

# Introducing Compiler Semantics into Large Language Models as Programming Language Translators: A Case Study of C to x86 Assembly

Anonymous ACL submission

## Abstract

Compilers are complex software containing millions of lines of code, taking years to develop. This paper investigates to what extent Large Language Models (LLMs) can replace hand-crafted compilers in translating high-level programming languages to machine instructions, using C to x86 assembly as a case study. We identify two challenges of using LLMs for code translation and introduce two novel data pre-processing techniques to address the challenges: numerical value conversion and training data resampling. While only using a 13B model, our approach achieves a behavioral accuracy of over 91%, outperforming the much larger GPT-4 Turbo model by over 50%. Our results are encouraging, showing that LLMs have the potential to transform how compilation tools are constructed.

## 1 Introduction

There is growing interest in using Large Language Models (LLMs) for software engineering tasks (Zhang et al., 2023b) like code retrieval (Li et al., 2022b,a), completion (Svyatkovskiy et al., 2020; Guo et al., 2023) and translation (Armengol-Estapé and O’Boyle, 2021; Armengol-Estapé et al., 2023). The training data of many LLMs, including CodeLlama (Rozière et al., 2022), Codex (Chen et al., 2021), and GPT4 (OpenAI et al., 2023) contains code examples. However, these models are not explicitly trained for code translation. Consequently, they are prone to errors during code translation (Armengol-Estapé et al., 2023). On the other hand, LLMs trained in natural language corpus have demonstrated impressive results in natural language understanding (Brown et al., 2020; Puk-sachatkun et al., 2020). As such, it is interesting to know if LLMs can learn to compile code.

This paper investigates the feasibility of using LLMs to translate a high-level programming language to machine instructions, a problem

known as *neural compilation* (Armengol-Estapé and O’Boyle, 2021). Traditionally, this is performed by a manually crafted compiler that usually takes many person-years of compiler engineers’ time to build. Recent developments in LLMs have shown promising results in leveraging pre-trained transformer models for tasks like decompilation (e.g., translating assembly code to C programs) (Armengol-Estapé et al., 2023) and program synthesis (Szafraniec et al., 2023). However, few works use LLMs as a compilation tool to translate a high-level programming language into low-level assembly instructions. Our work seeks to bridge this gap by taking C to x86 assembly as a case study.

A key challenge we face is managing the semantic gap between high-level languages optimized for human usability and low-level languages designed for hardware executions. This gap often manifests in a lack of direct correspondence between elements of the source and target languages. For instance, some commonly used data structures and programming constructs in C, such as *struct* and complex *for-loop*, do not have single equivalent x86 instructions. Similarly, C uses identifiers for variables, while assembly instructions use stack and memory addresses or registers. As a single line of C code can be translated into a varying number of assembly instructions, learning the translation from C to assembly would require different amounts of training samples depending on the complexity of the mapping, making it difficult to construct a balanced training corpus.

To overcome the aforementioned challenges, we leverage Low-Rank Adaptation (LoRA) (Hu et al., 2021) to fine-tune a pre-trained 13B CodeLlama model (Rozière et al., 2022). However, using the standard natural language training pipeline, our initial attempt yields a model with poor performance for C-to-assembly translations. After a close examination of the failure cases, we propose to introduce

082 compiler semantics as two key data pre-processing  
083 techniques to enhance the trained model: symbolic  
084 interpretation for numerical value conversion and  
085 switch-case normalization for switch-case inconsis-  
086 tency. Furthermore, we propose an automatic com-  
087 piler semantics guided refinement learning frame-  
088 work to improve the fine-tuned model iteratively.  
089 Our framework automatically resamples the distri-  
090 bution of semantic mapping samples and synthe-  
091 sizes the failure test cases in the validation set to  
092 improve the quality of the model training data.

093 We perform a large-scale evaluation on over 57k  
094 executable C programs and compare them against  
095 the state-of-the-art large language model GPT-4  
096 Turbo. We verify the correctness of the generated  
097 x86 assembly code by executing them against unit  
098 test cases. Experimental results show that our neu-  
099 ral compiler generates code that is more accurate  
100 than all competing baselines. Compared to GPT-4  
101 Turbo, our approach improves the translation accu-  
102 racy by over 50%, from 40.85% to 91.88%.

103 Our main contributions are:

- 104 • We propose an approach to introduce compiler  
105 semantics into the LLM as two new data pre-  
106 processing methods: symbolic interpretation  
107 and switch-case normalization. Experimen-  
108 tal results demonstrate that the two proposed  
109 methods allow the LLM to increase the num-  
110 ber of correct translations by over 30%.
- 111 • We implement an automatic refinement aug-  
112 mentation framework targeting the biased  
113 samples of different semantics in the corpora,  
114 where the long-tails under-fit. The framework  
115 resamples the semantics distribution by syn-  
116 thesizing incorrect cases, to obtain improved  
117 accuracy on the long tails.
- 118 • We can achieve 91.88% IO accuracy when  
119 translating C to x86 assembly and we believe  
120 it’s the highest accuracy when comparing with  
121 SOTA works.

## 122 2 Problem Statement

123 We target the problem of machine translating high-  
124 level programs (specifically, in the C language) into  
125 semantically equivalent low-level programs (in x86  
126 assembly) with limited bilingual parallel corpora.  
127 One approach is to use compilers, like GCC, as  
128 the oracle to generate semantic aligned assembly  
129 from C language corpora. We model the problem  
130 as follows:

**Definition 1** *There is a high level programming  
131 language  $\mathcal{L}_{high}$  and a low level programming lan-  
132 guage  $\mathcal{L}_{low}$ , each is an infinite set of valid pro-  
133 gram strings. There exists a unary relation  $\rightarrow$  from  
134  $\mathcal{L}_{high}$  to  $\mathcal{L}_{low}$ . Given two monolingual corpora  
135  $L_{high} \subset \mathcal{L}_{high}$  and  $L_{low} \subset \mathcal{L}_{low}$ , the problem is to  
136 learn a translator  $F$  such that  $\forall x \in \mathcal{L}_{high}, (\exists u \in$   
137  $\mathcal{L}_{low}, x \rightarrow u) \rightarrow (x \rightarrow F(x))$ .* 138

139 The main challenge of this problem is that al-  
140 though the semantic alignment between corpora  
141  $L_{high}$  and  $L_{low}$  can be provided by oracle com-  
142 piler GCC, the semantic gap between them is  
143 huge, where many attributes of  $L_{high}$  cannot be  
144 directly expressed in  $L_{low}$ . For example, like for-  
145 loop and if-else semantics, the translation must  
146 learn **a posteriori** to generate jump instructions  
147 and corresponding labels to express the original  
148 control flows. According to Rice’s Theorem of  
149 computability theory (Rice, 1953), there is no set  
150 of rules that can accurately model the relation  $\rightarrow$ ,  
151 because it is undecidable whether two programs  
152 are semantically-equivalent. Instead, we will use  
153 behavioral-equivalent to approximate.

## 154 3 Methodology

155 In order to translate from a high-level code to  
156 low-level code well, where we choose C and  
157 x86 respectively for illustration, we face many  
158 challenges since the semantic gap is enormously  
159 large, comparing to translation between high level  
160 codes, like C-to-CUDA (Wen et al., 2022), Java-to-  
161 Python (Rozière et al., 2020), etc. Our approach  
162 for translating high-level C code to low-level x86  
163 assembly code focuses on generating semantic-  
164 equivalent code in best effort, where we focus on  
165 non-optimized generation, and the x86 code fol-  
166 lows the oracle GCC to learn the translation pro-  
167 cess.

168 This section gives a brief overview of what are  
169 the challenges in our scenario, and our practice to  
170 overcome these challenges.

### 171 3.1 Dataset Preprocessing

172 First we need to generate a C-x86 aligned bilingual  
173 corpora. We majorly choose AnghaBench (Da Silva  
174 et al., 2021), ExeBench (Armengol-Estapé et al.,  
175 2022) to obtain a large C corpora codebase. Then,  
176 we perform standard data preprocessing on the  
177 codebase: we filtered all C code with larger than  
178 2048 token size, with multiple function definitions,  
179 and with non-standard library dependencies, then

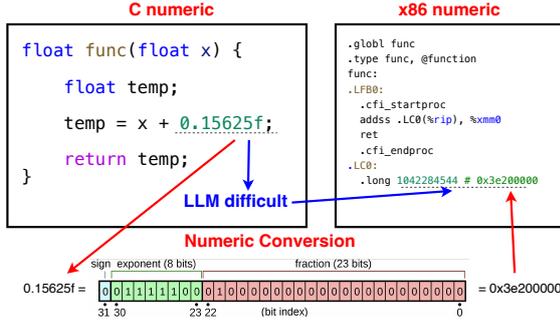


Figure 1: Numerical Conversion Feature Between C And x86

we use GCC-9.4.0 to compile each code with unified compiler options to obtain the corresponding x86 assembly corpora, which is naturally aligned with its C corpora.

After the initial preprocessing, we obtain a semantically aligned C-x86 bilingual corpora, which is already enough for the training process. However, after our first try on the training on this corpora, our model didn't learn well. After manually inspecting on the generation errors, we find the following challenges.

**Numerical Value Conversion.** A significant challenge in the translation between C and x86 assembly languages lies in the conversion of numerical values, which underscores the semantic differences between these languages. As depicted in Figure 1, In C, floating-point and double-precision values can be represented as literals, such as 1.0 or 3e-5. However, in assembly language, these numerical literals are not directly represented. Instead, they need to be converted to an internal representation following the IEEE-754 standard(iee, 1985) in most compiler implementations. This conversion process is rule-based and straightforward to implement. Yet, Large Language Models (LLMs) exhibit a notable weakness in this task, achieving a mere 3.8% accuracy on NumericBench, a large scale mathematical solving dataset derived from Math23K(Wang et al., 2017). This result underscores a critical limitation of LLMs in handling numerical computations.

To mitigate this limitation, we implement an effective data pre-processing method called symbolic interpretation, where we guide the LLM to generate symbolic expressions of the float/double values, which are subsequently processed by a rule-based interpreter. By delegating the actual numerical con-

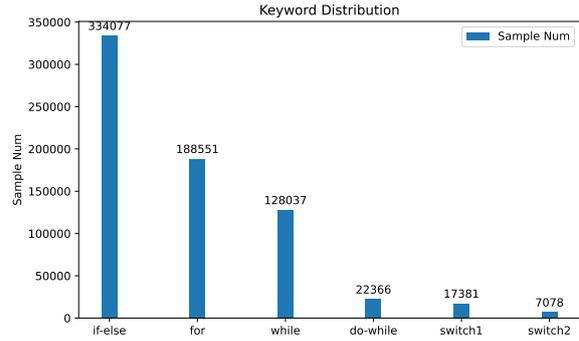


Figure 2: Long-tail Keyword Distribution of ExeBench

version to the interpreter, this method effectively circumvents the LLM's inherent weakness in numerical value conversions, thereby improving the overall accuracy of the translation process.

```

1 int x;
2 int main() {
3   switch(x) {
4     case 0: ...
5     case 1: ...
6     ...
7     case N: ...
8     default: ...
9   }
10  return 0;
11 }

```

Listing 1: C Switch1

```

1 int x;
2 int main() {
3   {
4     if(x == 0) ...
5     else if(x==1) ...
6     ...
7     else if(x==N) ...
8     else ...
9   }
10  return 0;
11 }

```

Listing 2: C Switch2

```

1 main:
2   movl x(%rip), %eax
3   cmpl $N, %eax
4   ja .Ldefault
5   leaq .LJT0(%rip), %rdx
6   movslq (%rdx,%rax,4), %rax
7   addq %rdx, %rax
8   notrack jmp %rax
9
10 .LJT0:
11   .long .L0-.LJT0
12   ...
13   .long .LN-.LJT0
14
15 .L0:
16   ...
17 .Ldefault:
18   ...

```

Listing 3: x86 Switch1

```

1 main:
2   movl x(%rip), %eax
3   cmpl $N, %eax
4   ja .Ldefault
5   cmpl $0, %eax
6   je .L0
7   cmpl $1, %eax
8   je .L1
9   ...
10  je .LN
11  jmp .Ldefault
12
13 .L0:
14   ...
15 .L1:
16   ...
17 .Ldefault:
18   ...

```

Listing 4: x86 Switch2

**Switch-case Statement Inconsistency.** Another kind of significant translation error lays on "switch-case" statement, where we observe that our baseline model generates inconsistently in two styles, where the original corpora messed them up. Listing 1 depicts the standard switch-case statement in C, and Listing 3 is its corresponding x86 assembly generated by GCC, where the cases are stored into a jump table, and using indirect jump instruction to control the jump target. However, switch-case statement can also be implemented by if-else logic, where Listing 2 depicts its semantic equivalent code in C, and Listing 4 is its x86 assem-

bly, where multiple comparison instructions and conditional jump instructions are used. By default, GCC generates the first type when cases are larger than threshold 4, and the second type otherwise, other compilers like Clang and MSVC also sharing this behavior with different thresholds. As depicted in Figure 2, we observe 7078 samples belong to the first and 17381 samples belong to the second in our initial training corpora, and their ratio on the whole corpora is also small, with 1.0% and 2.6% respectively. Comparing to other control keyword in C, which is clearly long-tailed.

To tackle the switch-case semantic inconsistency, we normalize the semantic of the switch-case statement to the if-else style in Listing 4, where we re-generate the x86 assembly from GCC compiler using compiler flag "-fno-jump-tables".

### 3.2 Dataset Augmentation

As already emphasized in the switch-case handling, the biased distribution of each semantic translation in the training corpora is a big challenge. Considering there are other long-tails besides switch-cases that also performs poorly, we need an automatic data augmentation method to improve the model’s accuracy on these long-tails. This is crucial and necessary because the LLM is only trained on limited corpora. If the input is few or even none in the corpora, it will translate poorly without any surprise.

Inspired by (Madaan et al., 2023), we construct an automatic refinement data augmentation framework as depicted in Figure 3, where the model is first trained on corpora from the previous method, and evaluated through multiple metrics, where we collect on the low-metric samples where we assume the model under-fits to learn them. Then we synthesize more samples from the incorrect samples to improve the distribution. we choose to use mistral-7B(Jiang et al., 2023) as the synthesizing LLM in our implementation, where we instruct the LLM to analyze, categorize, and generate ten times more similar samples.

With more long-tail samples been synthesized, we re-sample the corpora by adding synthesized samples to it, creating a re-sampled corpora that better represents the long-tail problems. Finally, we re-train the model on this re-sampled dataset. The whole above process can be iteratively executed, where more under-fitting long-tails can be discovered, re-sampled, and improved.

This refinement framework allows the model to

better learn how to handle these long-tailed samples, leading to improved accuracy in the generated low-level code. We provide examples illustrating its validity in the case studies.

### 3.3 Fine-Tuning

Machine translation has evolved significantly with the advent of neural machine translation (NMT), where models are trained on large corpora of text to learn the nuances of language translation. The general principle of machine translation, as pioneered by (Rozière et al., 2020), involves two key stages: pretraining and fine-tuning. Initially, models are pretrained on monolingual corpora to learn language features. Subsequently, they are fine-tuned on paired corpora to guide the translation between two languages.

We employ Low Rank Adaptation(Hu et al., 2021), one of the most popular Parameter-Efficient Fine-Tuning methods, to adapt LLMs to our translation task. LoRA modifies a small subset of the model’s weights by decomposing the weight changes into two smaller matrices, which are then fine-tuned. This approach allows us to bypass the initial pre-training phase typical in machine translation, as LLMs are already pretrained on extensive monolingual corpora. We use **codellama-13b**(Rozière et al., 2022) as our foundation model.

Similar to the construction of the training corpora, we construct the evaluation corpora solely on C, where we choose from the IO evaluation part of ExeBench(Armengol-Estapé et al., 2022) and Math23K(Wang et al., 2017), to evaluate the model’s translation accuracy, where the former represents general purpose code and the latter represents numerical computations. More detailed corpora components can be found in the following Evaluation Section.

## 4 Evaluation

To evaluate our proposed code translation methods, we perform a series of experiments on function-level C programs. We use the directly finetuned codellama-13b model as the baseline.

First, we perform end-to-end translation on a large evaluation dataset depicted in Table 1, which consists of 57,552 C functions, where we compare with the baseline model, the numerical conversion augmented model, the switch normalization augmented model, both applied model, and GPT-4-turbo. Then, as an ablation study, we compare

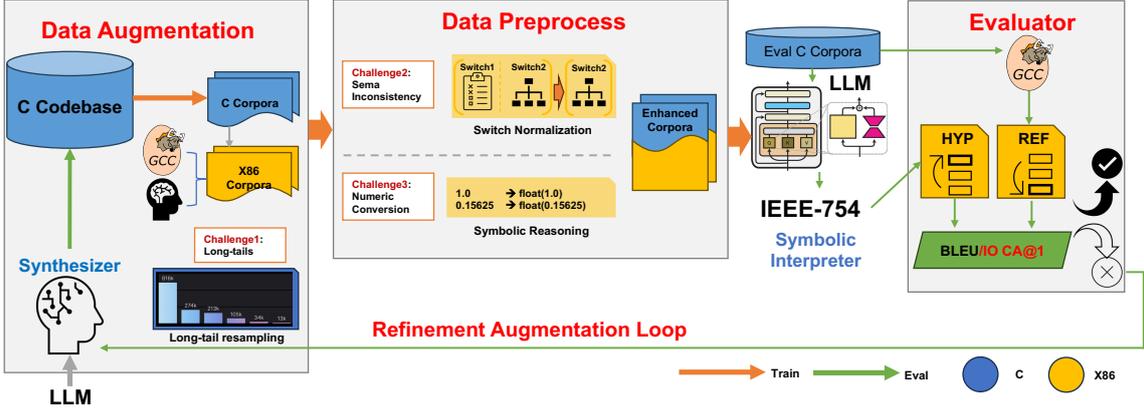


Figure 3: Data Augmentation Framework Overview

Datasets	Size	Tok (C)	Tok (x86)
Train	679665	107	391
Train-Num	40000	168	594
Eval	57552	110	
ExeBench	35704	108	
Numeric	21104	111	
Switch	744	237	

Table 1: Dataset Details

models accuracy within the numerical-specified subset, and the switch-specified subset.

#### 4.1 Dataset

Table 1 shows the details of the dataset we used in training and evaluation. We first finetune the codellama-13b foundation model to perform C-to-x86 code translation task, where we use dataset derived from ExeBench(Armengol-Estapé et al., 2022) and AnghaBench(Da Silva et al., 2021), two large scale dataset of compilable C functions, we first apply data cleaning, where we filtered oversized functions(we limit the size to 2048 tokens in our settings), and other features we are not going to cover like inline assembly. Finally we get a 680K size training dataset for baseline training. In the numerical value conversion preprocessing part, we establish a 40k numerical adjusted corpora to finetune the model. For evaluation part, we construct a 57K size dataset with I/O behavioral checks. As for the numerical conversion and switch-case generation challenges we found in our methods, we also categorize specified subsets for the evaluation, where a 21K numeric-specific subset and a 744 switch-specific subset are evaluated individually.

#### 4.2 Setup and Metrics

We set up the experiment on a Ubuntu 22.04 server with Intel Xeon Platinum 8358 CPU and 4 x A800 80GB GPUs. We begin with the codellama-13b-instruct checkpoint from huggingface hub as our foundation model. We then directly apply LoRA finetuning with the 680K training corpora to learn the C-to-x86 translation task, which we considered as the **Baseline** model. Later we apply the two data pre-processing methods, switch-case normalization or/and numerical value conversion, to adjust the training corpora, and re-train on the foundation model to get the **Switch** enhanced model, **Numeric** enhanced model and **ALL** enhanced model. We also use **GPT-4-Turbo**, the most advanced LLM, as the second baseline to compare with.

During the training process, we use **lora\_r** = 128, **lora\_alpha**=32, **lora\_dropout**=0.05 in the LoRA modules, where we attach all **Q, K, V, O** in the model for training. We use the sum of token-level cross-entropy loss with teacher-forcing as the loss function, which is on par with (Rozière et al., 2020). We use AdamW(Kingma and Ba, 2014) as the optimizer and apply a cosine learning rate that top at 1e-4 in training. The training process is performed fully in float16 precision, where we train the model for 1 epoch in 70 hours using 4xA800 80GB GPUs.

We evaluate the above models on the 57,552 functions evaluation dataset. We also construct the 21,104 size numeric-specified and the 744 size switch-specified subsets from the full dataset. Then we perform end-to-end evaluation on these datasets, which also serves as an ablation study. We examine each generated function in x86 assembly by linking it with the driver code that called the function to

396 obtain an executable, then performing Input/Out- 447  
397 put(IO) correctness checks. We use greedy genera- 448  
398 tion in the generation process, so the IO accuracy 449  
399 can also be viewed as CA@1 in other machine 450  
400 translation tasks. 451

### 401 4.3 End-to-End Evaluation 452

402 **Figure 4** summarizes the empirical end-to-end re- 453  
403 sults ablating different methods and comparing 454  
404 with GPT-4-Turbo. The baseline model, shows 455  
405 a fair overall result, which can reach 60% IO Ac- 456  
406 curacy. More detailed breakdowns of its wrong 457  
407 translations show it majorly falls into the following 458  
408 types: 459

409 **Generating wrong numerical values.** We cap- 460  
410 ture all the functions within the evaluation dataset, 461  
411 where there exists numerical value initialization, 462  
412 and categorize them into a numerical dataset, Nu- 463  
413 mericBench for breakdown. We find out that the 464  
414 baseline model can only generate 3.8% of Numeric- 465  
415 Bench correctly, and most of these happen-to-be- 466  
416 correct values are values with high frequency in 467  
417 the dataset, like 1.0 and 0.0. This breakdown in- 468  
418 deed reveals a crucial drawback of the LLM-based 469  
419 machine translation method. We then apply the 470  
420 symbolic interpretation method on the dataset pre- 471  
421 processing stage, which significantly improved the 472  
422 generation accuracy, rising from 3.8% to over 90%. 473

423 **Generating wrong labels and jump tables.** We 474  
424 evaluate the evaluation dataset and collect those 475  
425 with incorrect execution behaviors, where we find 476  
426 many in switch-case generations. After analyzing 477  
427 the generated assembly, we find out their transla- 478  
428 tion is very likely in an underfitting manner. We 479  
429 also find out the training dataset is inconsistent with 480  
430 the semantic of switch-case code generation, when 481  
431 cases numbers are above the threshold, they use 482  
432 indirect jump on the jump table in the generated 483  
433 assembly, while kept the if-else style in the others. 484  
434 This inconsistent behaviour is by default open for 485  
435 our oracle compiler gcc even in O0 optimization 486  
436 level, where dataset makers can hardly notice. 487

437 We further perform categorization of control- 488  
438 flow statements on the training dataset, which is 489  
439 clearly summarized in **Figure 2**, where the two 490  
440 types of switch-case generation are both rare in 491  
441 corpora, counting for 2.6% and 1.0% respectively. 492  
442 This categorization result depicts a long-tail dis- 493  
443 tribution in the training dataset, where the model 494  
444 under-fits the switch-case statement generation, 495  
445 and the inconsistency on switch-case statement gen- 496  
446 erations may further confuse the model.

To tackle this problem, we perform switch-  
case normalization, where we enable the GCC op-  
tion "-fno-jump-tables" to unify the generation be-  
haviours on switch-case, and re-train the model. As  
illustrated in **Figure 4**, the normalization of switch-  
case semantic improves the switch-case translation  
accuracy from 50.86% to 66.57%, which shows the  
effectiveness of the augmentation method.

**Other types of wrong generations**, which in-  
clude wrong generation of very long function log-  
ics, wrong generation of stack operation, wrong  
C-struct offset calculation, and wrong generation  
on rare samples, like AVX intrinsics, etc.

In the end-to-end evaluation, we tackle the first  
two kinds of errors. By augmenting with both nu-  
merical conversion and switch-case normalization,  
we successfully improves the overall I/O Accuracy  
to 91.88%, which improves drastically from the  
baseline model. To compare with, GPT-4-Turbo  
can only achieve 40.85% I/O Accuracy even with  
careful promptings applied.

## 468 5 Case Study 469

We conduct case studies to demonstrate how to  
overcome the challenges using data augmentation  
methods to learn C-to-x86 translation.

The first case study demonstrated a function  
that need float/double numerical value conversion.  
In x86 language, float/double immediate numbers  
can not be encoded in instructions directly, and  
modern compilers like GCC save them in binary  
format following the IEEE-754 standard. So as  
long as the program exists numerical initialization,  
there are numerical conversions during the trans-  
lation process, where LLMs perform poorly. As  
depicted in **Figure 5**, direct value conversion us-  
ing implicit IEEE-754 rule makes LLMs hard to  
predict, where the baseline models are very likely  
to generate wrong numbers. By delegating the  
numerical conversions from LLMs to rule-based  
interpreters, where we augment the model to gener-  
ate symbolic expressions instead of direct guessing,  
LLMs delegate the numerical conversion to rule-  
based interpreters, which can handle their conver-  
sions well, so that the numerical handling drawback  
of LLMs is efficiently mitigated.

The second case study depicted in **Figure 6**  
shows the challenge of switch-case generation,  
where the jump-table style generation are hard to  
learn for LLMs. The baseline model fails in the  
generation of jump table items, causing repeated

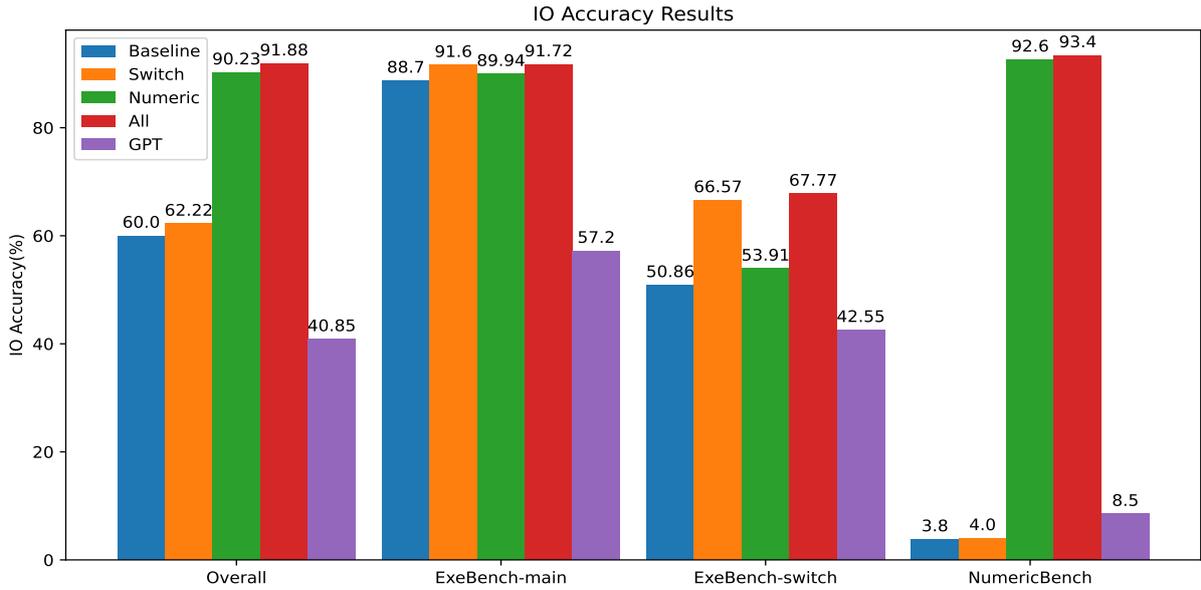


Figure 4: IO Accuracy Results

<pre> 1 float func() 2 { 3   float costA = 6.0; 4   float costB = 0.125; 5   float cash = 50.0; 6   float numA = 4.0; 7   float numB; 8   float temp; 9   temp = costA * numA; 10  temp = cash - temp; 11  numB = temp / costB; 12  return numB; 13 } </pre>	<pre> func: ... movss .LC0(%rip), %xmm0 movss %xmm0, -20(%rbp) ... .LC0: .long 1086324736 ; 6.0 .LC1: .long 1040187392 ; 0.125 .LC2: .long 1112014848 ; 50.0 .LC3: .long 1082130432 ; 4.0 </pre>	<pre> 1 2 3 4 5 6 7 8 9 10 11 12 13 </pre>
--	--	--

Figure 5: Case Study 1: Numerical Conversion

<pre> 1 int color_char_to_attr(char c) 2 { 3   switch (c) 4   { 5     ... 6     case 'R': 7       return (4); 8     case 'G': 9       return (5); 10    case 'B': 11      return (6); 12    ... 13  } 14  return (-1); 15 } </pre>	<pre> color_char_to_attr: ... subl \$66, %eax ; 'B' cmpl \$16, %eax ; 'R' - 'B' ja .L2 leaq .L4(%rip), %rdx movslq (%rdx,%rdi,4), %rax addq %rdx, %rax notrack jmp *%rax .section .rodata .L4: .long .L2-.L4 .long .L2-.L4 .long .L2-.L4 ; ... repeated pattern </pre>	<pre> 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 </pre>
--	--	--

Figure 6: Case Study 2: Switch Generation

patterns until the maximum generation length. By leveraging the if-else style data augmentation, the model has learned to treat switch-case statements as if-else style, where if-else corpus are on the head of keyword distribution with hundreds of thousand samples comparing to the rare long-tails, the deficient learning of switch-case generation is also mitigated.

The last case study shows how our refinement

framework improving the long-tails performance. As depicted in Figure 7, AVX instructions are the SIMD extension in x86 assembly language, and is encapsulated as AVX intrinsics to be used in C language.

Recalling Figure 3, we introduce the refinement framework to augment the incorrect generations, which is inspired by (Madaan et al., 2023). Initially, there are no AVX-related samples in the training corpora at all, where the model without any surprise translate incorrectly without apriori. Then the incorrect AVX sample is captured by the evaluator together with other incorrect samples. we then use LLM to analyze the C code, and synthesize more based on several rules as prompts to generate more C samples closely related to the incorrect cases. We use mistral-7B(Jiang et al., 2023) as the synthesizer LLM in our implementation. Finally, the synthesized augmented C corpora of incorrect samples is added back to the training dataset, where retraining/finetuning can be performed depending on the need.

Back to the case itself, a 10x synthesizing is sufficient enough to learn a new feature with simple semantic pattern, like the `_mm256_add_ps` intrinsic in the case, which simply generates a `vaddps` instruction. Such learning ability of aligning C and x86 semantics is very impressive, which shows the few-shot learning potential in the language translation task. Although more complex patterns need more cases to learn well, luckily, the refinement framework can be executed iteratively, which can

<pre> 1 void foo(float *x, *y, *o) { 2   xx = _mm256_load_ps(x); 3   yy = _mm256_load_ps(y); 4   zz = _mm256_add_ps(xx, yy); 5   _mm256_store_ps(o, zz); 6 } </pre>	<pre> foo:   vmovaps (%rdi), %ymm0   vmovaps (%rsi), %ymm1   vaddps %ymm1, %ymm0   vmovaps %ymm0, (%rdx)   ret </pre>	<pre> 1 2 3 4 5 6 </pre>
---	---	--------------------------

Figure 7: Case Study 3: AVX Intrinsics Learning

resample the corpora based on the generation accuracy, so that more complex cases can get more samples to be learned.

## 6 Related Work

**Code Translation** aims to translate a piece of code (usually a function or method) into another programming language. Early studies like (Nguyen et al., 2015) uses traditional statistical machine translation method. Neural-based method like (Chen et al., 2018) starts to be dominant, and capture the tree structure of programming languages. The emergence of pre-trained language models of code, such as CodeBERT(Feng et al., 2020) and CodeT5(Wang et al., 2021), has further improved the state of code translation. Large Language Models(LLMs)(OpenAI et al., 2023; Rozière et al., 2022) have continued this trend, showing promise in code translation task. However, the above approaches usually require fine-tuning on parallel corpora, which is often scarce.

**Data augmentation** techniques have been extensively used and found effective in machine translation tasks, which served as a solution to the scarcity of parallel corpora. Transcoder(Rozière et al., 2020) first propose back translation approach to learn unsupervised code translation, where the back-translation process also generates an automatic parallel corpora augmentation method. Transcoder-ST(Roziere et al., 2021), CodeXGlue(Lu et al., 2021), BabelTower(Wen et al., 2022) and CMTrans(Xie et al., 2023) also follow this approach, to obtain parallel corpora during the learning process. Besides direct generation, (Szafraniec et al., 2023) explores an IR-in-the-middle approach, while (Tang et al., 2023; Ahmad et al., 2023) both introduce an intermediate code summary stage, to improve the code translation accuracy.

To construct a balanced corpora in limited size in monolingual language is also challenging, it is naturally in a long-tailed distribution for different aspects of code semantics. where neural models tend to perform low accuracy on the tails due to

lack of samples. (Zhout et al., 2023) reveals that LLMs can perform between 30% to 254% worse in long-tailed cases, where the model under-fits them. Inspired by the survey of long-tailed learning(Zhang et al., 2023a), we establish a refinement augmentation method, where long-tailed C samples are recognized in the evaluation process via metrics, then analyzed, synthesized by another powerful LLM, compiled by GCC to obtain parallel samples, finally augmented the corpora with more long-tailed knowledge.

**Cross Level Code Translation.** On high-level code to low-level code translation researches, (Armengol-Estapé and O’Boyle, 2021) first gives a try of using neural machine translation on this scenario. (Guo and Moses, 2022) further studies on C-to-LLVM IR translation. However, they only perform limited investigations on the methods, and their results are still on the preliminary stage. There are more related works on the reverse process, to recover high-level code from low-level code(Fu et al., 2019; Cao et al., 2022; Armengol-Estapé et al., 2023). Unlike the difficulty on semantic mapping to low level code in our challenges, their challenges mainly are on optimization recovery and type inference, while the semantic recovery is relatively simpler.

## 7 Conclusion

Machine translation from high-level language to low-level machine instructions is difficult. Even using advanced LLMs can not reach high accuracy. By implementing symbolic interpretation and switch-case normalization, two novel data preprocessing methods, we overcome numerical value conversion and switch-case semantic inconsistency, two significant challenges in C-to-x86 language translation.

To improve the accuracy on long-tailed samples where the model under-fits to learn, we propose an automatic refinement augmentation framework to obtain improved accuracy on the long-tails by using synthesizing method on incorrect cases.

Finally we achieve state-of-the-art IO accuracy, over 91%, when translating C-to-x86 on a large-scale evaluation dataset. Comparing to LLM-only method(GPT-4-Turbo, 40.85%), and finetuning-only baseline method(59.87%), the methods show great efficiency.

629  
630  
631  
632  
633  
634  
635  
636  
637  
638  
639  
640  
641  
642  
643  
644  
645  
646  
647  
648  
649  
650  
651  
652  
653  
654  
655  
656  
657  
658  
659  
660  
661  
662  
663  
664  
665  
666  
667  
668  
669  
670  
671  
672  
673  
674  
675  
676  
677  
678  
679  
680  
681  
682  
683  
684

## Acknowledgements

## References

1985. [Ieee standard for binary floating-point arithmetic. ANSI/IEEE Std 754-1985](#), pages 1–20.

Wasi Uddin Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2023. [Summarize and generate to back-translate: Unsupervised translation of programming languages](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, pages 1520–1534. Association for Computational Linguistics.

Jordi Armengol-Estapé and Michael FP O’Boyle. 2021. Learning c to x86 translation: An experiment in neural compilation. *arXiv preprint arXiv:2108.07639*.

Jordi Armengol-Estapé, Jackson Woodruff, Alexander Brauckmann, José Wesley de Souza Magalhães, and Michael FP O’Boyle. 2022. Exebench: an ml-scale dataset of executable c functions. In *Proceedings of the 6th ACM SIGPLAN International Symposium on Machine Programming*, pages 50–59.

Jordi Armengol-Estapé, Jackson Woodruff, Chris Cummins, and Michael FP O’Boyle. 2023. Slade: A portable small language model decompiler for optimized assembler. *arXiv preprint arXiv:2305.12520*.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.

Ying Cao, Ruigang Liang, Kai Chen, and Peiwei Hu. 2022. Boosting neural networks to decompile optimized binaries. In *Proceedings of the 38th Annual Computer Security Applications Conference*, pages 508–518.

Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie

Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. 2021. [Evaluating large language models trained on code](#). 685  
686  
687  
688  
689  
690  
691  
692

Xinyun Chen, Chang Liu, and Dawn Song. 2018. Tree-to-tree neural networks for program translation. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, page 2552–2562, Red Hook, NY, USA. Curran Associates Inc. 693  
694  
695  
696  
697  
698

Anderson Faustino Da Silva, Bruno Conde Kind, José Wesley de Souza Magalhães, Jerônimo Nunes Rocha, Breno Campos Ferreira Guimaraes, and Fernando Magno Quinão Pereira. 2021. Anghabench: A suite with one million compilable c benchmarks for code-size reduction. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)*, pages 378–390. IEEE. 699  
700  
701  
702  
703  
704  
705  
706

Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. [Codebert: A pre-trained model for programming and natural languages](#). 707  
708  
709  
710  
711

Cheng Fu, Huili Chen, Haolan Liu, Xinyun Chen, Yuan-dong Tian, Farinaz Koushanfar, and Jishen Zhao. 2019. Coda: An end-to-end neural program decompiler. *Advances in Neural Information Processing Systems*, 32. 712  
713  
714  
715  
716

Daya Guo, Canwen Xu, Nan Duan, Jian Yin, and Julian Mcauley. 2023. [LongCoder: A long-range pre-trained language model for code completion](#). In *Proceedings of the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine Learning Research*, pages 12098–12107. PMLR. 717  
718  
719  
720  
721  
722  
723

Zifan Carl Guo and William S. Moses. 2022. [Enabling transformers to understand low-level programs](#). In *2022 IEEE High Performance Extreme Computing Conference (HPEC)*, pages 1–9. 724  
725  
726  
727

Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*. 728  
729  
730  
731  
732

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Léo Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. [Mistral 7b](#). 733  
734  
735  
736  
737  
738  
739

Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*. 740  
741  
742

743	Haochen Li, Chunyan Miao, Cyril Leung, Yanxian	Han, Jeff Harris, Yuchen He, Mike Heaton, Jo-	802
744	Huang, Yuan Huang, Hongyu Zhang, and Yanlin	hannes Heidecke, Chris Hesse, Alan Hickey, Wade	803
745	Wang. 2022a. <a href="#">Exploring representation-level aug-</a>	Hickey, Peter Hoeschele, Brandon Houghton, Kenny	804
746	<a href="#">mentation for code search</a> . In <i>Proceedings of the</i>	Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu	805
747	<i>2022 Conference on Empirical Methods in Natu-</i>	Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger	806
748	<i>ral Language Processing</i> , pages 4924–4936, Abu	Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie	807
749	Dhabi, United Arab Emirates. Association for Com-	Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser,	808
750	putational Linguistics.	Ali Kamali, Ingmar Kanitscheider, Nitish Shirish	809
751	Xiaonan Li, Daya Guo, Yeyun Gong, Yun Lin, Ye-	Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook	810
752	long Shen, Xipeng Qiu, Daxin Jiang, Weizhu Chen,	Kim, Christina Kim, Yongjik Kim, Hendrik Kirchner,	811
753	and Nan Duan. 2022b. <a href="#">Soft-labeled contrastive pre-</a>	Jamie Kiros, Matt Knight, Daniel Kokotajlo,	812
754	<a href="#">training for function-level code representation</a> . In	Łukasz Kondraciuk, Andrew Kondrich, Aris Kon-	813
755	<i>Findings of the Association for Computational Lin-</i>	stantinidis, Kyle Kopic, Gretchen Krueger, Vishal	814
756	<i>guistics: EMNLP 2022</i> , pages 118–129, Abu Dhabi,	Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan	815
757	United Arab Emirates. Association for Computa-	Leike, Jade Leung, Daniel Levy, Chak Ming Li,	816
758	tional Linguistics.	Rachel Lim, Molly Lin, Stephanie Lin, Mateusz	817
759	Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey	Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue,	818
760	Svyatkovskiy, Ambrosio Blanco, Colin Clement,	Anna Makanju, Kim Malfacini, Sam Manning, Todor	819
761	Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021.	Markov, Yaniv Markovski, Bianca Martin, Katie	820
762	Codexglue: A machine learning benchmark dataset	Mayer, Andrew Mayne, Bob McGrew, Scott Mayer	821
763	for code understanding and generation. <i>arXiv</i>	McKinney, Christine McLeavey, Paul McMillan,	822
764	<i>preprint arXiv:2102.04664</i> .	Jake McNeil, David Medina, Aalok Mehta, Jacob	823
765	Aman Madaan, Niket Tandon, Prakhar Gupta, Skyler	Menick, Luke Metz, Andrey Mishchenko, Pamela	824
766	Hallinan, Luyu Gao, Sarah Wiegrefe, Uri Alon,	Mishkin, Vinnie Monaco, Evan Morikawa, Daniel	825
767	Nouha Dziri, Shrimai Prabhunoye, Yiming Yang,	Mossing, Tong Mu, Mira Murati, Oleg Murk, David	826
768	et al. 2023. Self-refine: Iterative refinement with	Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak,	827
769	self-feedback. <i>arXiv preprint arXiv:2303.17651</i> .	Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh,	828
770	Anh Tuan Nguyen, Tung Thanh Nguyen, and Tien N.	Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex	829
771	Nguyen. 2015. <a href="#">Divide-and-conquer approach for</a>	Paino, Joe Palermo, Ashley Pantuliano, Giambat-	830
772	<a href="#">multi-phase statistical migration for source code (t)</a> .	tista Parascandolo, Joel Parish, Emy Parparita, Alex	831
773	In <i>2015 30th IEEE/ACM International Conference</i>	Passos, Mikhail Pavlov, Andrew Peng, Adam Perel-	832
774	<i>on Automated Software Engineering (ASE)</i> , pages	man, Filipe de Avila Belbute Peres, Michael Petrov,	833
775	585–596.	Henrique Ponde de Oliveira Pinto, Michael, Poko-	834
776	OpenAI, :, Josh Achiam, Steven Adler, Sandhini Agar-	rny, Michelle Pokrass, Vitchyr Pong, Tolly Pow-	835
777	wal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Ale-	ell, Alethea Power, Boris Power, Elizabeth Proehl,	836
778	man, Diogo Almeida, Janko Alvenschmidt, Sam Alt-	Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh,	837
779	man, Shyamal Anadkat, Red Avila, Igor Babuschkin,	Cameron Raymond, Francis Real, Kendra Rimbach,	838
780	Suchir Balaji, Valerie Balcom, Paul Baltescu, Haim-	Carl Ross, Bob Rotsted, Henri Roussez, Nick Ry-	839
781	ing Bao, Mo Bavarian, Jeff Belgum, Irwan Bello,	der, Mario Saltarelli, Ted Sanders, Shibani Santurkar,	840
782	Jake Berdine, Gabriel Bernadett-Shapiro, Christo-	Girish Sastry, Heather Schmidt, David Schnurr, John	841
783	pher Berner, Lenny Bogdonoff, Oleg Boiko, Made-	Schulman, Daniel Selsam, Kyla Sheppard, Toki	842
784	laine Boyd, Anna-Luisa Brakman, Greg Brockman,	Sherbakov, Jessica Shieh, Sarah Shoker, Pranav	843
785	Tim Brooks, Miles Brundage, Kevin Button, Trevor	Shyam, Szymon Sidor, Eric Sigler, Maddie Simens,	844
786	Cai, Rosie Campbell, Andrew Cann, Brittany Carey,	Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin	845
787	Chelsea Carlson, Rory Carmichael, Brooke Chan,	Sokolowsky, Yang Song, Natalie Staudacher, Fe-	846
788	Che Chang, Fotis Chantzis, Derek Chen, Sully Chen,	lipe Petroski Such, Natalie Summers, Ilya Sutskever,	847
789	Ruby Chen, Jason Chen, Mark Chen, Ben Chess,	Jie Tang, Nikolas Tezak, Madeleine Thompson, Phil	848
790	Chester Cho, Casey Chu, Hyung Won Chung, Dave	Tillet, Amin Tootoonchian, Elizabeth Tseng, Pre-	849
791	Cummings, Jeremiah Carrier, Yunxing Dai, Cory	ston Tuggle, Nick Turley, Jerry Tworek, Juan Fe-	850
792	Decareaux, Thomas Degry, Noah Deutsch, Damien	lipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya,	851
793	Deville, Arka Dhar, David Dohan, Steve Dowl-	Chelsea Voss, Carroll Wainwright, Justin Jay Wang,	852
794	ing, Sheila Dunning, Adrien Ecoffet, Atty Eleti,	Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei,	853
795	Tyna Eloundou, David Farhi, Liam Fedus, Niko	CJ Weinmann, Akila Welihinda, Peter Welinder, Ji-	854
796	Felix, Simón Posada Fishman, Juston Forte, Is-	ayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner,	855
797	abella Fulford, Leo Gao, Elie Georges, Christian	Clemens Winter, Samuel Wolrich, Hannah Wong,	856
798	Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh,	Lauren Workman, Sherwin Wu, Jeff Wu, Michael	857
799	Rapha Gontijo-Lopes, Jonathan Gordon, Morgan	Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qim-	858
800	Grafstein, Scott Gray, Ryan Greene, Joshua Gross,	ing Yuan, Wojciech Zaremba, Rowan Zellers, Chong	859
801	Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse	Zhang, Marvin Zhang, Shengjia Zhao, Tianhao	860
		Zheng, Juntang Zhuang, William Zhuk, and Barret	861
		Zoph. 2023. <a href="#">Gpt-4 technical report</a> .	862
		Yada Pruksachatkun, Jason Phang, Haokun Liu,	863
		Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang,	864



972 similar problems that LLMs need to adjust to. We  
973 also consider this as future work.