

Variable Name Recovery in Decompiled Binary Code using Constrained Masked Language Modeling

Pratyay Banerjee
Kuntal Kumar Pal
pbanerj6@asu.edu
kkpal@asu.edu
Arizona State University
Tempe, Arizona, USA

Fish Wang
Chitta Baral
fishw@asu.edu
chitta@asu.edu
Arizona State University
Tempe, Arizona, USA

ABSTRACT

Decompilation is the procedure of transforming binary programs into a high-level representation, such as source code, for human analysts to examine. While modern decompilers can reconstruct and recover much information that is discarded during compilation, inferring variable names is still extremely difficult. Inspired by recent advances in natural language processing, we propose a novel solution to infer variable names in decompiled code based on Masked Language Modeling, Byte-Pair Encoding, and neural architectures such as Transformers and BERT. Our solution takes *raw* decompiler output, the less semantically meaningful code, as input, and enriches it using our proposed *finetuning* technique, Constrained Masked Language Modeling. Using Constrained Masked Language Modeling introduces the challenge of predicting the number of masked tokens for the original variable name. We address this *count of token prediction* challenge with our post-processing algorithm. Compared to the state-of-the-art approaches, our trained VarBERT model is simpler and of much better performance. We evaluated our model on an existing large-scale data set with 164,632 binaries and showed that it can predict variable names identical to the ones present in the original source code up to 84.15% of the time.

CCS CONCEPTS

• **Security and privacy** → **Software reverse engineering**; • **Computing methodologies** → **Natural language processing**; **Machine learning**; **Model development and analysis**.

KEYWORDS

reverse engineering, decompilation, variable name recovery, masked language modeling, bert

ACM Reference Format:

Pratyay Banerjee, Kuntal Kumar Pal, Fish Wang, and Chitta Baral. 2020. Variable Name Recovery in Decompiled Binary Code using Constrained Masked Language Modeling. In *Proceedings of KDD '20*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/1122445.1122456>

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '20, Aug 22–27, 2020, San Diego, CA

© 2018 Association for Computing Machinery.
ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00
<https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

Despite the recent progress in the software open source movement, source code is oftentimes unavailable to security researchers, analysts, and even software vendors. This makes *software reverse engineering*, a family of techniques that aim at understanding the behavior of software without accessing its original source code, irreplaceable [12]. Reverse engineering techniques are essential to accomplishing a set of security-critical tasks, including malware analysis, vulnerability discovery, software defect diagnosis, and software repairing [2].

With the rapid development of modern decompilers (e.g., Hex-Rays Decompiler [11, 21] and Ghidra [15]), decompilation, or reverse compilation, is gradually becoming an essential technique for software reverse engineering tasks. Built atop advanced program analysis techniques and sophisticated heuristics, these decompilers can identify much information that is lost during compilation and binary stripping, such as function boundaries, function prototypes, variable locations, variable types, etc. Names of variables in a program usually embed much semantic information that is crucial for understanding the behaviors and logic of the program. Nevertheless, the state-of-the-art decompilers are all limited to recovering structural information and cannot recover meaningful variable names. Therefore, the usual first step taken by a decompiler user is painfully interpreting the decompiled code and renaming the unnamed variables one by one. This is, unfortunately, an insurmountable obstacle for reverse engineering stripped binary programs.

While researchers have taken steps in predicting variable names in high-level programming languages [1, 3, 4, 26, 41, 45], it is worth noting that inferring variable names in decompiled binary code poses a unique set of challenges. High-level programming languages like Java, Python, and JavaScript are syntactically rich: Variable types are preserved in these languages, while they are usually eliminated in binaries. In fact, type inferencing on binary programs still remains an open problem [8]. Moreover, modern compilers that generate machine code produce all sorts of irreversible changes [52], which diminish the power of variable name prediction techniques that are based on local data dependencies. DEBIN [19] and DIRE [30] are the pioneers in predicting variable names in stripped binary programs. DEBIN proposes a statistical model that yields poor performance and generalizability. DIRE produces much better results by employing an LSTM[23] and Gated Graph Neural Network-based model [31] that is trained on a huge

corpus of source code with sophisticated structural features extracted from Abstract Syntax Trees. Their approach is complex and error-prone.

Recently, there has been a substantial improvement in the field of natural language processing with the advent of techniques such as Transformers [46], Masked Language Modeling [14], and BERT [14]. With the help of these techniques, we propose a novel solution to predicting variable names in stripped binary programs using only the *raw* decompiled code as input. We propose *Constrained Masked Language Modeling*, a modification of the Masked Language Modeling task to train a BERT model to recover variable names. We learn a vocabulary of code-tokens using Byte-Pair Encoding and learn rich contextual representation using Masked Language Modeling.

Using Masked Language Modeling to predict variable names introduces a new challenge of predicting the *Count of Variable Name Tokens*, i.e., determining how many mask tokens are required to predict a variable name. We address this challenge by using a heuristics-based algorithm.

Both DIRE and our aforementioned solutions are computationally expensive. After careful analysis of the nature of the task, we argue that predicting variable names in the decompiled code does not require such a deep and computationally intensive neural architecture. To support our argument, we train a smaller BERT model and achieve a similar performance; Additionally, it is capable of taking longer input sequences of code.

We train and evaluate our models using the DIRE data set, a collection of 164,632 x86-64 binary programs generated from C projects crawled from GitHub [30]. Our evaluation results show that our models can predict variable names that are identical to the ones used in original source code up to 84.15% of the time, which surpasses our baselines, DEBIN and DIRE.

Our contributions are summarized below:

- We adapt Masked Language Modeling for the task of variable name recovery in decompiled code and propose Constrained Masked Language Modeling.
- We address the challenge of *Count of Variable Name Tokens* using a heuristic-based algorithm and evaluate its performance.
- We train two BERT models, VarBERT-base and VarBERT-small, that only use *raw* decompiled code as input. Models and code will be open-sourced.
- In the task of variable name recovery, we surpass DIRE, the state-of-the-art solution, on a large data set of 164,632 binary programs. We improve by **12.36%** in overall accuracy and **49%** in generalizability, that is, our model can successfully predict variables names that are not present in the training set.

In the spirit of open science, We will open source our source code and trained BERT models upon the publication of this paper.

2 BACKGROUND AND RELATED WORK

Our solution leverages techniques in machine learning and natural language processing to refine the results of binary code decompilation. Before diving into the details of our solution, in this section,

we will first present the necessary background as well as some work that is closely related to our proposed solution.

2.1 Binary Code Decompilation

Binary code decompilation is the process of converting compiled binary code to a higher-level representation, oftentimes C source code or C-like pseudocode. Compilation and binary stripping are lossy procedures where much semantic information is discarded since it is useless to CPUs. Such information, such as control flow structures, function names, function prototypes, variable types, and variable names, is very important to software reverse engineering tasks. While modern decompilers have made significant progress in recovering and inferring various types of lost information, recovering variable types in decompiled code remains challenging.

2.2 NLP and Program Analysis

Structured programs and natural languages have similar statistical properties [2, 13, 22]. This has led to the application of statistical models in program analysis and software engineering. Some of the notable applications are code completion [7, 35, 36], code synthesis [5, 18, 29, 32, 38, 54], obfuscation [33, 37], code fixing [16, 17], bug detection [40, 47], information extraction [10, 43, 51], syntax error detection [9] and correction [6], summarization [24, 25] and clone detection [49].

2.3 Statistically Predicting Variable Names

The probabilistic graphical model, Conditional Random Fields (CRF) has been useful in the prediction of syntactic names of identifiers and their semantic type information of JavaScript programs. Their system JSNICE [41] were able to correctly predict 63% of the names and 81% of the type information. Though our goal of prediction of variable names is same our approach and program domain differ considerably from theirs.

Another system AUTONYM [45] surpassed JSNICE in the prediction of variable names for Javascript codes using statistical machine translation. Jaffe et al. [26] also used statistical machine translation to generate meaningful names to the variables for the decompiled source codes written in C. They could achieve 16.2% accuracy in predicting meaningful names in a dataset of size 1.2TB of decompiled C code. They predicted variable names that are similar and convey the same meaning as the original names, whereas, we try to predict the original variable names from the decompiled code. In that respect, we keep a much stricter evaluation metric.

The NATURALIZE [1] framework has been quite successful (94% accuracy) in suggesting meaningful identifier names and format styles using n-gram probabilistic models. However, we use different models and training methods for the task.

With the help of machine learning techniques, researchers have made much progress in recovering variable names on decompiled binary code. However, existing techniques suffer from some major drawbacks. DEBIN uses probabilistic models like CRF [19]; Unfortunately, as shown in another paper, its accuracy and generalizability are poor, which renders DEBIN nearly useless in real-world settings [30]. DIRE leverages neural network on Abstract Syntax Trees (ASTs) extracted from collected source code [30]; While DIRE yields

a good prediction performance, the AST extraction process can be tedious and error-prone.

In this paper, we show that with advanced NLP techniques, it is possible to achieve a significantly better performance in accuracy and generalizability of variable name prediction in decompiled binary code. Moreover, our proposed solution does not require extracting ASTs from decompiled code. Instead, it directly uses *raw* decompiled code as input, which is simpler and more robust than DIRE.

2.4 Neural Models

The success of statistical models and progress in the development of stronger neural models led to the application of neural models for identifier name recovery. Our work is related to the following recent works in this area.

In a recent work, CONTEXT2NAME [4] attempted to assign meaningful names to the identifiers based on the context of minified JavaScript codes. They were able to successfully predict 47.5% of meaningful identifiers on 15,000 minified codes using recurrent neural networks. Our work differs from theirs in the sense that we use much advanced deep learning models, we predict the actual original names instead of assigning similar meaningful names and we predict variables for decompiled C code instead of minified JavaScript codes.

Few of the works attempted to predict function names from stripped binaries. Artuso et al. [3] used sequence to sequence networks in two settings (with or without pre-trained embeddings) to predict function names in a dataset created using stripped binaries compiled from 22,040 packages of Ubuntu apt repository with the precision and recall of 0.23 and 0.25 respectively. In another work, the NERO model by David et al. [12] predicted the procedure names in a dataset created from Intel 64-bit executables running on Linux, using LSTM with a precision and recall of 45.82 and 36.40 respectively. Our approach differs from both of the works in the sense that we focus on predicting the variable names instead of functions and our approach to the task.

DIRE [30] is the most recent work related to us. They were able to successfully predict the original variable names 74.3% of the time in the decompiled source code of 164,632 unique x86-64 C binaries mined from Github. They used gated graph neural networks(GGNN) [31] and bidirectional LSTM to encode the Abstract Syntax Tree (structural information) and the decompiled code(lexical information) followed by a decoder network with attention. For our work we use the same dataset published by them but our training method and inputs differ considerably.

2.5 BERT

Bidirectional Encoder Representations from Transformers is *designed to pre-train deep bidirectional representations from the unlabeled text by jointly conditioning on both left and right context in all layers* [14]. BERT is pre-trained in an unsupervised way on a huge collection of natural text for two tasks, first, the masked language modeling (MLM) and second, the next sentence prediction (NSP) to make the model understand the relation between tokens and sentences respectively. Since it's release it has become ubiquitous in almost all tasks in various domains.

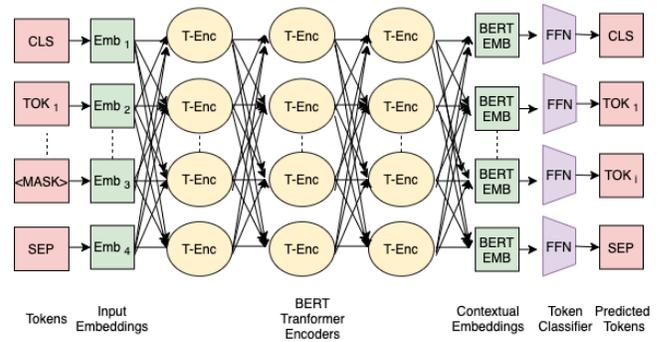


Figure 1: BERT Mask Language Modeling

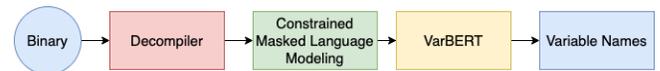


Figure 2: Overall End-to-end Approach

2.6 Vocabulary for Neural Models

All the neural models make use of a specific vocabulary set created from the training data to translate words or tokens to numeric encodings. It is created using a tokenizer on the natural language texts. The size of the vocabulary should not be too low to cover most of the words in training, validation and test data. It should not be too high either to create a massive vocabulary, which may lead to considerably large learned vector embeddings.

2.6.1 WordPiece and SentencePiece Embedding. To keep an optimal size of the vocabulary the WordPiece embedding [50] was introduced. In this embedding, each word is divided into a limited set of common sub-units(sub-words). The word-piece embedding is helpful in handling the rare words in the test samples. BERT uses WordPiece embeddings with a vocabulary of size 30,000. Unknown words are split into smaller units which are present in the vocabulary. Another type of embedding [28] can be directly trained from raw sentences without the need for pre-segmentation of text as compared to earlier approaches. It helps user to create a purely end-to-end and language-independent system.

2.6.2 Byte-Pair Encoding. Byte-Pair Encoding(BPE) [42] is a hybrid between the character and word level text representation. It can handle large vocabulary. The full word is broken down into smaller sub-unit words after performing statistical analysis on the training data. Another implementation of BPE [39] uses bytes instead of unicode characters as base sub-unit words.

3 APPROACH OVERVIEW

Our approach involves the following. We extract the decompiled *raw* code from the binaries. Each function is taken as an independent input instance. We create a corpus by using all such *raw* code to learn a vocabulary of most frequent code-tokens using the technique of Byte-Pair Encoding. We then proceed to learn representations of the code-tokens and a BERT model using the pre-training method of Masked Language Modeling. Finally, we

fine-tune the BERT model using our Constrained Masked Language Modeling technique, to predict only the variable names.

4 DATASET DESCRIPTION

In this work, we use the DIRE[30] dataset. The dataset consists of 3,195,962 decompiled x86-64 functions and their corresponding abstract syntax trees with proper annotations of the decompiled variable name and their corresponding original gold-standard names. The dataset was created by scraping open-source C codes from Github. The codes have been compiled keeping the debug information and then decompiled using Hex-Rays [21], a state-of-the-art industry decompiler. The decompiled variable names are the names assigned to the original variable names by the Hex-Rays. We worked on the reduced preprocessed dataset released by the authors with the same train-dev-test(80:10:10) splits. The number of train, validation and test files is 1,00,632, 12,579 and 12,579 respectively. In the dataset, there are 1,24,702, 1,24,179, and 10,11,054 functions which involve 6,74,854, 6,80,493, and 55,23,045 variables in validation, test and train files respectively.

4.1 Vocabulary

We use Byte-Pair Encoding to learn the vocabulary. We generate a corpus by first replacing all decompiler generated variables in the *raw* code with original variable names and then concatenating all decompiled functions. We use the HuggingFace Tokenizers tool to learn the Byte-Pair Encoding vocabulary. We learn two sets of vocabulary, one with 20,000 tokens, and another with 50,000 tokens to compare the results. Both are generated with allowed mers of 50,000.

5 MODEL DESCRIPTION

5.1 Masked Language Modeling

Traditional Language Modeling is the task of predicting the next word or token given the previous set of words or tokens. Formally, the task of Language modeling can be represented as learning the following probability for all words in a given vocabulary:

$$P(W_i = V_j | W_{i-1}, \dots, W_0) \quad (1)$$

, where W_i are the word or tokens, W_{i-1}, \dots, W_0 is the prior sequence, and V_j are the tokens present in the vocabulary.

In [14], they proposed a different language modeling task to train Transformer Encoders to learn rich representations for tokens. This task, Masked Language Modeling, is defined as follows. Let W_0, \dots, W_N be a sequence of tokens. A random set of m tokens from this sequence is chosen and replaced with a special mask token, $[mask]$. Then we learn the following probability for each of the masked tokens:

$$P([mask]_i = V_j | W_0, \dots, [mask]_{k..}, W_N) \quad (2)$$

, where $[mask]_i$ represents the i^{th} masked token, V_j represents the tokens in the vocabulary set V , $W_0, \dots, [mask]_{k..}, W_N$ is the sequence where m random tokens are masked and $[mask]_k$ is different masked token.

In Masked Language Modeling, not all tokens which are chosen to be replaced with $[mask]$ are replaced in the input. We choose 15% of the token positions at random for prediction. If a token is

selected, we replace 80% of selected tokens with $[mask]$, 10% of the tokens with another random token, and the rest 10% are kept unchanged. The set of masked tokens can change for a particular instance of data across multiple epochs. This is called *dynamic masking*. The model is trained to predict the original token with a cross-entropy loss.

It has been demonstrated in several applications the impact of such a modeling task in learning-rich token vector representations using unsupervised or self-supervised learning methods [14, 34].

5.2 Whole-Word Masking

In Masked Language Modeling, since we use a random sampling of tokens to mask, we may mask sub-parts of a word. As seen in the example, a big word that is not present in the vocabulary can be split into multiple tokens. A subset of these tokens may be chosen to be masked. To learn better representations, and ensure the model captures the entire word, we perform Whole-Word Masking, in which if a word is chosen to be masked, all the tokens of the word are masked and the model is trained to predict all the original token.

5.3 Constrained Masked Language Modeling

Constrained Masked Language Modeling is a variation of Masked Language Modeling where the tokens are not masked at random. We define Constrained Masked Language Modeling as follows:

Let W_0, \dots, W_N be a sequence of tokens. Let $C = \{A_0, \dots, A_c\}$ be a set of tokens which we define as the constrained set of tokens. Then in constrained masked language modeling, all tokens in W_0, \dots, W_N which belong to C are masked. C is a subset of V . Here all the tokens are replaced with a masked token, unlike Masked Language Modeling. Similar to Masked Language Modeling, we train the model to predict the masked token with a cross-entropy loss:

$$-\sum_{i=0}^N \sum_{c=1}^M y_{w_i, v_c} \log(p_{w_i, v_c}) \quad (3)$$

, where M is the size of the vocabulary, w_i is the current token, v_c is the target token, y is an indicator variable which is 1 if the target is v_c and 0 otherwise and p is the probability w_i is same as v_c .

5.4 Count of Token Prediction

In the task of Variable Name Recovery, during test time, we do not know the number of tokens the original variable should have. We define the Count of Token Prediction as the task to determine the count of mask tokens during test time. We solve this with the following heuristic defined in Algorithm 1. We evaluate our BERT models with both an Oracle model which gives us the number of tokens and our heuristics. The max allowed number of tokens is derived from a statistical analysis of the different variable names present in the Train dataset.

5.5 Technical Details

A Transformer architecture [46] contains multiple layers, with each block using multiple self-attention heads. We train two different versions of VarBERT each differing in the number of layers L , the number of self-attention heads A and the hidden dimension H . The vocabulary for both versions is the same. Training such a huge

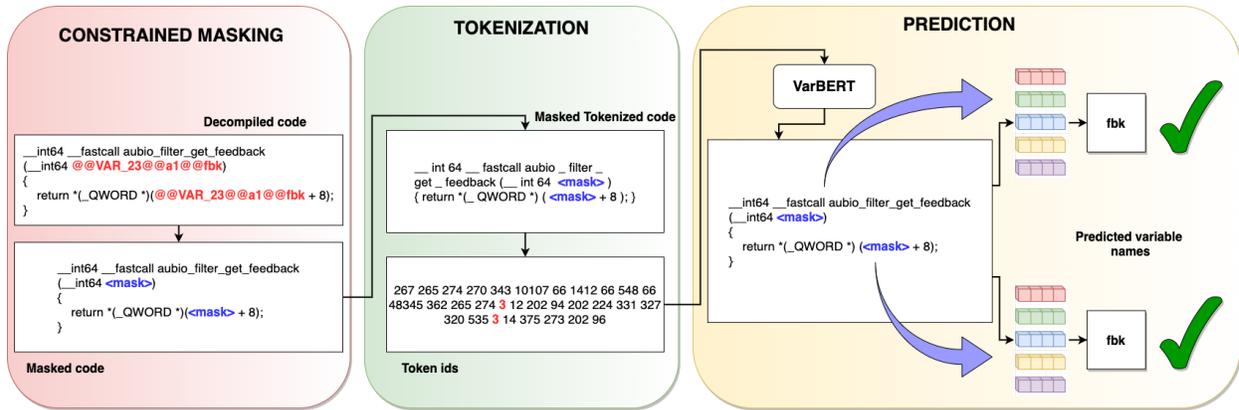


Figure 3: End-to-end Constrained Masked Language Modeling

```

Result: BestPredictedVariable
Count=1;
MaxProbability=0;
BestPredictedVariable="";
MaxAllowedToken;
Model;
while Count<MaxAllowedToken do
    Place count number of mask tokens;
    CurrVariable=Model.Predict;
    P=Mean of the probabilities for each token;
    if P>MaxProbability then
        BestVariable=CurrVariable;
        MaxProbability=P;
    end
end

```

Algorithm 1: Best Variable Selection Heuristics

architecture requires a lot of computing and training time. We define a fixed hyper-parameter budget of two training runs over the entire dataset for each of the versions to limit our training time and compute cost. We take initial hyper-parameters from RoBERTa [34].

5.5.1 *VarBERT-Base.* This version has $L = 12$, $A = 12$ and $H = 768$. This leads to total parameters to be trained to be nearly 125 million. The maximum number of input tokens this version can take is 512.

5.5.2 *VarBERT-Small.* We hypothesize in the task of Variable Name Recovery, a smaller model in terms of depth and parameters might perform reasonably well, and hence create this version which has $L = 6$, $A = 8$ and $H = 512$. We though increase the number of input tokens the model can take to 1024 to be able to train the model with longer code sequences and be more useful for the community. The number of trainable parameters for this model is 45 million, which is 2.5 times smaller than VarBERT-Base.

5.6 Training Methodology

For the task of Variable Name Recovery, we train our models with the following methodologies.

5.6.1 *Pre-Training.* The need and impact of Pre-Training on auxiliary tasks before the final intended task has been demonstrated in multiple previous natural language processing and computer vision work such as in [14, 34, 44, 53]. For the task of Variable Name Recovery, we define the Masked Language Modeling task as the pre-training task. We follow the same steps as done in RoBERTa [34] to train our model and learn rich representations for our code-tokens. We train, both with and without Whole-word Masking.

5.6.2 *Finetuning.* We finetune our models, for the Variable Name Recovery task using the Constrained Language Modeling task.

5.6.3 *Optimization.* We optimize our models using BERTAdam [14, 27] with following parameters: $\beta_1 = 0.9$ and $\beta_2 = 0.999$, $\epsilon = 1e - 6$ and L_2 weight decay of 0.01. We warmup over first 10,000 steps to a peak value of $1e - 4$ and then linearly decay. We set our dropout to 0.1 on all layers and attention weights. Our activation function is GELU [20].

We use HuggingFace Transformers and Facebook FairSeq tools to train our models.

6 EXPERIMENTS

6.1 Experimental Setup

Our models are evaluated on the DIRE dataset. Split of the dataset given in DIRE [30] is 80:10:10. We use only the *raw* code generated from the decompilers and replace the decompiler generated variable names with the original variable name given by the developer. We evaluate the impact of pre-training and whole-word masking. We also evaluate the impact of constrained masked language modeling using the pre-trained only model as a baseline. We evaluate the impact of the size of the vocabulary and the size of the models.

As defined above, we train the models with a hyper-parameter budget of two, so a bigger budget might result in even better performance. The hyperparameters for pre-training and finetuning are a batch size of 1024 and the number of epochs is restricted to 40. We use 4 Nvidia Volta V-100 16GB graphics cards to train our models. Approximately, the VarBERT-Base has a training time of 72 hours and VarBERT-Small has a training time of 38 hours.

6.2 Metrics

We evaluate our models using the following metrics. *Exact Match accuracy* of predicting the correct variable name. The final goal of such a model is to provide suggestions to system reverse engineers, and our models also provide a ranking of tokens, we measure the accuracy of the correct variable at different ranks, which are 1,3,5 and 10. In DIRE [30] they also measure the character error rate (CER) metric, which calculates the edit distance between the original and predicted names, then normalizes the length of the original name, as defined in CER [48]. We measure this for our heuristic-driven Count of Token Prediction algorithm. We also measure the **Perplexity** of both language modeling tasks. Perplexity measures how well a probabilistic language model predicts a target token.

The binaries in the dataset use C libraries. The different splits Train, Validation and Test contain binaries that share these libraries. The functions in these libraries have the same variable names and bodies. To better understand the generalizability of our proposed models and techniques, we also report the metrics on two sets, *Body not in the train* and *Overall*. These sets are defined as meta-tags in the DIRE dataset.

6.3 Baselines

Our baseline models for the Variable Name Recovery tasks are the following.

6.3.1 DEBIN. It uses statistical models such as CRFs and Extremely Randomized Trees with several handcrafted features to recover variable names, along with other debugging information from stripped binaries. Few of the handcrafted features are functions used, registers used, types, flags, instructions and relationships between functions, variables, and types. This acts as a weak baseline for our proposed models.

6.3.2 DIRE. It uses LSTMs, Gated Graph Neural Networks and Attention Mechanism with an Encoder-Decoder architecture to generate variable names. It also uses the *raw* code as one of the inputs. In addition to the decompiled code, it uses a GNN based structural encoder to encode Abstract Syntax Trees. This acts as a strong neural baseline for our proposed models.

Both the above models, use the *raw* code to recover variable names but also use additional features. In our proposed models, we only use the *raw* code.

7 DISCUSSION AND ANALYSIS

7.1 Dataset Analysis

We start with the initial analysis of the dataset. We measure the overlap between the variables present in the Train, Validation and Test set splits with our learned vocabulary. We also measure the length of the functions. Figure 4 and Figure 5 show the respective distributions. 216 and 167 functions from the Validation and Test set were truncated due to the limitation of maximum allowed sequence length of 1024. Most of the variables in the dataset are tokenized to a wide range of [1,7]. We use these statistics to define the *MaxAllowedToken* parameter in Algorithm 1.

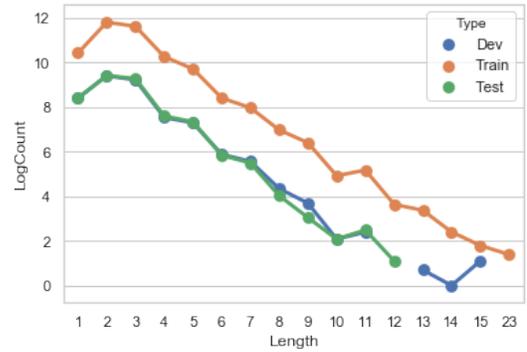


Figure 4: Distribution of length of gold tokenized variable names in Train, Validation and Test data. The counts are in Natural Log Scale.

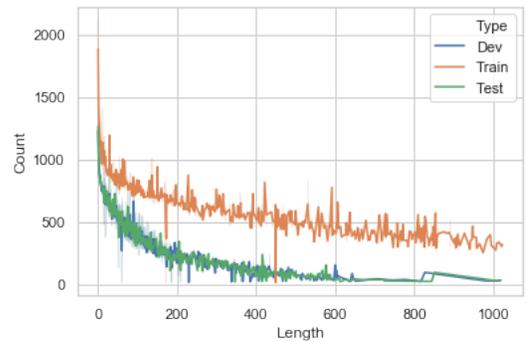


Figure 5: Distribution of length of function in Train, Validation and Test data

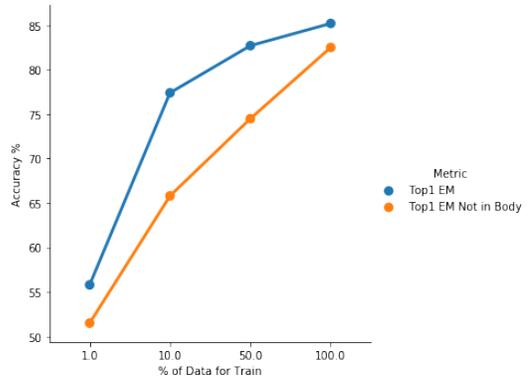


Figure 6: Impact of training corpus size on the performance of our model. The scores are for VarBERT-Small trained with CMLM, 50K Vocab, and Heuristics.

7.2 Model Analysis

Can we use just Masked Language Modeling for Variable Name Recovery? In Table 1, we compare the different modeling tasks. We observe that Constrained Masked Language Modeling

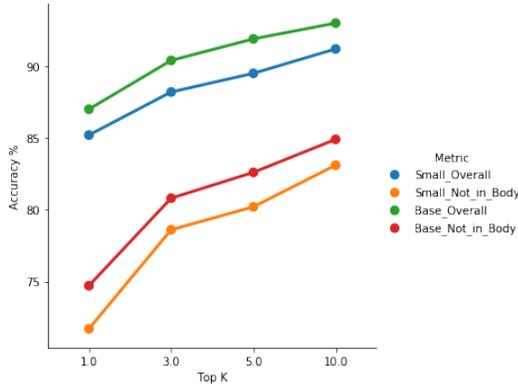


Figure 7: Impact of training corpus size on the performance of our model. The scores are for VarBERT-Small trained with CMLM, 50K Vocab, and Heuristics based Count of Token Prediction.

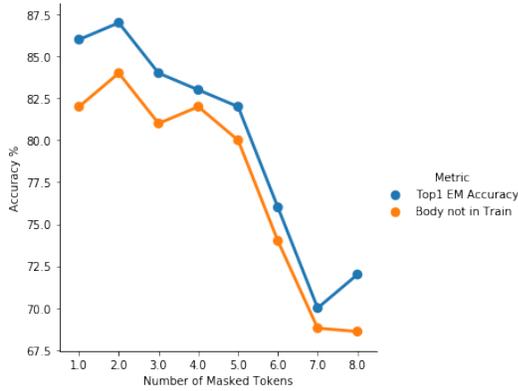


Figure 8: Accuracy of VarBert-Small for different lengths of target variable names.

with Whole-word Masking performs the best. It is also observed that just Masked Language Modeling performs significantly poor. This can be explained by the random 15% of tokens which are masked, and not all the variable name tokens. Masked Language Modeling with Whole-word Masking is slightly better. Pre-training also has a significant impact, as without Pre-training model Top1 Exact Match Accuracy are in the lower 50s, with a high perplexity. This is because the task of Constrained Masked Language Modeling is insufficient to learn rich vector representations for the entire vocabulary set.

Does increasing vocabulary size help? In Table 2 we compare our models with two different sets of vocabulary learned using Byte-Pair Encoding. We observe increasing vocabulary has a significant impact. With a smaller vocabulary, we can have a smaller model that trains faster, but we lose in performance by a significant margin.

How do VarBERT-base and VarBERT-small compared to each other? Our hypothesis of a smaller model with lesser parameters being able to perform reasonably well is demonstrated across all Tables. The VarBERT-base model does perform better but given

Metric	Model	MLM	MLM _w	CMLM	CMLM _w
Top1 EM % ↑	Small	0	1.2	83.5	85.2
	Base	5	5.5	86.5	87.0
Perplexity ↓	Small	1178	1033	1.07	1.067
	Base	932	899	1.58	1.081
CER % ↓	Small	98	97.8	18.32	16.4
	Base	97.2	96.5	16.54	15.65

Table 1: Comparison of Masked Language Modeling (MLM) and Constrained MLM (CMLM) on downstream Variable Name Recovery Validation Set. Vocab Size is 50k, and Token Prediction is through heuristics. W refers to being trained with whole-word masking.

Metric	Model Type	Vocab-20k	Vocab-50k
Top1 EM ↑	Small	78.4	85.2
	Base	80.2	87.0
Perplexity ↓	Small	1.832	1.067
	Base	1.778	1.081
CER ↓	Small	19.6	16.4
	Base	18.78	15.65

Table 2: Comparison of Vocabulary Size on downstream Variable Name Recovery Validation Set. Models are fine-tuned using CMLM, Token Prediction is through heuristics.

Metric	Model Type	Oracle	Heuristics
Top1 EM ↑	Small	91.1	85.2
	Base	92.6	87.0
Perplexity ↓	Small	1.054	1.067
	Base	1.077	1.081
CER ↓	Small	11.2	16.4
	Base	10.6	15.65

Table 3: Comparison of Count of Token Prediction algorithm on Variable Name Recovery Validation Set.

Data Splits	Model Type	Top1 EM ↑	CER ↓
Overall	Small	85.2	16.50
	Base	87.0	15.65
Body in Train	Small	86.7	15.40
	Base	89.9	14.34
Body not in Train	Small	71.8	26.70
	Base	74.7	24.20

Table 4: Comparison of the two trained models, VarBERT-Small and VarBERT-Base on the different splits of the Validation Sets. Models are trained with CMLM and Token prediction heuristics is used to predict count.

a restricted hyper-parameter budget, we observe the delta between the performances is within 2-3 % in Exact Match Accuracy. It should be noted VarBERT-small has half the number of layers compared to VarBERT-base but has double the max sequence length allowed. This makes VarBERT-small more useful as it can take input longer sequences of code.

Metric 1% Data	DEBIN	DIRE	Small
Train Time (h) ↓	2.4	13	3
Accuracy % - Overall ↑	2.4	32.2	55.8
Accuracy % - Body in Train ↑	3.0	40.0	56.8
Accuracy % - Body not in Train ↑	0.6	5.3	51.5

Table 5: Comparisons with Baselines Trained on 1% data. VarBERT-Small is pre-trained with MLM and finetuned with CMLM, with a vocab of 50k and Token Prediction is with heuristics. The scores are on the Validation Set.

Method	Top1 EM % ↑	CER % ↓	Body Not in Train % ↑
DIRE	74.0	28.3	35.30
VarBERT-Small	83.10	17.8	82.47
VarBERT-Base	86.36	15.4	84.15

Table 6: Comparison with the current state-of-the-art on the Test Set. Body Not In Train is also EM Accuracy.

How does our Heuristic Based Token Count Prediction Algorithm perform? From Table 3 it can be seen that if we know the correct number of masked tokens, the model can predict the original variable name with very high accuracy. The heuristics-based algorithm although performs reasonably well, it has a considerable room to improve.

How does our performance improve if we suggest users with Top K Variable Names? Figure 7 shows our performance improves by nearly 7-9% when we suggest Top K variable names. This shows the promise of using suggestion and ranking based models for variable name recovery.

How does our model improve with increasing train data? Figure 6 shows the learning curve of our VarBERT-Small model trained with Constrained Masked Language Modeling and a vocabulary size of 50K. We observe an interesting curve where the overall Top1 accuracy increases more than the *Body not in Train* set. This indicates with more data model learns to memorize the different variable names in shared library functions.

How does our models compare to existing baseline and state-of-the-art models? Table 5 and Table 6 show our model performs significantly well compared to both the models across all metrics. In Table 5, we compare the baselines with models trained on 1% of the Training data corpus. We choose this approach as DEBIN requires a considerable amount of training time and does not scale to the entire dataset. We compare with DIRE by training with the entire dataset in Table 6. VarBERT-small is more robust and data-efficient compared to both DEBIN and DIRE. Our pre-trained on Masked LM and finetuned on the Constrained MLM model generalizes significantly well, which is shown on the Exact Match Accuracy in the *Body not in Train* set. VarBERT-small is a more accurate and generalizable technique compared to the current state-of-the-art.

What is our performance for variables which are split into multiple tokens? Figure 8 shows the accuracy of our model across variables that are split into multiple tokens. It is expected for the model to perform well on variables that are not split as they have

rich representations learned during the pre-training task. It is interesting to note that the model performs reasonably well for variables split into two to five tokens. This shows the model has sufficient reasoning capabilities to predict such variable names.

Where does our model go wrong?

We analyzed the different errors the model make and we classify the errors broadly into the following categories: Error in the number of mask count prediction, Partial incorrect token prediction, and Off-by-few-chars errors. Error in the number of mask count prediction occurs when a different set of mask tokens count has a higher average probability. It is observed that a smaller number of tokens have a higher average probability. We can fix this issue if we learn this task instead of a heuristic-based algorithm.

Partial incorrect token prediction happens when a part of the tokens generated is wrong. This happens when there can be several generated combinations with the same prefix token. For example, variable names like *tokenIndex* and *tokenCount*. Our model predicts *tokenIndex* when instead, the original variable is *tokenCount*. This can be corrected by a suggestion based system and it gives the user the freedom to select which variable name is more suited.

Off-by-few-chars errors occur when the model predicts tokens which only differ by a few characters. For example, *token* and *tokens*. Both are present in our vocabulary, and the model predicts both in the top two positions together. There are several such pairs, and not restricted to only a single length of split-tokens.

What are the other sources of errors? The DIRE dataset comprises of scraped C code from GitHub. There was no filtering and quality control implemented to ensure that the dataset represents the set of target binaries that are reverse-engineered. Moreover, these binaries were not compiled with multiple different optimization and obfuscation compiler options. There is a possibility that our models may not perform on such code with such high accuracy but may work reasonably well compared to our baselines, and current state-of-the-art. To improve on such binaries is left as future work. Our model is tightly coupled with the output of the Hex-Rays decompiler. To support other decompilers, the model may need to be re-trained with a new training corpus generated from the new set of decompilers. We leave building an adaptable and multi-decompiler supporting model as future work.

8 CONCLUSION

Improving the understandability of decompiled binary code is crucial for software reverse engineering. In this paper, we have advanced the state of the art for variable name recovery in decompiled binary code. We adapted recent advances in the field of natural language processing, such as Masked Language Modeling and Transformers, and proposed a heuristic-based Count of Token Prediction algorithm. These techniques are crucial in improving the performance of variable name recovery in decompiled binary code.

In our evaluation, we showed the impact of each module. We trained two neural network models, VarBERT-base and VarBERT-small. These neural models take *raw* code as input, which makes them much simpler to build and use. Our evaluation of the DIRE data set shows that our techniques advance the state-of-the-art by 12.36% on overall accuracy and improve on generalizability by 49%.

REFERENCES

- [1] Miltiadis Allamanis, Earl T Barr, Christian Bird, and Charles Sutton. 2014. Learning natural coding conventions. In *Proceedings of the 22nd ACM SIGSOFT International Symposium on Foundations of Software Engineering*. 281–293.
- [2] Miltiadis Allamanis, Earl T Barr, Premkumar Devanbu, and Charles Sutton. 2018. A survey of machine learning for big code and naturalness. *ACM Computing Surveys (CSUR)* 51, 4 (2018), 1–37.
- [3] Fiorella Artuso, Giuseppe Antonio Di Luna, Luca Massarelli, and Leonardo Querzoni. 2019. Function Naming in Stripped Binaries Using Neural Networks. *arXiv preprint arXiv:1912.07946* (2019).
- [4] Rohan Bavishi, Michael Pradel, and Koushik Sen. 2018. Context2Name: A deep learning-based approach to infer natural variable names from usage contexts. *arXiv preprint arXiv:1809.05193* (2018).
- [5] Tony Beltramelli. 2018. pix2code: Generating code from a graphical user interface screenshot. In *Proceedings of the ACM SIGCHI Symposium on Engineering Interactive Computing Systems*. 1–6.
- [6] Sahil Bhatia and Rishabh Singh. 2016. Automated correction for syntax errors in programming assignments using recurrent neural networks. *arXiv preprint arXiv:1603.06129* (2016).
- [7] Marcel Bruch, Martin Monperrus, and Mira Mezini. 2009. Learning from examples to improve code completion systems. In *Proceedings of the 7th joint meeting of the European software engineering conference and the ACM SIGSOFT symposium on the foundations of software engineering*. 213–222.
- [8] Juan Caballero and Zhiqiang Lin. 2016. Type Inference on Executables. *Comput. Surveys* 48, 4 (2016), 1–35. <https://doi.org/10.1145/2896499>
- [9] Joshua Charles Campbell, Abram Hindle, and José Nelson Amaral. 2014. Syntax errors just aren't natural: improving error reporting with language models. In *Proceedings of the 11th Working Conference on Mining Software Repositories*. 252–261.
- [10] Luigi Cerulo, Michele Ceccarelli, Massimiliano Di Penta, and Gerardo Canfora. 2013. A hidden markov model to detect coded information islands in free text. In *2013 IEEE 13th International Working Conference on Source Code Analysis and Manipulation (SCAM)*. IEEE, 157–166.
- [11] F Chagnon. [n.d.]. IDA-Decompiler.
- [12] Yaniv David, Uri Alon, and Eran Yahav. 2019. Neural reverse engineering of stripped binaries. *arXiv preprint arXiv:1902.09122* (2019).
- [13] Premkumar Devanbu. 2015. New initiative: the naturalness of software. In *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, Vol. 2. IEEE, 543–546.
- [14] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805* (2018).
- [15] The ghidra decompiler. 2019. The ghidra decompiler. <https://ghidra-sre.org/>
- [16] Rahul Gupta, Aditya Kanade, and Shirish Shevade. 2018. Deep reinforcement learning for programming language correction. *arXiv preprint arXiv:1801.10467* (2018).
- [17] Rahul Gupta, Soham Pal, Aditya Kanade, and Shirish Shevade. 2017. Deepfix: Fixing common c language errors by deep learning. In *Thirty-First AAAI Conference on Artificial Intelligence*.
- [18] Tihomir Gvero and Viktor Kuncak. 2015. Synthesizing Java expressions from free-form queries. In *Proceedings of the 2015 ACM SIGPLAN International Conference on Object-Oriented Programming, Systems, Languages, and Applications*. 416–432.
- [19] Jingxuan He, Pesho Ivanov, Petar Tsankov, Veselin Raychev, and Martin Vechev. 2018. Debin: Predicting debug information in stripped binaries. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security*. 1667–1680.
- [20] Dan Hendrycks and Kevin Gimpel. 2016. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415* (2016).
- [21] SA Hex-Rays. 2013. Hex-Rays Decompiler.
- [22] Abram Hindle, Earl T Barr, Zhenqiong Su, Mark Gabel, and Premkumar Devanbu. 2012. On the naturalness of software. In *2012 34th International Conference on Software Engineering (ICSE)*. IEEE, 837–847.
- [23] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [24] Xing Hu, Yuhan Wei, Ge Li, and Zhi Jin. 2017. CodeSum: Translate program language to natural language. *arXiv preprint arXiv:1708.01837* (2017).
- [25] Srinivasan Iyer, Ioannis Konstas, Alvin Cheung, and Luke Zettlemoyer. 2016. Summarizing source code using a neural attention model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 2073–2083.
- [26] Alan Jaffe, Jeremy Lacomis, Edward J Schwartz, Claire Le Goues, and Bogdan Vasilescu. 2018. Meaningful variable names for decompiled code: A machine translation approach. In *Proceedings of the 26th Conference on Program Comprehension*. 20–30.
- [27] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [28] Taku Kudo and John Richardson. 2018. Sentencepiece: A simple and language independent subword tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226* (2018).
- [29] Nate Kushman and Regina Barzilay. 2013. Using semantic unification to generate regular expressions from natural language. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 826–836.
- [30] Jeremy Lacomis, Pengcheng Yin, Edward Schwartz, Miltiadis Allamanis, Claire Le Goues, Graham Neubig, and Bogdan Vasilescu. 2019. Dire: A neural approach to decompiled identifier naming. In *2019 34th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 628–639.
- [31] Yujia Li, Daniel Tarlow, Marc Brockschmidt, and Richard Zemel. 2015. Gated graph sequence neural networks. *arXiv preprint arXiv:1511.05493* (2015).
- [32] Wang Ling, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kočíský, Andrew Senior, Fumin Wang, and Phil Blunsom. 2016. Latent predictor networks for code generation. *arXiv preprint arXiv:1603.06744* (2016).
- [33] Han Liu. 2016. Towards better program obfuscation: optimization via language models. In *2016 IEEE/ACM 38th International Conference on Software Engineering Companion (ICSE-C)*. IEEE, 680–682.
- [34] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11629* (2019).
- [35] Cyrus Omar. 2013. Structured statistical syntax tree prediction. In *Proceedings of the 2013 companion publication for conference on Systems, programming, & applications: software for humanity*. 113–114.
- [36] Sebastian Proksch, Johannes Lerch, and Mira Mezini. 2015. Intelligent code completion with Bayesian networks. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 25, 1 (2015), 1–31.
- [37] Yewen Pu, Karthik Narasimhan, Armando Solar-Lezama, and Regina Barzilay. 2016. sk_p: a neural program corrector for MOOCs. In *Companion Proceedings of the 2016 ACM SIGPLAN International Conference on Systems, Programming, Languages and Applications: Software for Humanity*. 39–40.
- [38] Maxim Rabinovich, Mitchell Stern, and Dan Klein. 2017. Abstract syntax networks for code generation and semantic parsing. *arXiv preprint arXiv:1704.07535* (2017).
- [39] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners. *OpenAI Blog* 1, 8 (2019), 9.
- [40] Baishakhi Ray, Vincent Hellendoorn, Saheel Godhane, Zhaopeng Tu, Alberto Bacchelli, and Premkumar Devanbu. 2016. On the "naturalness" of buggy code. In *2016 IEEE/ACM 38th International Conference on Software Engineering (ICSE)*. IEEE, 428–439.
- [41] Veselin Raychev, Martin Vechev, and Andreas Krause. 2015. Predicting program properties from "big code". *ACM SIGPLAN Notices* 50, 1 (2015), 111–124.
- [42] Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909* (2015).
- [43] Abhishek Sharma, Yuan Tian, and David Lo. 2015. NIRMAL: Automatic identification of software relevant tweets leveraging language model. In *2015 IEEE 22nd International Conference on Software Analysis, Evolution, and Reengineering (SANER)*. IEEE, 449–458.
- [44] Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490* (2019).
- [45] Bogdan Vasilescu, Casey Casalnuovo, and Premkumar Devanbu. 2017. Recovering clear, natural identifiers from obfuscated JS names. In *Proceedings of the 2017 11th Joint Meeting on Foundations of Software Engineering*. 683–693.
- [46] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [47] Song Wang, Devin Chollak, Dana Movshovitz-Attias, and Lin Tan. 2016. Bugram: bug detection with n-gram language models. In *Proceedings of the 31st IEEE/ACM International Conference on Automated Software Engineering*. 708–719.
- [48] Weiyue Wang, Jan-Thorsten Peter, Hendrik Rosendahl, and Hermann Ney. 2016. Character: Translation edit rate on character level. In *Proceedings of the First Conference on Machine Translation: Volume 2, Shared Task Papers*. 505–510.
- [49] Martin White, Michele Tufano, Christopher Vendome, and Denys Poshyvanyk. 2016. Deep learning code fragments for code clone detection. In *2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 87–98.
- [50] Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, et al. 2016. Google's neural machine translation system: Bridging the gap between human and machine translation. *arXiv preprint arXiv:1609.08144* (2016).
- [51] Shir Yadid and Eran Yahav. 2016. Extracting code from programming tutorial videos. In *Proceedings of the 2016 ACM International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software*. 98–111.

- [52] Khaled Yakdan, Sebastian Eschweiler, Elmar Gerhards-padilla, and Matthew Smith. 2015. No More Gotos: Decompilation Using Pattern-Independent Control-Flow Structuring and Semantics-Preserving Transformations. In *Network and Distributed System Security*. 8–11.
- [53] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. In *Advances in neural information processing systems*. 5754–5764.
- [54] Pengcheng Yin and Graham Neubig. 2017. A syntactic neural model for general-purpose code generation. *arXiv preprint arXiv:1704.01696* (2017).

A SUPPLEMENTAL MATERIAL