On the Impact of Multiple Source Code Representations on Software Engineering Tasks - An Empirical Study

Karthik Chandra Swarna^{a,*}, Noble Saji Mathews^a, Dheeraj Vagavolu^a, Sridhar Chimalakonda^a

^aResearch in Intelligent Software & Human Analytics (RISHA) Lab, Department of Computer Science and Engineering, Indian Institute of Technology Tirupati, India

Abstract

Efficiently representing source code is crucial for various software engineering tasks such as code classification and clone detection. Existing approaches primarily use *Abstract Syntax Tree* (AST), and only a few focus on semantic graphs such as *Control Flow Graph* (CFG) and *Program Dependency Graph* (PDG), which contain information about source code that AST does not. Even though some works tried to utilize multiple representations, they do not provide any insights about the costs and benefits of using multiple representations. The primary goal of this paper is to discuss the implications of utilizing multiple code representations, specifically AST, CFG, and PDG. We modify an AST path-based approach to accept multiple representations as input to an attention-based model. We do this to measure the impact of additional representations (such as CFG and PDG) over AST. We evaluate our approach on three tasks: *Method Naming, Program Classification*, and *Clone Detection*. Our approach increases the performance on these tasks by **11%** (F1), **15.7%** (Accuracy), and **9.3%** (F1), respectively, over the baseline. In addition to the effect on performance, we discuss timing overheads incurred with multiple representations. We envision this work providing researchers with a lens to evaluate combinations of code representations for various tasks.

Keywords: Source Code Representation, Abstract Syntax Tree, Control Flow Graph, Program Dependence Graph, Code Embedding, Method Naming

1. Introduction

Due to the drastic rise in the availability of enormous volumes of source code in open-source projects and tools to extract and analyze them [1, 2, 3], researchers have explored many ways of solving software engineering problems such as *code classification* [4, 5], *code clone detection* [6, 7, 8], *code summarization* [9, 10], *method name prediction* [11], and *source code retrieval* [12, 13]. Representing source code is central to all of these software engineering tasks and is influential in determining the performance of the approaches [14, 15].

Many existing works utilize tree and graph-based representations such as *Abstract Syntax Tree* (AST) [7, 6, 14, 16, 17, 18, 19, 20], *Control Flow Graph* (CFG) [21, 22], or *Program Dependency Graph* (PDG) [23] to represent source code for various tasks. Since programming languages have strict and well-defined grammar, researchers have traditionally used formal methods to process and reason about source code. These formal approaches involve using algorithmic or mathematical techniques over the AST, CFG, or PDG [24, 6, 25]. Due to advancements in deep learning, researchers have started focusing on automatically learning program properties to make software more understandable and maintainable [15]. These learning-based techniques improved performance over traditional methods [15].

Though the learning-based techniques provide superior performance to traditional methods, most of these works are limited to specific aspects of source code. For example, as Table 1 highlights, AST dominates most of the research, and only a few works use representations such as CFG or PDG. Two critical factors for a learning-based technique are the learning model and the data used to train the model. Data in the current context refers to the code representations

^{*}Corresponding Author

Email addresses: cs17b026@iittp.ac.in (Karthik Chandra Swarna), ch19b023@iittp.ac.in (Noble Saji Mathews), cs17b028@iittp.ac.in (Dheeraj Vagavolu), ch@iittp.ac.in (Sridhar Chimalakonda)

used to train the model. We experimented with both these factors to arrive at an optimal model and input data to maximize performance. However, software engineering literature traditionally focused on optimizing the model rather than enhancing the input data to the model. One way to enhance the input data is by using a combination of representations, including various syntactic and semantic code representations. We can include multiple representations in the dataset because, unlike traditional methods, learning-based methods can selectively give importance to those program features that improve the model's performance and neglect the ones that do not. With learning-based approaches, we can offload the responsibility of choosing a suitable representation for a task to the model.

Representation	Work	Task(s)	Venue
	Deckard [6]	Code clone detection	ICSE
	White et al. [7]	Code clone detection	ASE
	Wei et al. [26]	Code clone detection	IJCAI
	TBCNN [4]	Code classification	AAAI
AST	ASTNN [14]	Code classification, Code clone detection	ICSE
	MTN [17]	Code classification, Code clone detection	TOSEM
	code2vec [27]	Method naming	POPL
	code2seq [28]	Code summarization, Code captioning	ICLR
	InferCode [18]	de [18] Code clone detection, Code classification,	
	Intercode [10]	Code clustering, Code search, Method naming	
CEC	Sun et al. [21]	Code clone detection	IFIP
CFG	Phan et al. [22]	Software defect prediction	ICTAI
PDG Li et al. [23]		Bug detection	OOPSLA
	Allamanis et al. [29]	Variable name prediction	ICLR
Combination of	Yamaguchi et al. [25]	Software vulnerability detection	SSP
representations	Long et al. [30]	Algorithm classification	AAAI

Table 1: An overview of the literature that utilizes AST / CFG / PDG

Few approaches have used different graph structures together to represent source code [25, 29]. The Code Property Graph (CPG) by Yamaguchi et al. [25] is one such approach where they integrated CFG, PDG, and AST into a joint *static* data structure to model and detect *software vulnerabilities*. The CPG helps express common code vulnerabilities for specific applications, which can be queried by graph traversal queries [31]. Similarly, Allamanis et al. [29] combine AST, control flow, and data flow edges into a single graph, where these graph edges are particularly designed for the variable name prediction task [27]. This approach needs execution trace information to create graphs, which is an additional overhead [17]. Although these works use semantic graphs and perform well on the intended tasks, they *do not* provide any insight into the effects of using semantic graphs along with AST. Besides, manually choosing the appropriate edges for a task can be laborious and cannot be easily extended for other tasks.

Thus, in this paper, we try to answer the question: "How does utilizing multiple source code representations affect the performance of diverse software engineering tasks?" We do this by extending an existing AST-based approach to semantic graphs like CFG and PDG, then analyzing the impact of these representations on the model's performance. The technique we use in this work involves representing a code snippet as a set of paths extracted from its AST. We selected this path-based technique since there is no need manually design program features for a specific task or language. Moreover, extending this approach to CFG and PDG would allow us to measure the performance boost provided by these semantic graphs. Code2vec [27] is a popular work that uses a similar technique wherein these AST paths are used to train an attention-based neural network. This neural network encodes the code's structural information from individual AST paths by generating method-level code embeddings, which can then be used for various downstream tasks. However, like most existing approaches, their works do not emphasize semantic graphs such as CFG and PDG.

Our work introduces a way to extract paths from semantic graphs CFG and PDG. We utilize paths from both

syntactic and semantic representations to work on software engineering tasks. We extend the *code2vec's* attention mechanism to work with paths extracted from multiple code representations, mainly AST, CFG, and PDG. Utilizing the attention mechanism to handle CFG and PDG paths may allow the model to capture the code semantics better, which may not be possible with AST paths alone. This *Aggregate Attention model* gives optimal weights to individual paths extracted from AST to compute a weighted average. It also parallelly calculates the weighted averages of CFG and PDG paths. These three individual weighted averages are concatenated to generate a code embedding. The attention mechanism allows the model to learn how much importance/weight it should give to each path. This mechanism enables the neural network to capture subtle similarities and differences between code snippets, which would be impossible if all paths were given equal weights without using attention.

The reason we take an existing AST-only approach and make it workable with multiple representations is to understand the impact that utilizing multiple representations has on the model's performance. Our goal is *not* to introduce a better learning model but to investigate the possibility of improving the model's performance by experimenting with multiple representations. We chose *code2vec* because it is neither task nor language-dependent and is easily extensible to various tasks like program classification and clone detection. Moreover, it is still widely used, and many recent works are built upon it [28, 32, 33, 34].

The core reasoning of this paper is to explore the idea of integrating semantic structures (such as CFG and PDG) with syntactic structures (AST) and study the effect on the performance of various software engineering tasks.

As a first step, we evaluated our approach on the task of Method Naming using a custom C language dataset of around 730K methods [35]. The comparative evaluation of *code2vec* and our model shows an increase of **11%** in the F1 score on the whole dataset. However, when evaluated on different projects in our dataset, we see an increase of F1 score up to 100%. This paper is an extension of our preliminary *New Ideas and Emerging Results* (NIER) paper [35], which briefly demonstrates this approach on the method naming task. This paper extends our approach to program classification and code clone detection tasks to find the generality across multiple tasks. We also have reimplemented our Path Extractor tool (refer to Section 3.2) to extend its functionality for these new tasks.

Our experiments show that our approach improves performance over the AST-only approach for these tasks. Our approach improves the accuracy by **15.7%** for program classification and the F1 score by **9.3%** for code clone detection. Moreover, our results show that diverse approaches built upon *code2vec* for various tasks [28, 32, 33, 34] can be enhanced by including paths from CFG and PDG. This paper discusses the additional timing overheads incurred and other implications of our approach. We also discuss threats to the validity of our results and some ways to overcome the shortcomings of our approach. Thus, we believe considering a combination of the syntactic and semantic structures can lead to a new direction in representing source code while also improving existing works that rely solely on AST, CFG, or PDG. We make all our code and data available on GitHub.¹

The main contributions of our work are as follows:

- We introduce a novel way of extracting and utilizing semantic paths from CFG and PDG to represent code.
- We employ an attention neural network to learn code embeddings from syntactic and semantic paths. We apply our approach to three tasks: *Method Naming*, *Program Classification*, and *Code Clone Detection*.
- We demonstrate that integrating semantic structures (such as CFG and PDG) with syntactic structures (AST) can predominantly improve the performance of these tasks by 11% (F1), 15.7% (Accuracy), and 9.3% (F1), respectively.

The remainder of this paper is organized as follows: Section 2 establishes the necessary background. Section 3 defines semantic paths and explains our attention pipeline. Section 4 describes three use cases, demonstrates our experiments, and reports the results. Section 5 breaks down the results and discusses the main research questions. Section 6 discusses the limitations and threats to the validity of our work. Section 7 presents the related work. Finally, we conclude the paper in Section 8 with potential implications for further work.

¹https://github.com/NobleMathews/mocktail-blue-lagoon

2. Background

2.1. Intermediate Code Representations

Researchers in program analysis and compiler design have developed multiple intermediate code representations over the years [36, 25]. These representations can also be used in other related fields as they represent the program properties well [25]. In this work, we use three of the well-known intermediate representations, namely - *Abstract Syntax Tree (AST), Control Flow Graphs (CFG)*, and *Program Dependence Graphs (PDG)*. We chose CFG and PDG to combine with AST because they are the most commonly used representations in literature after AST. We did not select the representations based on their suitability for tasks since each task may require a different semantic representation to achieve the best performance, which would lead to using a different set of representations for each task. But, one of our primary considerations was not to hand-pick the program features and instead let the model decide based on the input code snippets. Moreover, using semantics from multiple dimensions, such as control flow, data dependency, and control dependencies, can still give good results [37]. We now define each of these representations below:

Definition 1 (Abstract Syntax Tree). It is an ordered tree of source code tokens where the non-terminal nodes represent the statements or operators, and the terminal nodes represent the operands or variables. AST represents the syntax of a program. It is abstract because it does not represent a programming language's actual syntax, instead represents its structure.

Definition 2 (Control Flow Graph). It is a directed graph where the nodes represent predicates and statements, and the edges connecting them indicate the transfer of control between those nodes. It describes the order in which program statements execute. Each edge also has a label that indicates the necessary conditions to execute that path.

Definition 3 (Program Dependence Graph). It is a directed graph where the nodes represent predicates and statements, and the edges connecting them indicate control dependencies and data dependencies between those nodes. PDG's edges are of two types: control dependency edges representing how the result of a predicate influences the variable and data dependency edges that indicate how one variable influences another. Ferrante et al. [36] first introduced this representation for compiler optimization.

```
void random_function()
{
    while(!flag)
    if(isValid(x))
    flag = true;
}
```

Figure 1: A sample C function / method

Figure 1 shows a sample C language function called *random_function*, and Figure 2 shows the AST, CFG, and PDG of the code snippet in Fig. 1. Each node in the figure has two fields separated by a comma (,). The first field represents the *type* of the node, the second field represents the actual *code token* associated with the node.

2.2. Path-based Code Representation

In their work, Alon et al. [38] proposed representing a code snippet as a set of paths in its AST. They used these AST paths with Conditional Random Fields (CRF) and evaluated the approach on method naming. Additionally, they introduced the concept of *Path Contexts* in their work. Later, *code2vec* [27] used these AST paths with an attention-based neural network and gained better results than the CRF approach. We briefly explain the concepts of AST path and path contexts as defined by Alon et al. [38]:



Figure 2: Extracting AST, CFG, and PDG paths from the code snippet in Fig. 1

Definition 4 (AST Path). Let *n* represent a node in an AST, and it has two attributes - the type of node *d* and the code token *t*. A path between nodes in the AST that begins at a terminal node n_1 , goes through a series of intermediate non-terminal nodes $n_2, ..., n_k$, and ends at another terminal n_{k+1} is called an AST path of length *k*. The path *p* is represented by a sequence of the form: $d_1a_1d_2a_2...d_ka_kd_{k+1}$, where $d_1, d_2, ..., d_{k+1}$ are the types of the nodes $n_1, n_2, ..., n_{k+1}$ respectively, and $a_1, a_2, ..., a_k$ are the directions of movements between path nodes in the AST. Here, $a_i \in \{\uparrow, \downarrow\}$.

Definition 5 (AST Path Context). It is a tuple $\langle t_1, p, t_{k+1} \rangle$, where t_1 and t_{k+1} are the tokens associated with the AST nodes n_1 and n_{k+1} respectively, i.e., terminal nodes of path p.

The intuition behind this approach is that each AST path captures a unique structural template for a set of statements. Consider the code snippet in Fig. 1. We show the AST for this snippet (with red edges) and an example AST path in Fig. 2. This sample AST path represents the statement *flag* = *true*. Further, any other code snippet with a similar assignment statement in its body (say a = b) has a subtree in its AST that is identical to the subtree that represents the statement *flag* = *true*, regardless of the identifiers used. Hence, both assignment statements can be represented by the same AST path. By this rationale, a set of AST paths can identify the structure of a code snippet. Further, the tokens (t_1 and t_{k+1}) associated with the two terminal nodes of a path are called the *context words* since they help differentiate the same path that occurs in two different contexts.

2.3. Attention Mechanism

Attention in deep learning is a technique that puts more emphasis/attention on the inputs that help the model make better predictions and less emphasis on the inputs that do not [27]. An attention layer does this by assigning numerical weights to its inputs to calculate an aggregation (weighted average) of all the inputs. Attention models are heavily used in NLP tasks like Neural Machine Translation [39] and Speech Recognition [40]. The main advantage of the attention models is that for two similar data instances (but not the same), they generate aggregations that are close in the vector space. The subtle differences in similar inputs are captured through the difference in weights assigned to them. Consider x_1, \ldots, x_n to be the inputs to the attention layer. The goal of the attention layer is to learn an attention vector *a*, which can be used to compute the weight α_i for a given input x_i . Given the attention vector, weights can be calculated as the normalized dot product of inputs and the attention vector:

$$\alpha_i = \frac{\exp(a^T \cdot x_i)}{\sum_{j=1}^n \exp(a^T \cdot x_j)} \tag{1}$$

$$\widetilde{v} = \sum_{i=1}^{n} \alpha_i \cdot x_i \tag{2}$$

A weighted average vector \tilde{v} can then be calculated to generate a single representation for the data instance (Eq 2). The attention vector *a* will be learned over time to generate better representations for all the instances in the dataset. This type of attention where the weights are calculated using dot product is called *Luong attention* [39].

3. Approach

3.1. Defining Semantic Path Contexts

Semantic features such as control flows and data dependencies have been used in some works to solve Software Engineering problems [29, 41]. However, using paths from semantic representations is yet to be explored. AST paths do not capture semantic aspects such as control flow and program dependencies which may be required to achieve optimal performance. To overcome this issue, in this paper, we extend the concept of paths to CFG and PDG:

Definition 6 (CFG Path). Let *n* represent a node in a CFG, and it has two attributes - the type of node *d* and the code token *t*. A path between nodes in the CFG that begins at the START node n_1 , goes through a series of intermediate nodes n_2, \ldots, n_k , and ends at a node n_{k+1} is called a CFG path of length *k*. The last node n_{k+1} can be of two types:

- The END node,
- A previously visited intermediate node that represents a loop control structure (i.e. $n_{k+1} \in n_2, \ldots, n_k$.)

The path *p* is represented by a sequence of the form: $d_1 \downarrow d_2 \downarrow \dots d_k a_k d_{k+1}$, where d_1, d_2, \dots, d_{k+1} are the types of the nodes n_1, n_2, \dots, n_{k+1} , and the direction a_k depends upon the last node n_{k+1} . Here, a_k is \downarrow , if n_{k+1} is END node, \uparrow otherwise.

Definition 7 (CFG Path Context). It is a tuple $\langle t_1, p, t_{k+1} \rangle$, where t_1 and t_{k+1} are the tokens associated with the CFG nodes n_1 and n_{k+1} respectively, i.e., terminal nodes of path p.

The START and END nodes are special nodes that represent the start and end of program execution. Each of the CFG paths represents a possible control flow pattern during program execution. To represent loop structures in a CFG, we extract three different paths from it:

- 1. A path that ignores the loop and proceeds to the next node (i.e. $n_{k+1} = \text{END}$),
- 2. A path that goes through the loop only once and proceeds to the next node (i.e. $n_{k+1} = \text{END}$),
- 3. A path that goes through the loop only once and ends at the loop's start node (i.e. $n_{k+1} \in n_2, \ldots, n_k$.)

These three paths together represent possible executions of a loop. Hence, for this reason, we have two types of last node n_{k+1} in our definition of the CFG path. Any other constructs like conditionals do not cause loops in CFG and hence are straightforward to handle.

Definition 8 (PDG Path). Let *n* represent a node in a PDG, and it has two attributes - the type of node *d* and the code token *t*. An edge *e* in a PDG has a label *l* associated with it. A PDG path *p* is a sequence of nodes $n_1, n_2, ..., n_{k+1}$ where all of the edges along the path have the same label l_p . The path *p* is represented by a sequence of the form: $d_1a_1d_2a_2...d_ka_kd_{k+1}$, where $d_1, d_2, ..., d_{k+1}$ are the types of the nodes $n_1, n_2, ..., n_{k+1}$, and $a_1, a_2, ..., a_k$ are the directions of movements between path nodes in the PDG. Here, $a_i \in \{\uparrow, \downarrow\}$ and $l_p \in \{\text{CDG}, \text{DDG}\}$.

Definition 9 (PDG Path Context). It is a tuple $\langle t_1, p, t_{k+1} \rangle$, where t_1 and t_{k+1} are the tokens associated with the PDG nodes n_1 and n_{k+1} respectively, i.e., terminal nodes of path p.

Since PDG is a combination of Control Dependency Graph (CDG) and Data Dependency Graph (DDG), each PDG path can either represent a control dependence of statements or a data dependence. The labels of edges in a path decide the type of PDG path. Since PDG is a graph (and not a tree), moving up and down does not make sense in PDG. However, we chose to keep all the directions of movements as \downarrow . We do this to have a consistent model input format. The concept of CFG and PDG path contexts are similar to AST path contexts. Fig. 2 depicts the CFG and PDG for the code in Fig. 1 and an example path for each representation.



Figure 3: Our pipeline to utilize a combination of Source Code Representations. The specific settings for the final prediction layer are chosen based on the task at hand.

3.2. Extracting Path Contexts

We have developed and used a python-based tool to extract different path contexts from C programs. We first parse each code snippet in the dataset using a platform called Joern [25] to generate AST, CFG, and PDG. These graphs are exported as .dot files. We then extract paths from these graphs as per the definitions provided earlier. If a snippet has multiple methods, Joern generates all the graphs (AST/CFG/PDG) for each method separately. In such cases, we extract paths from all such graphs and use them to represent the complete code snippet. To account for the high variation in lengths of AST paths, we follow *code2vec's* policy and extract only those with a maximum length of 8 and width of 2. An AST path's width refers to the difference in leaf node indices when all the leaf nodes are indexed sequentially. Further, to avoid high variation in the number of path contexts across code snippets, we limit the number of AST, CFG, and PDG path contexts extracted from a snippet to 200, 10, and 100, respectively. If a graph has more paths than the maximum limit, we randomly sample the maximum number of paths allowed for that representation. Some of the files may have hundreds of paths, and others may have as low as 1, and this variation could create very sparse matrices while training the neural network. So we use these settings to avoid sparse matrices as they could adversely affect the performance. We have set the path limits for AST (200), CFG (10) and PDG (100) based on their average path counts in the dataset (refer Tables 2 and 4).

3.3. Parallel Attention Pipeline

Fig. 3 depicts different phases in our approach, and specifically, the third phase shows our parallel attention pipeline. To combine AST, CFG, and PDG paths, we extend the *code2vec* model [27], which takes only AST path contexts as input. Our model takes multiple types of path contexts as inputs in a parallel pipeline and finally generates a code vector. We use multiple *code2vec*'s parallelly instead of a single *code2vec* because AST, CFG, and PDG represent different aspects of a program and need to be learned separately, which would allow us to combine their learned vectors later. Consider a code snippet *C* presented to the model as a bag of path contexts extracted from it.

$$C = [\{a_1, \dots, a_{n_1}\}, \{c_1, \dots, c_{n_2}\}\{p_1, \dots, p_{n_3}\}],$$
(3)

Here a_i , c_j , and p_k are AST, CFG, and PDG path contexts respectively. Each of the AST path contexts is processed by the model as follows: The three components of a path context (two context words and a path) are passed through embedding layers to generate token and path embeddings of size D. These three embeddings are then concatenated and passed through a fully connected (Dense) layer to generate a *context vector* x of size D. The *tanh* activation function is used for this dense layer. The primary purpose of this dense layer is to learn to compress the concatenated vector (size 3D) to generate a context vector (size D). This layer learns to perform this compression in a way that gives more importance to a path when it appears with some tokens (context words) and less importance when it appears with others. This helps in differentiating the same path that appears in two different contexts. As a result of this dense layer, context vectors $x_1, x_2, \ldots, x_{n_1}$ are generated for the AST path-contexts $a_1, a_2, \ldots, a_{n_1}$ respectively. These context vectors are passed through an attention layer to compute a weighted average, as explained in Section 2.3. This process is done parallelly for path contexts of each type (AST, CFG, and PDG), and three weighted average vectors \tilde{a}, \tilde{c} , and \tilde{p} are generated. Then, average vectors are concatenated to generate the final code vector v. The code vector v is used to make predictions regarding the code snippet C using another prediction layer. The model is trained to minimize the prediction error. All the context, weighted average, and code vectors are trained and learned concurrently. The prediction layer and the error function are decided based on the task, and we mention the specific details about our experiments in Section 4. We have used Tensorflow's Keras API to build our parallel attention pipeline.

4. Use Cases

4.1. Method Naming

The main goal of method naming is to predict/suggest an accurate name, given the method body. The predicted name should accurately represent the semantics of the method. This is an important problem because having good method names makes the code more readable and maintainable, but poorly named methods can adversely affect the programmers' productivity [11]. Consider *v* the code vector generated by our model for the input code snippet *C* and y_1, y_2, \ldots, y_N are distinct method names found in the training dataset. Our aim now is to use *v* to predict which of the labels y_1, y_2, \ldots, y_N is the actual name for the method *C*. After our model produces *v*, we use a fully connected layer with softmax activation function to generate the prediction probability vector \hat{p} . Consider $W \in \mathbb{R}^{N \times d}$ as the weight matrix associated with this layer whose rows correspond to the labels y_1, y_2, \ldots, y_N . The probability vector \hat{p} is calculated as follows:

$$\hat{p}_i = \frac{\exp(W_i \cdot v)}{\sum_{i=1}^N \exp(W_j \cdot v)}$$
(4)

where \hat{p}_i is the probability of y_i being the method name, and W_i is the row *i* of *W*. We use the standard cross-entropy loss function for the training:

$$loss(t, \hat{p}) = -\sum_{i=1}^{N} t_i \log{(\hat{p}_i)},$$
(5)

where *t* is an *N*-dimensional one-hot encoded true label. Then the predicted method name \hat{y} will be the one with highest probability:

$$\hat{y} = \arg\max_{i} \, (\hat{p}_i) \tag{6}$$

4.1.1. Dataset Preparation

We opted for a C language dataset to assess our model's performance in the Method Naming task. Our intention is to initially investigate the influence of semantic representations on downstream tasks with an imperative language like C. Once this initial experiments are demonstrated, the work could be further extended to object-oriented languages, such as Java. We could not find any existing C language datasets for method naming, so we decided to collect our own.

As the first step of our dataset collection, we fetched a list of open-source C projects from GitHub and ranked them by popularity. Popularity was quantitatively assessed using a combination of stars and forks, following a standard zscore approach [9, 27]. Repositories containing less than 10,000 methods were excluded from our analysis, as higher method count generally indicates a project with a reasonably rich development history and complexity [42]. Including projects with more than 10K methods also allows us to measure how well the model can generalize when trained on a single repository. Conversely, repositories with over 300K methods were omitted to maintain tractability and prevent any single project from unduly skewing the dataset, thereby ensuring a balanced and manageable scope for analysis. Further, due to the prevalence of similar domains (e.g., operating systems, compilers, databases, and firmware, etc.) in our repository list, we excluded projects of a similar nature at random to retain a more diverse set. To ensure a valid and insightful comparison with *code2vec*, we construct a dataset of comparable size, consisting of roughly 700,000 methods. Ultimately, our curated dataset encompasses 16 open-source C projects, collectively containing 729,218 methods. A breakdown of the repositories considered, including their specific domain and the count of methods, is available in our GitHub repository².

We preprocess the dataset by dividing all the C source files into individual methods. We then replace any occurrences of the method's name in its body with a special token and store the method name separately. We do this to remove any additional help the model might get from path contexts that contain the method name as a context word. We then normalize all method names by converting them to lowercase and splitting them into subtokens. A special character like '|' separates the subtokens. The normalized name is the true label for the method naming task. Then we extract paths from all the methods as explained in Section 3.2. We also filter out invalid methods that do not have at least one path from each representation. In this case, we extract paths at the *method-level* since the task deals with each method individually. These sets of paths are used to train the model to generate *method-level embeddings*. We show the average path count statistics for the dataset in Table 2. In addition to the full dataset, we provide the statistics for 5 out of 16 projects since we use them individually for evaluation.

Dataset	Number of Methods	AST	CFG	PDG
Full C dataset	729,218	72.8	4.4	15.8
FFmpeg	15,790	131.0	5.7	32.5
SumatraPDF	16,356	90.6	5.6	20.8
KBEngine	21,949	78.5	5.5	20.1
QEMU	39,881	92.3	4.0	17.3
CatBoost	54,365	75.1	4.6	19.7

Table 2: Average Path Count (per method) for our Dataset

4.1.2. Training, Evaluation, and Results

To train and evaluate the model, we shuffle methods from all 16 projects and split them into 649,004 training, 54,691 test, and 25,523 validation methods. Also, we randomly select some projects from the dataset to train and evaluate the model by treating them as individual datasets. We do this to gauge how effectively the model could generalize within individual projects. We train and test the model on each dataset using different combinations of representations (i.e., AST, AST + CFG, AST + CFG + PDG.) We train our model to minimize the cross-entropy loss (Eq. 5.) We use the Adam optimization algorithm with a batch size of 1024 and a learning rate of 0.001. We use dropout regularization with a dropout value of 0.25 on the context vectors to avoid overfitting. Most of these settings are adopted from *code2vec* to make our model's performance comparison with *code2vec* as fair as possible.

The method naming task is a multi-class classification problem, and for this reason, we use the well-known F1 score as the metric [27, 9]. We briefly describe the Precision, Recall, and F1 score metrics.

Precision: It is the ratio of true positives to all positive predictions. Measures the accuracy of a model's predictions.

Recall: It is the ratio of true positives to all actual positives. Measures a model's ability to find all positive instances.

F1 Score: The harmonic mean of Precision and Recall. Provides a balanced performance measure for classification.

To measure our model's performance on method naming, we consider the quality of the predicted method name. We calculate the precision and recall over sub-words within the predicted method name. The intuition is that the quality of a predicted method name is primarily determined by the sub-words used to construct it. When a prediction has a high recall, we can infer that model can predict most of the sub-words of the true label (actual method name). When a prediction has high precision, we can say that most of the sub-words in the predicted label are also in the true label.

²https://github.com/NobleMathews/mocktail-blue-lagoon

	F1 scores				
Dataset	<i>code2vec</i> (AST) AST + CFG AST + P		AST + PDG	AST + CFG + PDG	
Full C dataset	47.5	50.5	51.3	52.7	
FFmpeg	38.5	45.8	46.1	47.3	
SumatraPDF	13.7	27.8	29.4	33.2	
KBEngine	22.9	37.3	38.9	41.0	
QEMU	26.5	34.2	34.5	37.3	
CatBoost	37.7	46.3	47.8	49.7	

Table 3: Results for Method Name Prediction Task

RQ1. How do combinations of representations perform on method naming compared to AST?

We use *code2vec* as the baseline for all our experiments. Since *code2vec* only uses AST paths, it as a baseline allows us to measure the performance gain achieved using semantic code representations. We did not use any advances after *code2vec* as a baseline since our work is fundamentally an extension of it, and our main goal is to measure the impact of semantic representations but **not** to design a better learning model (which most of the advances after *code2vec* try to achieve). We summarize our results for the method naming task in Table 3. We can see that by adding additional representations like CFG and PDG, the F1 score increased in all of the cases. For this task, the performance boost given by PDG paths is more than the boost given by the CFG paths. Overall, CFG and PDG increased the F1 score from **47.5 to 52.7 (11% increase)** on the full dataset. Furthermore, we observe that the model can capture program properties very well within a project. For example, for the *SumatraPDF* project, the performance boost is more than 100%. The performance gain for individual projects is more significant than for the entire dataset. Thus, the CFG and PDG paths help the model perform well within a project rather than on the full dataset. This behaviour is expected as the combined dataset has functions from different projects with diverse programming styles and coding patterns.

Moreover, as Table 2 shows, the average AST, CFG, and PDG paths per method in the full dataset is 72.8, 4.4, and 15.8. There are many AST paths per method, even though we limit their number based on the length and width during the path extraction phase. In contrast, though CFG and PDG paths are not limited based on length, the average number of CFG and PDG paths are only 4.4 and 15.8, respectively. This introduces a huge problem of data sparsity on the CFG and PDG pipelines. However, for individual projects, the average number of CFG and PDG paths is much higher, and hence the performance gain is also higher (e.g., more than 100% increase in F1 for SumatraPDF project compared to 11% increase for the whole dataset.)

4.2. Program Classification

Program classification is a task that is aimed to classify the given code snippet into one of the many classes based on its functionality. It is an important task that is primarily helpful in maintaining huge collections of software [4, 43]. Consider v the code vector generated by our model for the input code snippet C and N be the total number of classes. The true label t is an N-dimensional one-hot encoded vector. We use a fully connected layer with the softmax activation function as the prediction layer. It takes the code vector v as input to generate the prediction probability vector \hat{p} . We use the standard cross-entropy loss function (Eq. 5) for the training. Then the predicted class label \hat{y} can be calculated as the dimension with the highest value (Eq. 6.)

4.2.1. Dataset Preparation

We use the Open Judge (OJ) dataset first introduced by Mou et al. [4] and has been used in several other works on Program Classification [14, 17, 18]. The dataset is a collection of solutions to 104 different programming problems submitted to an online open judge (OJ). Each problem has 500 different C language solutions, making the dataset a

collection of $104 \ge 52,000$ files divided into 104 classes. We aim to classify the programs that solve the same problem into the same class.

We prepare the dataset by extracting paths from each source file (Section 3.2) and including only those files with at least one path from each representation. Unlike method naming, this task operates with individual files rather than methods. Hence, the set of paths extracted from each program is a *file-level* representation. We can use these file-level representations to train our parallel attention pipeline to generate *file-level embeddings*. We show the average path count statistics for the OJ dataset in Table 4. One can observe that the average number of paths of each type is significantly higher than in the method naming dataset since each file now can have multiple methods.

Table 4: Average	Path Count	(per file)) for the OJ	Datase
		MI		

Number of Files	AST	CFG	PDG
52,000	175.7	29.3	74.6

4.2.2. Training, Evaluation, and Results

Before training the model, we create Train-Test-Validation splits for the OJ dataset. We split each class (of 500 files) into a 70:20:10 ratio, creating a dataset with 36,400 train samples, 10,400 test samples, and 5,200 validation samples. We train the network using the Adam optimization algorithm with a batch size of 1024 and a learning rate of 0.001, the same model settings and configurations as method naming. In this context, a fully connected layer with the softmax activation function acts as the prediction layer for program classification. We use the classification accuracy on the test dataset as the metric to evaluate the performance as done in other related works [4, 14, 17]. It is calculated as the percentage of classifications done correctly.

RQ2. How do combinations of representations perform on program classification compared to AST?

Table 5 compares the performance of our approach to the baseline AST. The results follow a similar pattern to that of method naming. We can see that as we include CFG and PDG, the accuracy increases with each addition. The PDG paths significantly contribute to the performance gain than the CFG paths, which may be attributed to the much higher number of PDG paths than the CFG paths (Table 4). When all three representations are included, the model shows a performance increase of **15.7**%.

Approach	Test Accuracy (%)
AST (code2vec)	73.65
AST + CFG	75.79
AST + PDG	84.88
AST + CFG + PDG	85.23

Table 5: Results of Program Classification Task

4.3. Code Clone Detection

Code clone detection task aims to detect whether a given pair of code snippets are similar. This paper deals with a specific kind of clone pairs called *functional clone pairs* or **Type-4** clones [44]. We are not dealing with other types of clone pairs that might be lexically or syntactically similar but lack functional similarity. Type-4 clones typically pose the most intricate detection challenges due to their inherent complexity and are usually the hardest to detect [44]. For instance, a pair of code snippets are considered Type-4 clones, irrespective of differences in code tokens and structural composition, as long as they have the same functionality (e.g., the same solution implemented using

different algorithms). Owing to our current goal of demonstrating the combination of source code representations for downstream tasks, the detection of Type-1 through Type-3 clones is not within the current scope.

Consider v_1 and v_2 the code vectors generated by our model for two input code snippets. If the input pair is a true clone pair (i.e., they solve the same problem), its ground truth label (*t*) will be 1, and a false clone pair will have ground truth label -1. Using this convention, we can calculate their functional similarity using the cosine similarity measure:

$$sim(v_1, v_2) = \frac{v_1 \cdot v_2}{|v_1| \cdot |v_2|} \tag{7}$$

The cosine similarity score determines how similar two vectors are in an N-dimensional space, and we later use it to determine a clone pair. To determine if two code snippets constitute a clone pair, we require the vectors of both snippets simultaneously. Consequently, our model cannot be directly used for supervised training aimed at minimizing the error in detecting clone pairs, as it requires modifying the model with additional layers. Instead, we adopted an *unsupervised* code clone detection approach [6, 45, 18]. It is a commonly used approach where the main idea is to pre-train the model on a different task with the same dataset and use the learned model weights to generate and save the code vectors for all code snippets in the dataset. Then, we form code vector pairs to calculate cosine similarity and determine whether the input code snippets are clones based on whether the similarity score exceeds a threshold (say θ):

$$\hat{y} = \begin{cases} 1, & sim(v_1, v_2) > \theta \\ 0, & sim(v_1, v_2) \le \theta \end{cases}$$

$$\tag{8}$$

This paper uses program classification as the pretraining task. We chose program classification as the pretraining task because the primary goal in detecting functional clones and classifying programs based on functionality is the same: to capture the functionality of the code snippets.

4.3.1. Dataset Preparation

We use the same OJ dataset for the clone detection task as well. The dataset is a collection of solutions to 104 different programming problems submitted to an online open judge (OJ). Each solution is stored as a separate file in the dataset. We prepare the dataset as explained in Section 4.2.1 and use it to pretrain the model on the program classification task. However, we use a modified version of the OJ dataset to evaluate the unsupervised code clone detector (Eq 7 and 8.) Any two solutions/files from the same category of 500 solutions form a clone pair (i.e., its ground truth label is 1). This is based on the fact that all files belonging a category solve the same problem, and thus, have the same functionality. Conversely, any two files that belong to two different categories solve different problems and *do not* form a clone pair (i.e., its ground truth label is 0). Since we have 104 categories with 500 solutions each, which generates 28M code pairs, which is still hard to process. We then randomly sample 50K true clone pairs and 50K false clone pairs to create a final dataset for code clone detection. This modified version of the OJ dataset is called *OJClone*, first introduced by Wei and Li [26] and later adopted by other works [14, 18]. Similar to program classification, this task operates at the *file-level*, aiming to evaluate functional similarity between files in the dataset rather than individual functions.

4.3.2. Training, Evaluation, and Results

Because we are dealing with unsupervised code clone detection, we first train the model on program classification exactly as in Section 4.2.2 and then save the vectors for all code snippets. Exporting the code vectors makes it easy to create code pairs and to reproduce the results. As explained before, we form 50K clone and 50K non-clone code vector pairs and use cosine similarity to evaluate the model (Eq. 7 and 8.) We chose the threshold (θ) as **0.4** for our clone detector. We chose this value as it allowed our clone detector to achieve optimal performance in all four cases. Since the problem is formulated as a binary classification problem (clone pair or not), we evaluate our approach using Precision, Recall, and F1-score. To calculate these metrics, we treat true clone pairs as positive and false clone pairs as negative samples.

RQ3. How do combinations of representations perform on code clone detection compared to AST?



Figure 4: Similarity score distribution for the OJClone dataset predicted by four different approaches.

Approach	Precision	Recall	F1
AST (code2vec)	0.93	0.81	0.86
AST + CFG	0.96	0.87	0.91
AST + PDG	0.95	0.86	0.90
AST + CFG + PDG	0.96	0.92	0.94

Table 6: Results of Code Clone Detection Task

Fig. 4 depicts the cosine similarity distribution for all clone and non-clone pairs. We can see that CFG and PDG cause the similarity distribution of positive samples to shift towards 1, while the similarity distribution of negative samples shifts in the opposite direction. A threshold value of **0.4** produced the best results in all three cases. We compare the performance of our approach with AST in Table 6. Once again, the results show that including semantic code representations improves the performance significantly. Though the improvement in precision is only 0.03, this is expected since AST alone achieves an impressive precision of 0.93. However, our approach's main enhancement is the recall, which has increased from 0.81 to 0.92. Overall, including CFG and PDG improves the F1 by **9%**. An interesting observation is unlike the previous two tasks, CFG provides a slightly better performance than PDG. Even though the PDG paths are significantly more in number than CFG paths (refer Table 4), PDG did not provide superior performance gain. This outcome can mean that control flows can capture the behaviour of programs much better than data dependencies. This is understandable since PDG is more inclined towards capturing how a variable depends on another but not how a variable behaves based on other variables.

5. Discussion

When evaluating how good an approach is to solve a problem, there can be multiple factors to consider, of which we have already discussed an important one: *performance*. Another crucial factor would be the *additional effort* needed compared to existing approaches. As previously discussed, a path-based approach does not manually craft program features but lets the model pick appropriate features for a task. Moreover, the path extraction process can be extended to other languages with minor adjustments, such as replacing the parser and modeling interactions between any language-specific constructs as paths. Since the manual effort is significantly reduced and offloaded to the model, measuring the additional processing overheads incurred in the data preparation and training phases is essential. We analyze the processing overheads in this section and examine our results for further insights. While we acknowledge that there could be other ways of measuring the impact of using multiple source code representations for Software

Engineering tasks, we limit the scope of this paper to *performance* (using metrics such as Accuracy, F1, etc.) and *additional overhead incurred* (in terms of computational time overhead).

RQ4. What are the additional processing overheads incurred by including multiple representations?

Table 7 compares the throughputs of four different approaches - AST, AST + CFG, AST + PDG, AST + CFG + PDG. We compare the throughputs at three stages of our pipeline: Path Extraction (Dataset creation), Model Training, and Model Inference. We use a system with two Intel(R) Xeon(R) CPU E5-2640 v4 chips during the path extraction phase. This system effectively has 40 CPUs, which our path extractor utilizes concurrently to extract paths from multiple source files. We have used a system with three *GeForce GTX 1080 Ti* GPUs for training and evaluating the model.

Table 7 shows that all three throughputs decrease as we include additional representations. The path extraction throughputs for the method naming task are much higher than for program classification since the former is a method-level task, and thus, samples have significantly fewer paths per representation than the latter. The path extraction throughput decreased by 35.2% and 34.4% for both the tasks when CFG and PDG were included. The dataset creation times can be hugely improved by making the path extraction algorithm more efficient and traversing the graphs in-memory (without exporting them to a .dot file).

	Path Extraction Throughput		Model Throughput (samples/second)			
	(samples/minute)		During Training		During Inference	
Representations used	Method Naming	Program Classification	Method Naming	Program Classification	Method Naming	Program Classification
AST	2272	29	268	846	858	4110
AST + CFG	1827	23	179	535	694	3827
AST + PDG	1538	22	164	461	619	2854
AST + CFG + PDG	1470	19	126	317	563	2722

Table 7: Comparing the time taken by each approach in different phases of the pipeline

The model throughput during the training phase is quite affected when all three representations are used. The training throughput decreased by 53% and 62.5% for method naming and program classification tasks. The throughput for program classification is more affected because the average number of CFG and PDG paths at the file level is much higher than at the method level (refer to Tables 2 and 4), and hence it takes more time for the model to process a code snippet for program classification. The program classification task's training throughputs are much greater than the method naming task. One might think this is a contradiction since a file-level task has more paths per sample than a method-level task, but other factors also affect model throughput, for example, model size and dataset size. The dataset for method naming is much larger than program classification, which increases its model size and thus decreases the training throughput.

Though the model's inference throughput was also affected when we used multiple representations, the throughput is still acceptable. To get a much clearer picture, we should look at the total time taken during inference, including time taken for path extraction and model inference time. For method naming, the average total inference time per sample comes out to be 0.027 seconds using AST and 0.043 seconds using all three representations (59% increase). Similarly, average inference times for program classification are 2.07 and 3.16 seconds, respectively (52.6% increase). Most of the increase in time comes from the path extraction phase. As mentioned, one can decrease path extraction times using better algorithms, avoiding expensive file operations, etc.

RQ5. What are the implications of our results? and how does each representation contribute to the model's performance?

Our experiments show that each semantic representation, when used with AST, only increases the model's performance but never deteriorates it. This observation indicates that AST does not provide a complete picture, and we might need additional representations to capture program features (semantics) effectively. However, some representations may be much effective than others for a given task and dataset. Consider Table 8, where we compare the performance gain provided by CFG and PDG for all three tasks. We measure performance gain as the increase in the performance metric when CFG/PDG is included. We can observe that the combination AST + PDG is more effective than AST + CFG for method naming and classification tasks. This may be due to the higher number of paths in PDG than CFG, which leads to PDG providing more features for the model to generalize. Although the combination AST + CFG gives a slightly better performance for clone detection than AST + PDG, the results are close. This variation suggests that selecting the most suitable combination of code representations is largely task-dependent.

Table 8: Comparison of performance gain provided by CFG and PDG

Task	AST	AST + CFG	AST + PDG
Method Naming (F1)	47.5	50.5	51.3
Classification (Accuracy)	73.65	75.79	84.88
Clone Detection (F1)	0.86	0.91	0.90

Furthermore, our results indicate that using all three representations may not always yield a significant performance boost compared to a more focused subset. For example, in the case of the program classification task, AST + PDG proves to be an optimal combination, delivering nearly equivalent performance to AST + CFG + PDG, but with higher training and inference throughputs of 461 and 2854 samples/second, compared to 317 and 2722 samples/second for AST+CFG+PDG. This suggests that customizing code representations for specific software engineering tasks has the potential to reduce processing times while maintaining good performance. While we have observed promising results within the context of our study, further research on a broader range of code bases and tasks is essential to validate these observations and establish generalized conclusions.

6. Threats to Validity

6.1. External Validity

Language Generalization: Our study primarily focuses on the C programming language, and the conclusions drawn may not directly generalize to other languages. Hence, there is a need to adapt our approach for each language, taking into account their distinct characteristics.

Dataset Bias: The diversity of the dataset used for training and evaluation can significantly impact the results. While we tried to diversify the dataset with projects from different domains, it still may not fully represent the broader landscape of software projects. Further exploration using larger datasets, broader in scope, and encompassing languages beyond C is crucial to draw generalizable conclusions.

6.2. Internal Validity

Data Sparsity: Though our approach significantly boosts the performance for all three tasks, the data sparsity problem discussed in Section 4.1.2 can potentially impact the model negatively when used on a larger scale. One possible solution to address this sparsity is to devise more unique ways to formulate and extract semantic features from these representations, such as paths that capture interactions between files and modules. We can also use additional semantic features like *data flow* in addition to control flows (CFG) and program dependencies (PDG). One downside of this approach is the increased effort for pre-processing and training, which can be decided as a trade-off for a respective downstream task. We are essentially extracting more information from a program in the form of paths and training the model rather than extracting limited information and requiring more training data to achieve the same performance.

Parameter Settings: To ensure a fair comparision with code2vec, we have adopted code2vec's network hyperparameters, including learning rates, batch sizes, and dropout rates. Extensive hyperparameter tuning is essential to achieve optimal results with the aggregate attention model. Additionally, limiting maximum path lengths and widths to 8 and 2 (refer 3.2) may potentially exclude relevant information and affect the quality of extracted features. Similarly, using threshold values on the number of paths extracted from a code snippet (200 AST paths, 10 CFG paths, and 100 PDG paths) derived from average path counts within the dataset may not universally accommodate diverse code structures, and can potentially impact representation quality. Furthermore, while avoiding sparse matrices in neural network training, the strict enforcement of path count limits may inadvertently lead to information loss in code snippets with inherently larger path counts. A detailed sensitivity analysis could be performed to optimise and fine-tune parameters such as path lengths, widths and maximum path counts. Such analysis would involve assessing the impact of these parameter choices on the performance of the Aggregate attention model on downstream tasks. However, in the current work we limit ourselves towards exploring the impact of using multiple code representations in software engineering downstream tasks.

6.3. Construct Validity

Automation Trade-Off: Another drawback to our approach is that since we are not manually designing program features to represent code, some unusual cases like *inline assembly* may not be appropriately handled while extracting paths. Accommodating such unique scenarios through automation is a trade-off that requires continued exploration and refinement.

7. Related Work

7.1. Source Code Representation

Traditionally, researchers have expressed source code as a sequence of tokens to address important software engineering tasks [46, 47, 45]. *SourcererCC* [45] creates a partial index for source code by using code tokens and then uses it to detect code clones. For the task of bug localization, Zhou et al. [47] treat source code files as a text corpus to find the similarity between each file and the bug report.

Representing structural information of source code using Abstract Syntax Trees (AST) has emerged as a critical approach, capturing both lexical and syntactic information. ASTs have been used extensively in the literature [7, 6, 14, 16, 18, 19, 20]. For example, Deckard [6] introduces an algorithm for identifying similar sub-trees of two ASTs to detect code clones. Researchers have also tried to capture syntactic information using deep learning models like RNN [7] or Tree-LSTM [26] on ASTs to detect code clones. Mou et al. [4] use a custom Tree-based Convolutional Neural Network (TBCNN) on ASTs to learn vector representations of code snippets. In their work, Zhang et al. [14] extract the sub-trees from AST and feed them to an AST-based neural network (ASTNN) to generate code vectors that can capture the sequential dependency of code statements. Li et al. [23] utilize the global program dependencies of source code along with the local ASTs to predict bugs. Alon et al. [38, 27] introduce the AST path-based approach for representing source code and various learning models to generate code vectors using AST paths. Code2seq [28] adopts a similar strategy to *code2vec* for the task of Neural Machine Translation. More recently, Wang et al. [16] introduced heterogeneous program graphs by including additional type information for nodes and edges in an AST and used GNNs to learn program properties. In another work, Wang et al. [17] use a modular tree-based network to detect the semantic difference in programs based on their ASTs. InferCode [18] use the subtrees of an AST as the training labels and train a TBCNN in a self-supervised way. This way, the generated code vectors are not tied to a specific task. Xiao et al. [48] used a new notion of *path context* and introduced the path context augmented network (PCAN) to learn code vectors.

Researchers also explored the usage of *Data Flow Graphs* to capture source code's structure [49, 50]. *Graph-CodeBERT* [49] introduced a transformer-based model that uses the data flow in programs to learn the representations. Instead of taking the syntactic-level code structure like an AST, these approaches use the data flow to capture the inherent code structure.

Some of the works have combined different graph structures; for example, Allamani et al. [29] represent the program as a directed graph of code tokens with different labeled edges like syntax tree, control flow, and data flow for predicting variable and method names. The Code Property Graph (CPG) by Yamaguchi et al. [25] is the most relevant work to our approach where they create a *static* combination of AST, CFG, and PDG for the task of vulnerability detection. They show its effectiveness by finding 18 previously undiscovered vulnerabilities in the Linux kernel's

source code. Zhang et al. [51] constructed a code knowledge graph and used a bi-attention layer neural network to detect bugs. More recently, Long et al. [30] introduced a multiview graph using data-flow, control-flow, read-write graphs to obtain multiple perspectives about source code and employed a GGNN to extract information from them.

Many works have also built upon the *code2vec* model for various downstream tasks [32, 33, 34]. Compton et al. [32] extend the *code2vec* model to Java classes by aggregating different method embeddings found in a class. Shi et al. [33] use the *code2vec* model on pairs of AST paths as input for the task of defect prediction. They show an increase in performance over the state-of-the-art model by 17%. In another work, Shi et al. [34] use the *code2vec* model for discovering misconceptions in computing assignments. Our results show that the works built upon *code2vec* can be enhanced by including semantic paths.

7.2. Method Naming

Allamanis et al. [11] is the first work that explicitly proposes a solution to the method naming problem to the best of our knowledge. They use a log-bilinear neural network to map the method names to a high-dimensional continuous space such that semantically similar names are closer to each other in the space. As the method's name usually indicates its semantics, this problem is a special case of code summarization. With this intuition, some of the works on code summarization use method naming to evaluate the approach. For example, Allamanis et al. [9] introduce a convolutional attention neural network that takes code tokens as input to detect program features in a context-dependent way and produce a method name for that code snippet. Alon et al. [38] use AST paths as inputs to Conditional Random Fields (CRF) to predict a method name for the input code snippet. They improved the performance in their subsequent work, *code2vec* [27], using AST paths as input to an attention-based neural network. Later, *code2seq* [28] followed a similar path-based approach with a different model and achieved a significant performance boost compared to previous works. Recently, Wang et al. [16] used GNNs on their custom program graphs to predict method names.

7.3. Program Classification

Several works have tried to classify source code based on the programming language used [52, 53], authorship [54], domain [43], or as in our case, its functionality [4, 14, 17]. Mou et al. [4] are one of the first works that use a learning model (TBCNN) to classify programs based on functionality. Later, Zhang et al. [14] proposed ASTNN to overcome the problem of vanishing gradients from which previous deep learning techniques suffered. Wang et al. [17] proposed a model (MTN) with multiple neural modules to deal with different AST semantic units. Their approach showed an accuracy improvement of 1.8% over TBCNN. More recently, *InferCode* [18] used *self-supervised* learning to improve the performance of TBCNN by 4%.

7.4. Code Clone Detection

Researchers have been actively studying code clone detection due to its applications in software engineering [44]. Many works have suggested a wide range of approaches, from token-based [46, 45] or tree-based techniques [24, 6] to supervised [26, 55, 14, 56] and unsupervised [7, 18] deep learning methods. Baxter et al. [24] detect code clones by generating ASTs for input code snippets and comparing their subtrees. This approach is one of the first attempts to solve this problem without using string matching. *CCFinder* [46] creates a regularized token sequence from the programs to detect duplicate code. White et al.'s [7] deep learning based technique is one of the earlier attempts to automatically learn program features using a Recursive Neural Network on ASTs to detect code clones. Wei and Li [26] propose a deep learning framework that detects clones by learning features using an AST-based LSTM. Several supervised learning techniques like Oreo [55], ASTNN [14], MTN [17], and Fang et al.'s *fusion learning* [56] have also shown significant improvements over traditional techniques.

8. Conclusion and Future Work

Moving away from the predominantly common approach of using AST for downstream tasks, in this work, we explored the idea of integrating AST with semantic code representations (CFG / PDG) to measure the impact on Software Engineering tasks. Towards this goal, we extended an AST path-based technique and adapted it to include CFG and PDG path contexts. This allowed us to measure the impact of using multiple code representations. We

evaluated our approach on the method naming task with different sets of data, first with the full C dataset of 16 projects and then with some individual projects. By including CFG and PDG path contexts, we demonstrate that the model outperforms *code2vec* by 11% on the full dataset and up to 100% on individual projects. The performance boost observed for individual projects is much more significant than for the entire dataset, potentially owing to the higher variation in the number of AST, CFG, and PDG paths. To test whether this approach has the potential to be generalized for multiple software engineering tasks, we also evaluated the approach on program classification and code clone detection. The combination AST+CFG+PDG outperforms *code2vec* in these two tasks by 15.7% (Accuracy) and 9.3% (F1), respectively. We also measured the impact of multiple representations by measuring the additional computational overhead incurred.

While our initial findings are promising, we see immense scope for further analysis towards generalization. One interesting direction is to extend the approach to more programming languages. Specifically, the approach could be extended to object-oriented languages such as Java to capture the control flow and dependencies between objects and classes. Analyzing the impact of different language characteristics on the effectiveness of semantic representations and model performance can provide valuable insights. To make our approach more robust and widely applicable, it is important to work with more diverse and extensive datasets. This can potentially provide a better understanding of the impact of semantic representations in different software development scenarios, taking into account various coding styles and practices. Another research direction is to explore other semantic representations such as *Data Flow Graph*.

Though our primary intent with this work is not to propose an efficient learning model, our approach can still be explored for other tasks, such as *bug localization* and *code generation*. The proposed approach itself can be improved by conducting comprehensive hyperparameter optimization and sensitivity analyses to fine-tune model configurations. This can help achieve the optimal performance with our approach. The proposed approach can also be leveraged to improve the performance of several works built upon *code2vec*. Moreover, the way these representations are utilized can be modified to suit the task at hand. Also, the proposed path-based approach itself could be further investigated to find the *optimal combinations of code representations* for various downstream tasks. Another direction of future work in our approach is to solve the data sparsity problem during training, as discussed in Section 5.

We anticipate that this study can motivate researchers to replicate existing approaches by integrating syntactic and semantic representations. In addition, we see that the proposed approach can lay the groundwork to spin off multiple novel code representations for various sub-domains of Software Engineering while efficiently leveraging advances in AI and Programming Language Processing.

References

- S. Bajracharya, J. Ossher, C. Lopes, Sourcerer: An infrastructure for large-scale collection and analysis of open-source code, Science of Computer Programming 79 (2014) 241–259.
- [2] D. Spadini, M. Aniche, A. Bacchelli, Pydriller: Python framework for mining software repositories, in: Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, 2018, pp. 908– 911.
- [3] S. M. Reza, O. Badreddin, K. Rahad, Modelmine: a tool to facilitate mining models from open source repositories, in: Proceedings of the 23rd ACM/IEEE International Conference on Model Driven Engineering Languages and Systems: Companion Proceedings, 2020, pp. 1–5.
- [4] L. Mou, G. Li, L. Zhang, T. Wang, Z. Jin, Convolutional neural networks over tree structures for programming language processing, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 30, 2016.
- [5] F. Barchi, E. Parisi, G. Urgese, E. Ficarra, A. Acquaviva, Exploration of convolutional neural network models for source code classification, Engineering Applications of Artificial Intelligence 97 (2021) 104075.
- [6] L. Jiang, G. Misherghi, Z. Su, S. Glondu, Deckard: Scalable and accurate tree-based detection of code clones, in: 29th International Conference on Software Engineering (ICSE'07), IEEE, 2007, pp. 96–105.
- [7] M. White, M. Tufano, C. Vendome, D. Poshyvanyk, Deep learning code fragments for code clone detection, in: 2016 31st IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2016, pp. 87–98.
- [8] G. Shobha, A. Rana, V. Kansal, S. Tanwar, Code clone detection—a systematic review, Emerging Technologies in Data Mining and Information Security (2021) 645–655.
- [9] M. Allamanis, H. Peng, C. Sutton, A convolutional attention network for extreme summarization of source code, in: International conference on machine learning, PMLR, 2016, pp. 2091–2100.
- [10] S. Liu, Y. Chen, X. Xie, J. Siow, Y. Liu, Retrieval-augmented generation for code summarization via hybrid gnn, in: International Conference on Learning Representations, 2021.
- [11] M. Allamanis, E. T. Barr, C. Bird, C. Sutton, Suggesting accurate method and class names, in: Proceedings of the 2015 10th joint meeting on foundations of software engineering, 2015, pp. 38–49.
- [12] X. Gu, H. Zhang, S. Kim, Deep code search, in: 2018 IEEE/ACM 40th International Conference on Software Engineering (ICSE), IEEE, 2018, pp. 933–944.

- [13] X. Ling, L. Wu, S. Wang, G. Pan, T. Ma, F. Xu, A. X. Liu, C. Wu, S. Ji, Deep graph matching and searching for semantic code retrieval, ACM Transactions on Knowledge Discovery from Data (TKDD) 15 (5) (2021) 1–21.
- [14] J. Zhang, X. Wang, H. Zhang, H. Sun, K. Wang, X. Liu, A novel neural source code representation based on abstract syntax tree, in: 2019 IEEE/ACM 41st International Conference on Software Engineering (ICSE), IEEE, 2019, pp. 783–794.
- [15] M. Allamanis, E. T. Barr, P. Devanbu, C. Sutton, A survey of machine learning for big code and naturalness, ACM Computing Surveys (CSUR) 51 (4) (2018) 1–37.
- [16] W. Wang, K. Zhang, G. Li, Z. Jin, Learning to represent programs with heterogeneous graphs, arXiv preprint arXiv:2012.04188 (2020).
- [17] W. Wang, G. Li, S. Shen, X. Xia, Z. Jin, Modular tree network for source code representation learning, ACM Transactions on Software Engineering and Methodology (TOSEM) 29 (4) (2020) 1–23.
- [18] N. D. Bui, Y. Yu, L. Jiang, Infercode: Self-supervised learning of code representations by predicting subtrees, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 1186–1197.
- [19] Y. Li, S. Wang, T. N. Nguyen, Fault localization with code coverage representation learning, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 661–673.
- [20] S. Kim, J. Zhao, Y. Tian, S. Chandra, Code prediction by feeding trees to transformers, in: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, 2021, pp. 150–162.
- [21] X. Sun, Y. Zhongyang, Z. Xin, B. Mao, L. Xie, Detecting code reuse in android applications using component-based control flow graph, in: IFIP international information security conference, Springer, 2014, pp. 142–155.
- [22] A. V. Phan, M. Le Nguyen, L. T. Bui, Convolutional neural networks over control flow graphs for software defect prediction, in: 2017 IEEE 29th International Conference on Tools with Artificial Intelligence (ICTAI), IEEE, 2017, pp. 45–52.
- [23] Y. Li, S. Wang, T. N. Nguyen, S. Van Nguyen, Improving bug detection via context-based code representation learning and attention-based neural networks, Proceedings of the ACM on Programming Languages 3 (OOPSLA) (2019) 1–30.
- [24] I. D. Baxter, A. Yahin, L. Moura, M. Sant'Anna, L. Bier, Clone detection using abstract syntax trees, in: Proceedings. International Conference on Software Maintenance (Cat. No. 98CB36272), IEEE, 1998, pp. 368–377.
- [25] F. Yamaguchi, N. Golde, D. Arp, K. Rieck, Modeling and discovering vulnerabilities with code property graphs, in: 2014 IEEE Symposium on Security and Privacy, IEEE, 2014, pp. 590–604.
- [26] H. Wei, M. Li, Supervised deep features for software functional clone detection by exploiting lexical and syntactical information in source code., in: IJCAI, 2017, pp. 3034–3040.
- [27] U. Alon, M. Zilberstein, O. Levy, E. Yahav, code2vec: Learning distributed representations of code, Proceedings of the ACM on Programming Languages 3 (POPL) (2019) 1–29.
- [28] U. Alon, S. Brody, O. Levy, E. Yahav, code2seq: Generating sequences from structured representations of code, arXiv preprint arXiv:1808.01400 (2018).
- [29] M. Allamanis, M. Brockschmidt, M. Khademi, Learning to represent programs with graphs, arXiv preprint arXiv:1711.00740 (2017).
- [30] T. Long, Y. Xie, X. Chen, W. Zhang, Q. Cao, Y. Yu, Multi-view graph representation for programming language processing: An investigation into algorithm detection, arXiv preprint arXiv:2202.12481 (2022).
- [31] H. Perl, S. Dechand, M. Smith, D. Arp, F. Yamaguchi, K. Rieck, S. Fahl, Y. Acar, Vccfinder: Finding potential vulnerabilities in open-source projects to assist code audits, in: Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 2015, pp. 426–437.
- [32] R. Compton, E. Frank, P. Patros, A. Koay, Embedding java classes with code2vec: Improvements from variable obfuscation, in: Proceedings of the 17th International Conference on Mining Software Repositories, 2020, pp. 243–253.
- [33] K. Shi, Y. Lu, J. Chang, Z. Wei, Pathpair2vec: An ast path pair-based code representation method for defect prediction, Journal of Computer Languages 59 (2020) 100979.
- [34] Y. Shi, K. Shah, W. Wang, S. Marwan, P. Penmetsa, T. Price, Toward semi-automatic misconception discovery using code embeddings, in: LAK21: 11th International Learning Analytics and Knowledge Conference, 2021, pp. 606–612.
- [35] D. Vagavolu, K. C. Swarna, S. Chimalakonda, A mocktail of source code representations, in: 2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE), IEEE, 2021, pp. 1296–1300.
- [36] J. Ferrante, K. J. Ottenstein, J. D. Warren, The program dependence graph and its use in optimization, ACM Transactions on Programming Languages and Systems (TOPLAS) 9 (3) (1987) 319–349.
- [37] J. K. Siow, S. Liu, X. Xie, G. Meng, Y. Liu, Learning program semantics with code representations: An empirical study, arXiv preprint arXiv:2203.11790 (2022).
- [38] U. Alon, M. Zilberstein, O. Levy, E. Yahav, A general path-based representation for predicting program properties, ACM SIGPLAN Notices 53 (4) (2018) 404–419.
- [39] M.-T. Luong, H. Pham, C. D. Manning, Effective approaches to attention-based neural machine translation, arXiv preprint arXiv:1508.04025 (2015).
- [40] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, Y. Bengio, End-to-end attention-based large vocabulary speech recognition, in: 2016 IEEE international conference on acoustics, speech and signal processing (ICASSP), IEEE, 2016, pp. 4945–4949.
- [41] M. Tufano, C. Watson, G. Bavota, M. Di Penta, M. White, D. Poshyvanyk, Deep learning similarities from different representations of source code, in: 2018 IEEE/ACM 15th International Conference on Mining Software Repositories (MSR), IEEE, 2018, pp. 542–553.
- [42] S. Omari, G. Martinez, Enabling empirical research: A corpus of large-scale python systems, in: Proceedings of the Future Technologies Conference (FTC) 2019: Volume 2, Springer, 2020, pp. 661–669.
- [43] M. Linares-Vásquez, C. McMillan, D. Poshyvanyk, M. Grechanik, On using machine learning to automatically classify software applications into domain categories, Empirical Software Engineering 19 (3) (2014) 582–618.
- [44] C. K. Roy, J. R. Cordy, A survey on software clone detection research, Queen's School of Computing TR 541 (115) (2007) 64-68.
- [45] H. Sajnani, V. Saini, J. Svajlenko, C. K. Roy, C. V. Lopes, Sourcerercc: Scaling code clone detection to big-code, in: Proceedings of the 38th International Conference on Software Engineering, 2016, pp. 1157–1168.
- [46] T. Kamiya, S. Kusumoto, K. Inoue, Ccfinder: a multilinguistic token-based code clone detection system for large scale source code, IEEE

Transactions on Software Engineering 28 (7) (2002) 654-670.

- [47] J. Zhou, H. Zhang, D. Lo, Where should the bugs be fixed? more accurate information retrieval-based bug localization based on bug reports, in: 2012 34th International Conference on Software Engineering (ICSE), IEEE, 2012, pp. 14–24.
- [48] D. Xiao, D. Hang, L. Ai, S. Li, H. Liang, Path context augmented statement and network for learning programs, Empirical Software Engineering 27 (2) (2022) 1–26.
- [49] D. Guo, S. Ren, S. Lu, Z. Feng, D. Tang, S. Liu, L. Zhou, N. Duan, A. Svyatkovskiy, S. Fu, et al., Graphcodebert: Pre-training code representations with data flow, arXiv preprint arXiv:2009.08366 (2020).
- [50] P. Vytovtov, K. Chuvilin, Unsupervised classifying of software source code using graph neural networks, in: 2019 24th Conference of Open Innovations Association (FRUCT), IEEE, 2019, pp. 518–524.
- [51] J. Zhang, R. Xie, W. Ye, Y. Zhang, S. Zhang, Exploiting code knowledge graph for bug localization via bi-directional attention, in: Proceedings of the 28th International Conference on Program Comprehension, 2020, pp. 219–229.
- [52] J. K. Van Dam, V. Zaytsev, Software language identification with natural language classifiers, in: 2016 IEEE 23rd International Conference on Software Analysis, Evolution, and Reengineering (SANER), Vol. 1, IEEE, 2016, pp. 624–628.
- [53] S. Gilda, Source code classification using neural networks, in: 2017 14th International Joint Conference on Computer Science and Software Engineering (JCSSE), IEEE, 2017, pp. 1–6.
- [54] G. Frantzeskou, S. MacDonell, E. Stamatatos, S. Gritzalis, Examining the significance of high-level programming features in source code author classification, Journal of Systems and Software 81 (3) (2008) 447–460.
- [55] V. Saini, F. Farmahinifarahani, Y. Lu, P. Baldi, C. V. Lopes, Oreo: Detection of clones in the twilight zone, in: Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering, 2018, pp. 354–365.
- [56] C. Fang, Z. Liu, Y. Shi, J. Huang, Q. Shi, Functional code clone detection with syntax and semantics fusion learning, in: Proceedings of the 29th ACM SIGSOFT International Symposium on Software Testing and Analysis, 2020, pp. 516–527.