	22.9	29.3	30.9	14.0	29.9	34.9
	UM 2	30.6	30.8	25.8	22.6	12.9	32.1	32.1	26.5	21.9	27.4	33.5	15.8	27.4	33.0
	UM 3	31.9	33.0	23.9	25.9	13.7	31.4	34.1	26.5	25.3	29.5	31.5	18.5	28.7	34.8
	UM 4	30.5	30.4	26.7	24.9	12.8	31.3	33.0	29.0	23.2	26.5	34.6	16.2	28.0	34.6
4B	Nat	47.9	51.1	39.6	37.9	12.5	53.4	53.0	27.0	38.7	48.5	52.9	16.7	43.4	50.7
	UM 1	42.4	46.6	42.3	38.3	14.6	50.6	47.9	33.8	42.0	44.0	46.2	20.1	44.6	50.6
	UM 2	44.3	41.2	40.6	34.9	16.0	40.9	44.2	35.9	38.8	42.0	48.9	24.1	43.1	44.6
	UM 3	44.8	44.4	43.3	37.3	21.3	49.9	50.8	40.0	43.2	45.8	48.8	27.9	49.8	51.5
	UM 4	47.9	43.5	37.7	36.1	20.3	46.1	47.3	39.1	42.2	41.7	46.3	23.4	44.8	46.1
8B	PaLM	16.8	19.7	14.7	14.3	1.1	19.9	20.9	2.0	13.2	17.8	21.0	2.9	9.6	22.5
	PaLM-C	27.1	30.1	19.4	20.9	2.5	29.8	31.0	2.4	20.7	29.6	39.5	7.3	13.4	32.5
62B	PaLM	43.9	40.8	26.9	31.4	6.9	48.3	46.2	8.3	36.4	41.6	44.7	13.8	24.3	44.6
	PaLM-C	49.2	50.0	37.6	38.7	9.0	57.0	56.7	12.1	41.1	46.9	64.1	16.9	31.7	54.8

Table 9. BC-TP3 $pass@25$ values for the different models and training distributions where the source language is Python. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	R	Rust	TS
1B	Nat	8.2	16.5	9.6	14.6	0.3	23.4	13.8	1.1	3.7	10.6	0.8	10.6	19.8
	UM 1	2.3	3.2	2.8	9.0	2.6	6.5	9.3	2.0	1.1	1.6	0.2	5.5	9.8
	UM 2	5.2	1.3	1.7	5.4	1.7	8.2	9.8	4.3	1.2	1.1	0.4	6.2	11.3
	UM 3	4.5	2.1	2.3	4.4	3.7	12.6	11.6	1.2	0.4	0.4	0.1	5.4	9.2
	UM 4	3.4	4.6	5.9	5.2	5.4	10.8	11.4	8.6	1.7	4.6	0.7	8.3	17.0
2B	Nat	8.3	18.2	8.3	11.1	5.7	24.5	24.4	3.8	1.4	18.3	3.4	15.3	15.7
	UM 1	15.8	11.9	6.3	8.7	5.2	23.2	10.9	5.4	4.6	11.1	2.1	13.7	4.6
	UM 2	20.7	20.0	19.1	11.7	6.9	26.3	32.9	8.1	3.2	20.6	1.7	18.6	27.8
	UM 3	16.9	7.7	4.4	7.7	5.9	21.2	25.7	5.3	2.3	16.2	4.4	11.2	13.9
	UM 4	24.3	18.8	15.3	14.0	9.6	32.1	28.1	13.9	3.9	17.3	3.5	21.8	21.5
4B	Nat	29.1	31.9	16.6	14.6	7.7	42.2	39.5	17.3	11.9	40.1	3.7	24.1	32.9
	UM 1	28.9	30.0	30.0	22.0	8.8	37.6	49.2	22.5	18.2	40.7	6.9	32.5	41.7
	UM 2	35.5	31.0	30.2	23.7	13.0	43.7	49.5	24.6	17.3	46.1	10.6	37.3	39.0
	UM 3	35.2	24.8	25.5	16.2	13.0	34.3	41.9	16.4	11.9	33.7	10.6	35.2	38.8
	UM 4	25.5	29.7	23.9	19.5	12.1	38.5	40.5	18.6	8.3	26.7	9.8	32.8	29.0
8B	PaLM	19.4	22.6	19.0	17.2	2.8	26.7	26.6	4.0	17.0	31.7	1.9	10.7	25.9
	PaLM-C	25.9	26.2	17.9	16.7	2.0	30.1	34.1	5.9	22.6	40.3	3.2	11.8	29.3
62B	PaLM	38.9	35.2	27.2	24.8	6.1	43.0	48.4	10.6	28.3	48.2	7.2	18.0	42.6
	PaLM-C	41.8	38.7	31.2	26.7	7.2	45.2	55.8	11.3	33.8	56.5	11.4	20.5	48.7

Table 10. BC-Transcoder $pass@25$ values for the different models and training distributions where the source language is Python. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	C++	Dart	Go	HS	Java	JS	Julia	Lua	PHP	R	Rust	TS
1B	Nat	14.0	18.4	5.4	10.3	2.1	17.3	15.3	2.8	7.4	8.7	2.9	10.3	14.4
	UM 1	3.1	1.3	1.1	5.7	3.3	4.9	5.5	0.9	1.0	1.3	1.8	6.7	7.6
	UM 2	4.3	2.7	2.2	3.5	3.9	8.2	11.0	3.3	0.4	2.9	2.6	6.9	10.4
	UM 3	5.8	5.1	0.8	2.1	3.9	6.8	9.4	1.7	0.9	1.2	2.6	6.3	7.2
	UM 4	5.1	3.8	2.6	1.0	6.6	7.6	11.4	2.6	1.3	5.4	3.6	7.8	10.6
2B	Nat	20.9	34.7	11.0	17.5	6.3	30.0	37.0	5.0	5.9	29.4	7.6	11.9	24.3
	UM 1	21.4	22.3	10.8	12.5	5.2	27.7	23.0	2.5	5.6	15.9	6.4	15.2	10.9
	UM 2	29.4	36.1	20.7	20.0	6.9	31.5	43.6	9.1	4.2	28.7	5.8	19.3	29.1
	UM 3	18.6	12.8	4.1	4.9	7.8	22.6	29.0	0.7	1.7	14.9	5.6	13.0	13.9
	UM 4	28.3	29.7	19.1	18.2	9.0	30.0	39.9	12.0	12.2	21.7	7.6	20.1	27.2
4B	Nat	68.4	82.5	34.0	45.5	9.0	80.2	77.6	13.5	23.7	75.9	12.7	38.4	66.1
	UM 1	59.8	75.8	40.2	56.2	11.6	70.5	80.6	16.0	37.9	73.8	11.3	53.9	74.5
	UM 2	58.6	66.7	36.9	57.1	14.2	64.4	76.4	21.2	31.1	69.3	19.7	51.2	68.0
	UM 3	64.6	77.2	39.1	50.2	14.5	73.4	79.0	8.4	24.8	69.1	21.8	58.9	74.8
	UM 4	59.3	72.5	25.4	51.5	11.4	65.0	72.7	13.7	27.1	47.2	19.4	54.2	62.8
8B	PaLM	26.8	48.6	21.0	27.7	3.6	29.7	51.8	2.5	22.4	44.2	5.9	15.0	42.1
	PaLM-C	44.0	52.0	26.7	29.6	5.4	45.9	65.6	9.0	39.7	58.0	9.7	17.5	54.4
62B	PaLM	70.6	78.5	32.7	50.9	8.4	65.1	80.3	15.6	53.4	79.4	17.1	27.5	76.6
	PaLM-C	77.1	83.7	39.8	57.4	8.8	72.2	82.6	20.3	62.2	84.0	23.7	26.6	79.3

Table 11. BC-Transcoder *pass@25* values for the different models and training distributions where the source language is C++. Used $T = 0.8$ and sampled 50 programs per problem. Nat is the natural distribution. UM is Unimax distribution. PaLM-C is the PaLM-Coder distribution. HS is Haskell, JS is JavaScript, Py is Python, and TS is TypeScript.

Size	Dist.	C#	Dart	Go	HS	Java	JS	Julia	Lua	PHP	Py	R	Rust	TS
1B	Nat	24.2	22.4	10.4	1.5	21.2	22.0	1.9	4.7	14.9	15.4	2.3	6.4	27.9
	UM 1	8.1	21.4	8.1	2.6	14.4	12.3	0.6	2.5	12.0	12.4	2.4	7.6	16.2
	UM 2	16.3	18.2	5.7	3.5	13.9	16.0	0.3	0.5	12.8	10.4	1.2	6.7	14.6
	UM 3	21.9	18.5	8.3	3.9	17.5	17.7	0.2	0.2	9.1	13.2	1.6	7.6	16.4
	UM 4	17.0	23.7	10.3	7.0	18.0	17.6	3.4	4.7	18.3	11.9	3.5	8.8	17.9
2B	Nat	38.6	29.7	19.6	6.6	45.2	49.9	5.5	12.5	48.7	40.5	7.5	14.9	45.1
	UM 1	30.9	26.5	19.3	7.6	38.8	33.0	3.2	11.4	35.1	31.8	3.8	16.9	28.0
	UM 2	40.3	27.3	18.0	9.8	46.1	50.4	5.3	10.9	52.8	34.4	4.5	16.1	48.4
	UM 3	34.1	28.4	18.9	9.9	33.9	44.5	6.4	12.2	35.9	34.1	5.3	14.7	40.5
	UM 4	41.3	29.9	25.5	12.2	41.1	49.2	14.4	12.2	41.3	30.1	7.2	19.5	44.5
4B	Nat	71.3	33.2	60.7	10.7	81.9	77.3	20.8	36.3	80.5	79.9	13.3	38.4	76.0
	UM 1	69.6	38.1	63.9	12.7	77.9	77.8	16.1	38.6	76.4	74.7	12.0	52.5	76.5
	UM 2	66.3	37.0	60.8	15.3	73.8	77.6	27.3	33.4	71.0	73.5	18.7	50.7	75.2
	UM 3	75.2	34.8	54.4	14.3	78.3	79.0	12.1	34.8	77.6	76.7	20.9	56.4	79.4
	UM 4	70.7	33.0	59.7	15.0	73.7	74.1	14.3	33.7	61.1	72.8	16.1	47.0	74.7
8B	PaLM	50.5	31.7	32.1	4.5	48.0	60.5	8.5	24.4	62.2	42.3	5.1	15.8	58.2
	PaLM-C	54.8	34.7	37.9	5.3	60.5	68.9	7.8	39.6	68.6	64.3	4.9	20.3	65.3
62B	PaLM	72.0	35.8	55.2	8.4	74.7	77.4	23.0	51.3	82.5	73.0	16.8	25.9	75.4
	PaLM-C	76.2	41.5	58.9	8.4	79.5	80.9	31.6	56.5	84.3	83.1	23.5	27.5	78.4

Table 12. % changes in $pass@k$ compared to the models trained on the natural distribution for High Resource languages. For BC-HumanEval(HE), $k = 100$. For BC-TP3(TP3), BC-Transcoder Python(TC-Py), and BC-Transcoder C++(TC-C++), $k = 25$. The **cells** represent the worst value for that language for that size and dataset. The **cells** represent the best value for that language for that size and dataset.

DS	Size	Dist.	Java	Python	C++	PHP	TS	JS	Go	Mean
HE	1B	UM 1	-29.0	-15.0	-30.3	-12.4	-20.3	-4.1	-19.6	-18.7
		UM 2	-19.6	-9.2	-27.0	-13.4	-14.7	0.2	-20.1	-14.8
		UM 3	-12.1	-9.7	-41.7	-14.6	-16.7	-3.3	3.8	-13.5
		UM 4	-18.9	-11.9	-28.2	-10.8	-12.1	-5.1	-17.2	-14.9
	2B	UM 1	-15.2	-20.3	-15.9	-17.1	-7.9	-8.3	-5.6	-12.9
		UM 2	-15.2	-13.7	-18.6	-22.2	-12.7	-12.7	-18.6	-16.2
		UM 3	-16.9	-18.7	-12.7	-16.3	-8.0	-7.3	-6.7	-12.4
		UM 4	-17.2	-10.6	-19.5	-24.8	-8.6	-10.4	-10.4	-14.5
	4B	UM 1	-5.3	-12.6	-8.9	-9.4	-0.1	-9.5	1.1	-6.4
		UM 2	-23.4	-7.5	-19.5	-13.4	-11.9	-16.4	-8.0	-14.3
		UM 3	-6.6	-7.7	-13.1	-5.7	1.5	-4.0	-1.7	-5.3
		UM 4	-13.7	-12.5	-14.9	-13.9	-9.0	-10.6	-4.7	-11.3
TP3	1B	UM 1	-72.3	N/A	-80.5	-84.8	-50.8	-32.5	-38.6	-59.9
		UM 2	-65.2	N/A	-92.0	-89.8	-43.1	-28.8	-62.7	-63.6
		UM 3	-46.3	N/A	-87.3	-96.2	-53.5	-16.1	-69.8	-61.5
		UM 4	-53.9	N/A	-72.3	-56.9	-14.1	-17.1	-64.3	-46.5
	2B	UM 1	-5.5	N/A	-34.6	-39.4	-70.4	-55.4	-21.2	-37.7
		UM 2	7.3	N/A	9.4	12.4	77.5	35.1	5.4	24.5
		UM 3	-13.4	N/A	-57.6	-11.6	-11.0	5.5	-30.2	-19.7
		UM 4	31.0	N/A	3.1	-5.6	37.3	15.3	26.5	17.9
	4B	UM 1	-10.8	N/A	-5.8	1.4	26.7	24.6	50.8	14.5
		UM 2	3.5	N/A	-2.9	14.9	18.6	25.3	62.5	20.3
		UM 3	-18.6	N/A	-22.2	-16.2	18.0	6.2	11.3	-3.6
		UM 4	-8.7	N/A	-6.9	-33.4	-11.8	2.6	34.2	-4.0
TC-C++	1B	UM 1	-31.9	-19.3	N/A	-19.1	-41.9	-44.3	-22.0	-29.7
		UM 2	-34.6	-32.7	N/A	-13.6	-47.7	-27.6	-45.1	-33.5
		UM 3	-17.4	-14.5	N/A	-38.6	-41.0	-19.9	-20.5	-25.3
		UM 4	-15.2	-22.9	N/A	23.2	-35.7	-20.2	-1.4	-12.0
	2B	UM 1	-14.3	-21.3	N/A	-28.0	-37.8	-33.8	-1.4	-22.8
		UM 2	1.9	-15.0	N/A	8.5	7.3	1.0	-8.0	-0.7
		UM 3	-25.0	-15.7	N/A	-26.2	-10.1	-10.9	-3.3	-15.2
		UM 4	-9.1	-25.6	N/A	-15.2	-1.4	-1.3	30.3	-3.7
	4B	UM 1	-4.9	-6.6	N/A	-5.2	0.7	0.6	5.3	-1.7
		UM 2	-9.9	-8.0	N/A	-11.8	-1.0	0.4	0.1	-5.0
		UM 3	-4.4	-4.0	N/A	-3.6	4.5	2.1	-10.4	-2.6
		UM 4	-10.1	-8.9	N/A	-24.1	-1.7	-4.1	-1.7	-8.4
TC-Py	1B	UM 1	-71.5	N/A	-92.7	-85.4	-47.3	-63.9	-44.9	-67.6
		UM 2	-52.7	N/A	-85.2	-66.8	-27.9	-27.9	-66.1	-54.4
		UM 3	-60.8	N/A	-72.0	-86.5	-49.9	-38.3	-80.0	-64.6
		UM 4	-56.0	N/A	-79.1	-38.4	-26.5	-25.2	-90.3	-52.6
	2B	UM 1	-7.6	N/A	-35.8	-45.7	-55.0	-38.0	-28.9	-35.1
		UM 2	5.3	N/A	4.0	-2.3	20.0	17.6	14.0	9.8
		UM 3	-24.7	N/A	-63.2	-49.4	-42.7	-21.7	-72.0	-45.6
		UM 4	0.0	N/A	-14.6	-25.9	12.0	7.6	3.7	-2.9
	4B	UM 1	-12.1	N/A	-8.1	-2.9	12.6	3.8	23.6	2.8
		UM 2	-19.6	N/A	-19.1	-8.7	2.8	-1.5	25.6	-3.4
		UM 3	-8.4	N/A	-6.4	-9.0	13.1	1.8	10.5	0.3
		UM 4	-19.0	N/A	-12.1	-37.9	-5.0	-6.3	13.4	-11.1

Table 13. % change of $pass@k$ compared to the models trained on the natural distribution for low resource languages. For BC-HumanEval(HE), $k = 100$. For BC-TP3(TP3), BC-Transcoder Python(TC-Py), and BC-Transcoder C++(TC-C++), $k = 25$. The red cells represent the worst value for that language for that size and dataset. The green cells represent the best value for that language for that size and dataset.

DS	Size	Dist.	Dart	Lua	Rust	C#	R	Julia	HS	Mean
HE	1B	UM 1	-2.8	33.8	4.6	68.5	100.0	191.9	205.9	86.0
		UM 2	-4.1	38.3	19.0	98.1	161.7	254.7	238.7	115.2
		UM 3	-6.4	30.8	15.5	87.1	162.7	218.4	308.9	116.7
		UM 4	-3.9	46.5	14.8	115.2	166.9	272.6	294.5	129.5
	2B	UM 1	15.6	-1.6	12.7	59.4	28.6	145.2	147.7	58.2
		UM 2	21.6	-5.9	3.4	71.2	44.9	172.9	161.5	67.1
		UM 3	12.2	8.9	8.2	78.6	68.9	173.5	177.9	75.4
		UM 4	25.6	-0.4	5.6	70.8	48.6	198.7	160.5	72.8
	4B	UM 1	7.0	8.6	3.0	-11.5	20.4	25.3	16.8	9.9
		UM 2	2.6	0.2	-0.6	-7.5	44.0	32.9	27.6	14.2
		UM 3	9.5	11.5	14.7	-6.5	66.9	48.2	70.3	30.7
		UM 4	-4.7	9.0	3.2	-0.1	40.3	44.8	62.2	22.1
TP3	1B	UM 1	-71.1	-70.0	-48.4	-72.0	-70.6	80.5	660.5	58.4
		UM 2	-82.1	-69.1	-41.6	-36.3	-50.0	297.0	389.8	58.2
		UM 3	-75.7	-89.1	-49.1	-44.9	-83.3	9.0	992.0	94.1
		UM 4	-38.4	-53.9	-21.1	-58.6	-16.7	693.3	1504.5	287.0
	2B	UM 1	-23.2	221.6	-10.1	90.3	-38.1	40.2	-9.4	38.8
		UM 2	131.9	128.3	22.0	149.1	-49.5	109.9	21.6	73.3
		UM 3	-46.8	63.5	-26.4	103.5	30.6	39.0	3.5	23.9
		UM 4	85.7	172.4	43.1	192.1	4.2	260.5	68.6	118.1
	4B	UM 1	80.3	53.2	34.9	-0.9	85.6	29.9	13.9	42.4
		UM 2	81.6	45.7	55.0	21.7	187.5	42.4	67.3	71.6
		UM 3	53.3	0.1	46.2	20.7	187.8	-5.2	67.5	52.9
		UM 4	43.9	-29.8	36.1	-12.5	166.7	7.3	56.1	38.3
TC-C++	1B	UM 1	-4.5	-47.2	18.3	-66.7	7.0	-69.8	67.3	-13.6
		UM 2	-18.6	-88.7	3.8	-32.8	-46.4	-84.8	130.4	-19.6
		UM 3	-17.1	-94.9	18.8	-9.5	-28.5	-89.9	157.1	-9.1
		UM 4	6.2	1.5	36.9	-29.7	53.8	82.3	357.7	72.7
	2B	UM 1	-10.8	-8.4	13.8	-20.1	-49.6	-41.9	14.9	-14.6
		UM 2	-8.3	-12.5	8.6	4.2	-40.6	-3.4	47.6	-0.6
		UM 3	-4.5	-2.3	-1.2	-11.9	-29.8	17.7	48.6	2.4
		UM 4	0.5	-2.6	31.0	7.0	-4.2	163.2	84.2	39.9
	4B	UM 1	14.8	6.4	36.6	-2.4	-10.0	-22.9	18.5	5.9
		UM 2	11.4	-7.8	31.9	-6.9	40.6	30.9	42.8	20.4
		UM 3	4.8	-4.1	46.8	5.6	57.5	-42.1	33.8	14.6
		UM 4	-0.8	-7.1	22.4	-0.8	21.2	-31.2	40.5	6.3
TC-Py	1B	UM 1	-79.8	-86.1	-34.9	-78.0	-35.9	-69.9	55.1	-47.1
		UM 2	-60.0	-94.2	-33.3	-69.6	-10.9	14.5	83.7	-24.3
		UM 3	-85.1	-88.4	-38.3	-58.8	-9.0	-39.7	83.0	-33.7
		UM 4	-52.0	-82.1	-24.5	-63.5	25.6	-7.7	210.6	0.9
	2B	UM 1	-1.3	-4.8	27.8	2.3	-16.5	-48.8	-17.0	-8.3
		UM 2	88.9	-28.6	62.2	40.5	-23.4	83.9	9.3	33.3
		UM 3	-62.2	-72.1	9.3	-10.9	-25.9	-86.5	25.1	-31.9
		UM 4	74.6	106.5	69.4	35.2	-0.8	142.2	44.0	67.3
	4B	UM 1	18.3	60.2	40.3	-12.5	-11.2	18.4	28.7	20.3
		UM 2	8.5	31.4	33.1	-14.3	54.9	57.4	58.0	32.7
		UM 3	15.0	4.7	53.4	-5.6	71.0	-38.0	61.0	23.1
		UM 4	-25.4	14.4	41.0	-13.3	52.0	1.3	26.7	13.8

Table 14. Number of Questions passed for BC-HumanEval(HE) and TP3. BC-HE has 161 total problems and TP3 has 370 total problems. S is the size of the model, and D is the distribution it was trained on. P is the PaLM distribution while PC is the PaLM-Coder distribution. Languages are sorted from high to low resource. **Green** values are the best values for that language, while **red** values are the worst.

N	S	D	Java	Py	C++	PHP	TS	JS	Go	Dart	Lua	Rust	C#	R	Julia	HS
HE	1B	N	46	44	44	32	44	38	31	28	18	27	13	8	9	5
		U1	33	38	32	30	34	36	23	27	24	26	26	17	23	13
		U2	38	39	33	28	38	38	21	26	26	32	28	20	30	17
		U3	43	41	25	29	38	37	31	25	24	29	28	20	25	19
		U4	41	40	32	31	39	34	23	26	28	32	32	19	32	18
	2B	N	69	70	70	69	71	70	53	40	43	52	33	21	18	9
		U1	58	56	60	55	64	61	53	46	42	58	53	27	47	21
		U2	60	61	56	54	60	61	43	51	40	51	56	31	50	25
		U3	58	55	64	57	67	62	49	46	49	54	59	35	51	27
		U4	58	64	57	53	66	62	46	51	46	53	58	30	54	25
	4B	N	95	96	93	89	94	98	69	71	70	81	88	35	50	23
		U1	91	82	84	80	96	87	67	76	79	81	76	38	61	26
		U2	74	90	71	77	80	81	66	74	68	79	80	45	65	28
		U3	94	89	80	85	95	92	65	81	77	93	80	53	72	39
		U4	84	82	78	77	84	86	64	67	78	81	88	46	70	38
	8B	P	37	41	39	35	46	41	29	28	26	18	30	6	5	2
		PC	57	74	60	56	65	58	40	37	39	27	55	15	5	6
	62B	P	91	81	76	76	85	85	61	50	68	49	88	26	16	14
		PC	104	119	92	85	105	108	71	72	77	62	92	32	25	17
	TP3	1B	N	122		89	61	102	73	78	55	18	62	45	6	6
U1			41		20	11	52	50	53	17	7	31	15	1	14	14
U2			54		8	6	60	49	32	9	8	38	34	3	26	10
U3			72		14	3	49	58	26	14	3	33	29	1	8	21
U4			62		28	30	84	61	32	34	9	46	22	5	47	29
2B		N	127		94	95	81	127	56	43	10	76	43	16	20	26
		U1	120		66	57	23	56	49	35	25	73	87	10	33	29
		U2	132		105	107	137	158	65	93	19	98	107	9	45	36
		U3	110		48	89	77	124	48	26	14	66	95	23	32	32
		U4	153		99	84	104	133	73	77	18	110	119	17	67	49
4B		N	190		149	182	150	181	72	81	64	123	140	16	81	37
		U1	177		144	182	185	211	109	139	89	158	141	33	99	43
		U2	199		153	208	178	217	118	143	86	175	165	54	113	61
		U3	162		120	162	176	189	77	119	60	167	169	50	80	64
		U4	181		143	126	134	188	95	114	41	156	123	43	87	60
8B		P	130		106	149	123	121	85	93	88	53	100	9	20	14
		PC	148		126	182	140	161	80	86	109	61	129	17	32	11
62B		P	189		161	213	192	218	115	132	129	88	181	31	49	26
		PC	204		175	247	218	243	124	145	156	100	192	50	51	33

Table 15. Number of Questions passed for Transcoder. There are a total of 524 questions, and N represents the source language. S is the size of the model, and D is the distribution it was trained on. P is the PaLM distribution while PC is the PaLM-Coder distribution. Languages are sorted from high to low resource. **Green** values are the best values for that language, while **red** values are the worst.

N	S	D	Java	Py	C++	PHP	TS	JS	Go	Dart	Lua	Rust	C#	R	Julia	HS
Py	1B	N	118		124	62	96	103	70	41	52	68	99	20	24	15
		U1	41		13	11	51	40	41	10	10	43	28	14	9	25
		U2	62		25	25	73	78	28	19	4	46	33	17	27	26
		U3	51		43	11	52	66	17	8	9	43	46	19	17	30
		U4	58		34	43	70	81	8	21	12	51	40	25	22	42
	2B	N	191		225	197	160	231	114	76	47	78	140	46	36	41
		U1	182		154	115	80	160	89	78	43	102	144	41	22	36
		U2	205		233	190	188	271	133	130	33	132	192	35	61	45
		U3	152		100	103	98	196	42	33	14	94	132	40	7	50
		U4	195		196	152	172	248	119	123	73	134	185	46	70	59
	4B	N	449		457	434	388	437	272	206	161	244	384	80	85	56
		U1	408		427	420	424	445	327	237	239	330	354	70	100	75
		U2	380		385	402	396	429	337	222	202	307	344	121	133	90
		U3	417		430	397	417	431	300	229	159	347	369	132	54	95
		U4	383		412	304	367	409	306	161	174	321	346	119	84	82
	8B	P	192		291	270	246	301	168	134	143	99	191	35	22	25
		PC	280		314	336	324	371	175	169	233	115	267	62	57	34
	62B	P	379		438	441	429	444	303	199	308	171	400	101	99	56
		PC	421		459	463	442	457	332	237	359	157	432	142	127	55
	C++	1B	N	143	100		112	182	143	71	137	33	50	163	18	18
U1			104	78		92	104	82	49	125	20	50	66	18	5	20
U2			98	64		88	95	102	39	120	5	45	122	10	3	25
U3			121	82		65	112	112	57	123	2	48	162	12	2	28
U4			120	77		123	112	112	65	143	32	56	121	25	25	48
2B		N	278	245		295	269	285	127	171	86	97	226	48	41	42
		U1	242	202		224	183	207	129	153	75	111	196	26	25	51
		U2	285	218		311	282	299	121	153	68	105	244	31	37	63
		U3	225	218		224	239	264	124	161	80	96	213	35	47	64
		U4	260	190		247	263	288	163	174	78	131	255	46	94	80
4B		N	448	446		446	423	433	348	194	217	235	393	81	133	65
		U1	437	410		422	419	425	365	224	234	315	391	73	112	79
		U2	424	416		396	418	428	349	213	213	308	382	117	168	92
		U3	435	434		428	433	431	322	202	212	334	418	129	86	88
		U4	415	414		363	412	413	350	188	216	285	399	104	103	95
8B		P	283	253		352	328	335	191	176	151	94	288	34	58	28
		PC	350	370		379	367	382	228	197	236	126	319	37	58	33
62B		P	424	407		451	411	420	323	202	300	157	405	100	141	49
		PC	441	462		454	427	441	336	240	326	163	420	137	191	49

Table 16. Metrics for HR languages on BC-HumanEval for all models. Δ is the mean change of each of the displayed languages when compared to the natural. % Failed tests is the percent of predictions that did not have any errors, but failed a test. % Error is the percent of predictions that had either a runtime or compilation error. % Timed Out is the percent of predictions that timed out. The time out was set to 10 for all languages except for Java and TS, which was 15. % Passed is the percent of predictions that passed all test cases. % Passed One is the percent of predictions that passed at least one test case, but failed. % Tests Passed is the mean percent of test cases passed per problem for all predictions.

Metric	<i>D</i>	Java	Py	C++	PHP	TS	JS	Go	Δ
% Error	N	25.32	19.36	17.80	8.61	21.53	11.66	49.02	
	U1	28.85	17.45	19.83	10.25	21.53	11.00	47.23	0.40
	U2	34.08	18.16	19.80	8.65	20.87	9.67	50.13	1.15
% Failed Test	N	59.94	65.12	64.94	78.45	63.92	74.60	42.02	
	U1	58.12	71.33	63.03	79.32	64.41	76.44	44.85	1.22
	U2	54.34	68.89	66.97	80.25	66.59	77.56	42.34	1.13
% Passed	N	13.45	14.60	12.70	10.12	11.71	12.29	8.15	
	U1	11.57	10.68	11.29	8.41	11.69	11.64	7.50	-1.46
	U2	10.16	11.93	11.05	8.37	11.29	11.30	6.96	-1.71
% Passed One	N	47.26	46.20	43.77	46.70	45.32	49.87	28.82	
	U1	44.68	42.83	42.80	43.60	45.95	48.47	30.39	-1.32
	U2	41.69	43.92	43.38	43.02	46.13	47.87	28.69	-1.89
% Tests Passed	N	33.46	33.77	31.07	28.84	30.71	32.78	20.03	
	U1	30.45	28.69	29.25	26.29	31.42	31.78	20.21	-1.79
	U2	27.44	29.58	28.64	25.49	30.44	30.61	18.75	-2.81
% Timed Out	N	1.29	0.93	4.57	2.82	2.84	1.45	0.80	
	U1	1.45	0.54	5.86	2.02	2.37	0.92	0.42	-0.16
	U2	1.42	1.02	2.18	2.74	1.25	1.47	0.57	-0.58

Table 17. Metrics for LR languages on BC-HumanEval for all models. Δ is the mean change of each of the displayed languages when compared to the natural. % Failed tests is the percent of predictions that did not have any errors, but failed a test. % Error is the percent of predictions that had either a runtime or compilation error. % Timed Out is the percent of predictions that timed out. The time out was set to 10 for all languages except for Java and TS, which was 15. % Passed is the percent of predictions that passed all test cases. % Passed One is the percent of predictions that passed at least one test case, but failed. % Tests Passed is the mean percent of test cases passed per problem for all predictions.

Metric	<i>D</i>	Dart	Lua	Rust	C#	R	Julia	HS	Δ
% Error	N	62.06	31.31	51.61	43.80	70.08	68.90	85.70	
	U1	56.05	23.39	48.20	44.40	54.24	50.51	70.80	-9.41
	U2	54.64	20.28	42.62	41.11	52.07	47.10	69.75	-12.27
% Failed Test	N	28.71	57.98	38.51	45.42	26.66	25.72	11.67	
	U1	34.37	66.01	41.50	46.84	42.05	42.28	24.69	9.01
	U2	35.26	69.21	45.62	48.56	43.26	44.92	25.52	11.10
% Passed	N	8.74	8.60	8.74	9.94	2.99	4.75	1.81	
	U1	9.19	9.23	9.47	7.97	3.46	6.57	3.08	0.49
	U2	9.27	8.74	10.73	8.86	3.98	6.83	3.57	0.92
% Passed One	N	25.51	40.39	28.48	36.26	15.62	25.07	8.29	
	U1	30.95	42.98	30.94	36.57	23.29	34.48	16.90	5.21
	U2	31.58	35.78	33.45	36.41	26.42	33.11	17.55	4.95
% Tests Passed	N	19.43	24.31	20.53	25.49	8.70	13.79	5.13	
	U1	22.31	26.00	22.31	23.45	12.03	19.59	9.53	2.55
	U2	22.32	22.79	24.36	23.92	13.68	19.25	10.48	2.77
% Timed Out	N	0.49	2.10	1.13	0.85	0.28	0.63	0.82	
	U1	0.39	1.38	0.82	0.80	0.25	0.64	1.43	-0.09
	U2	0.83	1.78	1.04	1.46	0.69	1.14	1.16	0.26

Table 18. Metrics for HR languages on TP3 for all models. Δ is the mean change of each of the displayed languages when compared to the natural. % Failed tests is the percent of predictions that did not have any errors, but failed a test. % Error is the percent of predictions that had either a runtime or compilation error. % Timed Out is the percent of predictions that timed out. The time out was set to 10 for all languages except for Java and TS, which was 15. % Passed is the percent of predictions that passed all test cases. % Passed One is the percent of predictions that passed at least one test case, but failed. % Tests Passed is the mean percent of test cases passed per problem for all predictions.

Metric	<i>D</i>	Java	C++	PHP	TS	JS	Go	Δ
% Error	N	60.94	49.05	59.66	62.13	60.44	92.53	
	U1	65.04	56.15	58.08	52.67	47.27	86.54	-3.17
	U2	52.71	31.20	50.03	56.74	47.73	82.85	-10.58
% Failed Test	N	25.09	16.37	29.14	12.31	28.45	3.54	
	U1	23.17	17.19	32.55	19.78	38.94	7.71	4.07
	U2	32.95	19.17	37.92	17.92	39.45	12.63	7.52
% Passed	N	9.40	6.54	10.39	7.33	10.91	3.91	
	U1	7.67	6.10	8.62	9.56	13.52	5.66	0.44
	U2	8.30	4.12	9.80	7.73	11.68	4.36	-0.42
% Passed One	N	28.57	15.97	27.25	12.20	31.76	3.77	
	U1	27.17	16.99	29.46	19.59	43.19	8.45	4.22
	U2	37.95	17.98	34.96	17.14	40.92	12.97	7.07
% Tests Passed	N	23.31	14.69	24.17	13.66	26.44	5.81	
	U1	20.82	14.76	23.62	19.67	34.89	9.80	2.58
	U2	26.81	13.21	27.44	16.55	31.70	10.72	3.06
% Timed Out	N	4.57	28.03	0.81	18.23	0.20	0.02	
	U1	4.13	20.56	0.76	17.98	0.28	0.09	-1.34
	U2	6.04	45.51	2.24	17.61	1.14	0.16	3.47

Table 19. Metrics for LR languages on TP3 for all models. Δ is the mean change of each of the displayed languages when compared to the natural. % Failed tests is the percent of predictions that did not have any errors, but failed a test. % Error is the percent of predictions that had either a runtime or compilation error. % Timed Out is the percent of predictions that timed out. The time out was set to 10 for all languages except for Java and TS, which was 15. % Passed is the percent of predictions that passed all test cases. % Passed One is the percent of predictions that passed at least one test case, but failed. % Tests Passed is the mean percent of test cases passed per problem for all predictions.

Metric	<i>D</i>	Dart	Lua	Rust	C#	R	Julia	HS	Δ
% Error	N	90.12	93.45	84.47	80.85	97.34	89.43	89.60	
	U1	79.83	85.22	77.93	80.02	95.20	83.65	89.39	-4.86
	U2	80.96	86.36	72.00	71.00	92.17	81.19	84.96	-8.09
% Failed Test	N	4.78	5.42	11.54	13.10	2.00	4.23	7.96	
	U1	12.24	10.25	16.07	14.06	3.60	6.73	7.98	3.13
	U2	12.73	9.82	21.21	21.30	6.38	9.20	11.34	6.13
% Passed	N	5.07	0.94	3.77	5.87	0.62	3.51	1.31	
	U1	7.83	4.22	5.84	5.76	1.20	5.92	1.74	1.63
	U2	6.09	3.11	6.18	7.14	1.32	6.10	2.75	1.65
% Passed One	N	5.28	5.51	11.77	14.87	0.95	7.15	7.83	
	U1	13.96	11.43	16.51	16.31	1.83	11.11	7.97	3.68
	U2	14.57	9.88	22.00	24.70	5.34	14.85	11.45	7.06
% Tests Passed	N	7.76	3.55	9.59	13.23	1.10	6.87	5.18	
	U1	14.74	9.62	14.10	13.74	2.11	11.03	5.77	3.40
	U2	13.33	7.78	17.00	19.01	3.92	12.77	8.40	4.99
% Timed Out	N	0.02	0.20	0.22	0.18	0.03	2.82	1.12	
	U1	0.11	0.31	0.16	0.16	0.01	3.70	0.89	0.11
	U2	0.22	0.71	0.60	0.56	0.13	3.51	0.96	0.30