

# Advising OpenMP Parallelization via a Graph-Based Approach with Transformers

Tal Kadosh<sup>1,2</sup>, Nadav Schneider<sup>2,3</sup>, Niranjana Hasabnis<sup>4</sup>, Timothy Mattson<sup>4</sup>, Yuval Pinter<sup>1</sup>, and Gal Oren<sup>5,6</sup>

<sup>1</sup> Computer Science Dept, Ben-Gurion University of the Negev, Israel

<sup>2</sup> Israel Atomic Energy Commission, Tel-Aviv, Israel

<sup>3</sup> Electrical & Computer Eng' Dept, Ben-Gurion University of the Negev, Israel

<sup>4</sup> Intel Labs

<sup>5</sup> Scientific Computing Center, Nuclear Research Center – Negev, Israel

<sup>6</sup> Computer Science Dept, Technion – Israel Institute of Technology

talkad@post.bgu.ac.il, nadavsch@post.bgu.ac.il,  
niranjana.hasabnis@intel.com, timothy.g.mattson@intel.com,  
uvp@cs.bgu.ac.il, galoren@cs.technion.ac.il

**Abstract.** There is an ever-present need for shared memory parallelization schemes to exploit the full potential of multi-core architectures. The most common parallelization API addressing this need today is OpenMP. Nevertheless, writing parallel code manually is complex and effort-intensive. Thus, many deterministic source-to-source (S2S) compilers have emerged, intending to automate the process of translating serial to parallel code. However, recent studies have shown that these compilers are impractical in many scenarios. In this work, we combine the latest advancements in the field of AI and natural language processing (NLP) with the vast amount of open-source code to address the problem of automatic parallelization. Specifically, we propose a novel approach, called OMPIFY, to detect and predict the OpenMP pragmas and shared-memory attributes in parallel code, given its serial version. OMPIFY is based on a Transformer-based model that leverages a graph-based representation of source code that exploits the inherent structure of code. We evaluated our tool by predicting the parallelization pragmas and attributes of a large corpus of (over 54,000) snippets of serial code written in C and C++ languages (*Open-OMP-Plus*). Our results demonstrate that OMPIFY outperforms existing approaches — the general-purposed and popular ChatGPT and targeted PRAGFORMER models — in terms of F1 score and accuracy. Specifically, OMPIFY achieves up to 90% accuracy on commonly-used OpenMP benchmark tests such as NAS, SPEC, and PolyBench. Additionally, we performed an ablation study to assess the impact of different model components and present interesting insights derived from the study. Lastly, we also explored the potential of using data augmentation and curriculum learning techniques to improve the model's robustness and generalization capabilities. The dataset and source code necessary for reproducing our results are available at <https://github.com/Scientific-Computing-Lab-NRCN/OMPify>.

**Keywords:** NLP · Code Completion · OpenMP · Shared Memory Parallelism · Transformers · S2S Compilers · Code Representations

## 1 Introduction

There is an ever-growing need to develop parallel applications these days. The ever-growing demand for computing power is leading to various types of complex architectures, including shared-memory multi-core architectures. A part of the demand arises from the recent HPCaaS paradigm that has become widespread and available to a broader community of developers [4]. The services offered as HPCaaS usually depend on the CPU core count and the duration of compute usage. Furthermore, the number of cores per CPU node has increased over the years — for example, from dozens of physical cores available in GCP’s C2 family [2] to hundreds of physical cores available in GCP’s future C3 family [9].

Despite the growing need to write parallel programs, introducing shared-memory parallelization into code remains challenging due to numerous pitfalls. Besides the fact that parallelizing serial code requires extensive knowledge of the code structure and semantics, it also requires the programmer to avoid parallelization pitfalls, such as the need to synchronize simultaneous reads and writes to the same variables (leading to race conditions), as well as making sure that the workload is distributed evenly across the threads and across the system resources (load balancing). In addition, it also requires a high degree of human expertise to comprehend fine details and abstract correlations between variables and different code segments [1]. It is then unsurprising that the number of parallel programming experts is relatively tiny compared to the growing community of users who can benefit from parallel programs.

The complexity of writing parallel programs is partly addressed by source-to-source (S2S) compilers [14,15,16], which are compilers that translate code from one programming language to another while preserving the code semantics. These compilers analyze the code for data dependencies that could prevent parallelization and automatically insert appropriate parallelization APIs (such as OpenMP *pragmas*) into it. Nevertheless, these compilers have several major drawbacks [23,34,35], such as long execution times and limited robustness to the input, even when optimized on runtime [28]. More importantly, these compilers require manual development and maintenance efforts, for instance, to support a new programming language or a new specification of parallel programming APIs.

We observed that the recent advances and successes of deep-learning-based Natural Language Processing (NLP) models, such as Transformer architecture and attention mechanism [38], masked language modeling (MLM) [17], could help in addressing the limitations of S2S compilers. These models are commonly called large language models (LLMs) because they capture the characteristics of languages. There are already examples of applying these LLMs to programming-related tasks. For example, Codex (based on GPT) [13], a state-of-the-art model that powers GitHub’s CoPilot [3], has shown an interesting application for generating code from natural language prompts. Another example is Google’s ML-enhanced code completion tool [5] that can predict possible completions of incomplete code fragments. An internal study conducted by Google has shown the promising potential of these technologies in reducing programming efforts. These

technologies are commonly deployed in programming editors and integrated development environments (IDEs) to provide immediate feedback to developers.

Although AI-based programmer assistance tools already exist, to our knowledge, PragFormer [22] is the only AI-based programmer assistance tool that can advise programmers in parallel programming. Specifically, PragFormer uses a Transformer-based architecture to predict if a given serial code could be parallelized (using OpenMP *pragma*) and if private or reduction clause could be applied to it. Specifically, it formulates this problem as multiple binary classification problems, where one problem tackles the need of determining if OpenMP *pragma* could be applied, while the other two tackle the need of determining the need of private and reduction clauses respectively.

While PragFormer has definitely shown interesting perspective towards automated parallel programming, in our experiments with PragFormer, we identified several of its limitations. One of its key limitations is the problem formulation; conceptually, if a serial code cannot be parallelized, then there is no need of determining private/reduction clause. As such, we found that these three are not independent problems, and rather formulating the problem as a multi-label classification problem seems much more intuitive. We address this and a few other limitations in PragFormer to propose a new model, named OMPIFY, that improves upon PragFormer on several fronts. Our experimental evaluation on a corpus of 54,000 *for-loops* mined from GitHub revealed that OMPIFY outperforms PragFormer and several state-of-the-art AI models for code in assisting programmers in parallel programming.

The rest of this article is organized as follows. Section 2 describes related work and provides the necessary background of our work. Section 3 presents the research objectives. Section 4 describes OMPIFY and illustrates our proposed method. Section 5 evaluates our method against previous methods. Finally, section 6 concludes this article and suggests possible extensions of this work.

## 2 Related Work

Initially, the approaches for translating serial code into parallel heavily relied on heuristics and rule-based methods, which often had limited capabilities and robustness (§2.1). However, with the rapid advancement of deep learning techniques in the field of NLP, along with the easy availability of open-source code, there have been some approaches to apply deep learning techniques to source code (§2.2). These approaches, however, process source code as text (similar to NLP) and fail to fully exploit the potential of other code representations (§2.3). By incorporating multiple code representations that capture different aspects of source code, multimodal learning techniques can overcome the limitations of these approaches.

### 2.1 Rule-based Methods

Several S2S compilers, including Cetus [15] and Par4All [14], have emerged in the last decade or so to insert OpenMP pragmas into code automatically. These

tools rely on program analysis-based techniques to analyze and identify potential constraints (e.g., loop-carried dependencies) that may restrict the code from being parallelized. The general workflow of S2S compilers can be summarized as follows:

1. Create an abstract syntax tree (AST) [29], which is a tree representation of the code’s syntactic structure. ASTs are constructed using source code parsers, such as *ANother Tool for Language Recognition* (ANTLR) [32] or *pycparser* [12], etc.
2. Apply data dependence algorithms [18] to ASTs.
3. Produce appropriate OpenMP directives based on the data dependence graph.

There are multiple drawbacks associated with the approach of generating ASTs and applying data dependence algorithms. Firstly, creating an AST with a parser can be a challenging task with limited robustness to input due to each programming language’s unique syntactic structures that have evolved over the years. Thus, many S2S compilers cannot handle the diverse syntax of programming languages. Moreover, not all parsers are publicly available. As a result, some S2S compilers may fail to produce an AST and analyze the input code. Secondly, data-dependence algorithms can be time-consuming, particularly for large-scale code, since these algorithms are strongly dependent on the size of the AST, which in turn is influenced by the length of the code. Additionally, studies by Harel et al. [23] and Prema et al. [34] have shown that S2S compilers may produce sub-optimal results and even degrade program performance in some cases.

## 2.2 Unimodal Machine-Learning Driven Methods

Rule-based methods also suffer from another important limitation – tools relying on these methods require manual programming efforts to add new rules to maintain and update them. However, recent AI-based programming assistance tools have demonstrated that it is possible to reduce manual effort by instead learning the rules from data. Specifically, with the powerful computing devices and vast availability of open-source code as data, these AI-based tools can learn programming rules such as syntax, typing rules [24], etc. Continuing this trend, in recent years, several Transformer-based models have been proposed for various programming-related tasks [20,31]. Typically, these models are pre-trained on massive code corpora containing multiple programming languages (PLs) and then applied to various programming problems [30], such as program completion, code search, bug finding, etc., as downstream tasks. One of the common pre-training tasks is masked language modeling (MLM) [17].

Previous work by Harel et al. [22] showed the possibility of applying Attention-based models (Transformer) to determine if code can be parallelized with OpenMP. In their work, they introduced PRAGFORMER, which is a transformer model based on DeepSCC [41], which itself is a *RoBERTa* model fine-tuned on a corpus of 225k code snippets written in 21 programming languages (such as Java, Python, C, and C++) collected from Stack Overflow.

In the parlance of AI-based models, PragFormer formulates the parallel programming assistance problem as *Code Language Processing for Parallelization (CLPP)* task. Specifically, it breaks this task down into three sub-problems: given a serial code (*for-loop*), determine (1) if it can be parallelized (using OpenMP *pragma*), (2) if private clause would be applicable to OpenMP *pragma*, and (3) if reduction class would be applicable to OpenMP *pragma*. It then approaches these three sub-problems independently and formulates them as three separate binary classification problems as below:

1. *pragma classification*: Classifying the need for OpenMP *parallel for pragma*.
2. *private clause classification*: Classifying the need for a *private* clause (specifying a variable to be private to each thread in a parallel region).
3. *reduction clause classification*: Classifying the need for a *reduction* clause (specifying an operator and a variable to reduce across all threads in a parallel region).

Although PragFormer shows great potential in using Transformer architecture to solve the shared-memory parallelization task, it still suffers from several deficiencies. Primarily, PRAGFORMER is based on *RoBERTa*, which is essentially a model for Natural Language (NL) understanding. Applying an NL model to a code-related task is sub-optimal compared to models pre-trained directly on code [30]. Additionally, PRAGFORMER regards source code as a sequence of tokens, ignoring the inherent structure of the code. Intuitively, structural information of code, such as variable dependence information, etc., should provide crucial code semantic information that could improve the code understanding process. Furthermore, the approach of separating the classifications is unintuitive since the tasks of predicting the need for OpenMP *pragmas* and data-sharing attribute clauses are highly correlated — there will not be a private or reduction clause if there is no need for OpenMP *pragma* at all.

### 2.3 Multimodal Machine-Learning Driven Methods

While the unimodal ML methods accept source code in only one representation (most commonly as a sequence of tokens), multimodal ML methods realize that other code representations may offer richer semantic information that could improve the accuracy of the models on programming-related tasks. Specifically, multimodal ML models also accept source code in other representations such as AST, control-flow graph (CFG), data-flow graph (DFG), etc. Consequently, many pre-trained machine learning models have been developed, with each model incorporating different code formats into the training process.

Feng et al. presented CODEBERT [19], a bimodal Transformer model trained on programming languages alongside natural languages. In their experiments, they used *CodeSearchNet* [26] dataset, which includes 6.4M code snippets from 6 programming languages (Python, Java, JavaScript, Go, Ruby, and PHP). They compared CODEBERT trained on samples from natural languages and programming languages, CODEBERT trained only on natural languages, and

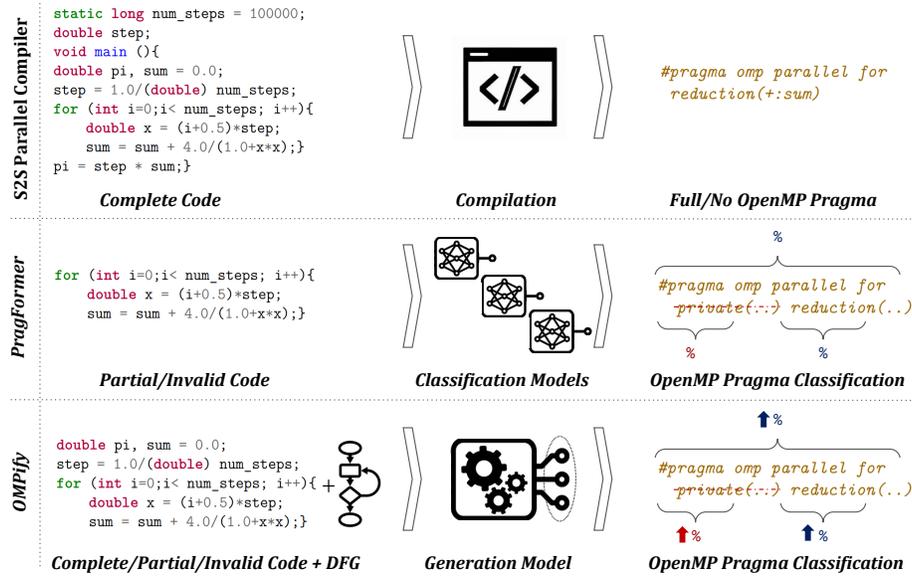


Fig. 1: Differences between S2S compilers, PRAGFORMER and OMPIFY.

ROBERTA, and showed the superiority of CODEBERT trained on both natural languages and programming languages on several programming-related tasks. Guo et al. created GRAPHCODEBERT [20], a multimodal Transformer model trained on *CodeSearchNet* dataset and input programs as natural language text alongside programming languages and DFG. They showed that DFG enhances the code understanding process compared to CODEBERT. Another multimodal model that exploits the structural aspect of the code is SPT-CODE [31] that was proposed by Niu et al. They presented SPT-CODE, which is trained on natural language, programming language, and AST, from the *CodeSearchNet* dataset. While CODEBERT and GRAPHCODEBERT are models that use Transformer encoders, SPT-CODE is uses Transformer encoder-decoder architecture. Experiments have shown the superiority of SPT-CODE in code generation tasks.

Drawing inspiration from some of the design choices of multimodal models, we have designed the OMPIFY model also as a multimodal model. Figure 1 summarizes the difference between OMPIFY and related works.

### 3 Research Objectives

This paper draws inspiration from PragFormer and approaches the parallel programming assistance problem as a *Code Language Processing for Parallelization (CLPP)* task. Nevertheless, we improve upon PragFormer by posing the following research questions that are designed to evaluate the limitations of PragFormer discussed earlier in Related Work (§2.2).

**RQ1:** *Which code representations impact the CLPP task?*

Given the discussion of different code representations, this question focuses on assessing the influence of various code modalities on code comprehension, particularly in CLPP tasks. We will evaluate the effectiveness of the previously mentioned multimodal models on our dataset.

**RQ2:** *Does the scope of the for-loop from input serial code matter for performance on CLPP task?*

Conceptually, the semantics of a *for-loop* from serial code heavily relies on its *context*. This question assesses the impact of different context lengths on the performance of various multimodal models on CLPP tasks.

**RQ3:** *Can code augmentation improve model’s performance on CLPP task?*

We will investigate the potential of code augmentation techniques, specifically variable name replacement, to improve the performance of existing models on CLPP tasks.

**RQ4:** *Will multi-label classification based formulation for solving CLPP task perform better than PragFormer’s multiple binary-classification based formulation?*

While PRAGFORMER employed three binary classification-based models to predict the requirement for an OpenMP *pragma* and whether it should include a work-sharing construct, we hypothesize that these predictions are interdependent and potentially benefit each other. We evaluate this hypothesis by formulating a multi-label classification problem and developing a single generative model for the CLPP task. We then compare our model against PragFormer and other state-of-the-art multimodal models.

## 4 OMPify

This section describes the model architecture, its input, the code representations, and the fine-tuning process. OMPIFY predicts the need for both OpenMP *pragma* and shared-memory attributes (*private* and *reduction*) simultaneously, allowing the model to learn inter-dependencies between these tasks.

### 4.1 Model

OMPIFY (Figure 2) is a Transformer-based, multimodal model that utilizes a graph-based representation of source code. OMPIFY is based on GRAPHCODEBERT [20], a pre-trained model for programming languages that considers the inherent structure of the code by accepting source code along with its DFG. OMPIFY is composed of GRAPHCODEBERT and a fully connected layer. This architecture allows OMPIFY to perform multi-label classification, where each task is individually classified.

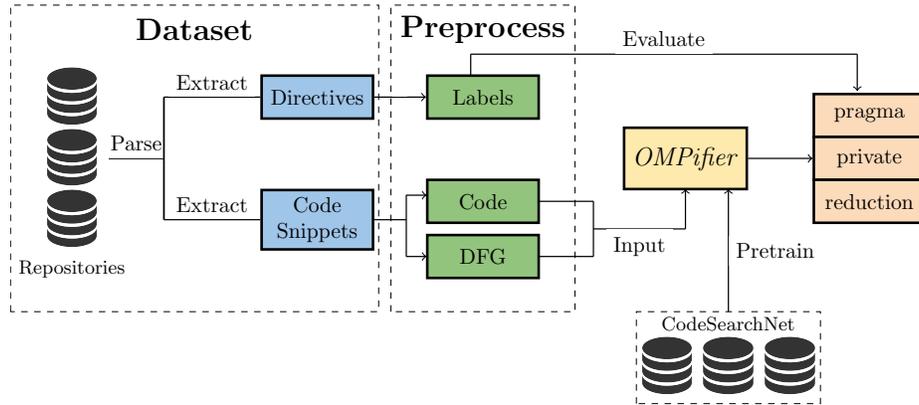


Fig. 2: Overview of OMPiFY training process.

## 4.2 Model Input

The model’s inputs are two code modalities: the actual source code as a sequence of tokens and the serialized DFG.

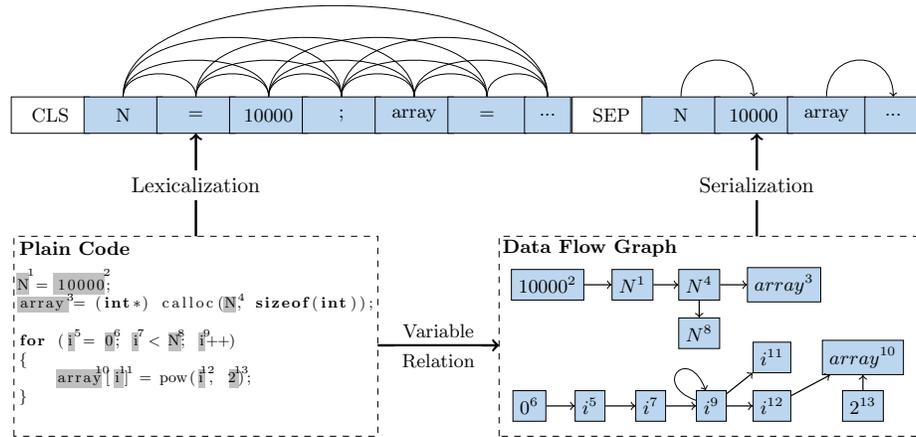


Fig. 3: The input format for a C code snippet.

- **Code Tokens.** As shown in Figure 3, the first part of the input to OMPiFY consists of a sequence of code tokens. Feeding source code as a sequence of strings to Transformer-based models does not work well in practice. A common approach, in this case, is to train a tokenizer that maps strings of tokens into unique IDs (numbers), such that strings that are semantically closer have their IDs closer also. We used the tokenizer provided by GRAPHCODEBERT to generate a sequence of token IDs.

- **Serialized DFG.** The second part of the input, which is the serialized DFG, is created by converting the code into an AST using TreeSitter parser <sup>7</sup>. We extract variables and their data dependence relationships from an AST to generate a DFG. The DFG nodes are serialized in the program ordered and serve as the model input.
- **Attention Mask.** Figure 3 provides an overview of the token connections, showcasing the interconnections between the tokens. Whereas the code tokens attend to each other during the self-attention mechanism [39], when dealing with DFG, we aim to disregard attention between variables that are not connected. To achieve this, we employed the masked attention approach described by Guo et al. [20]. During the computation of self-attention, we utilized an attention mask, denoted as  $M$ , which contains zeros in positions  $(i, j)$  where we want the tokens to attend to each other, and contains  $-inf$  to prevent certain relations. The self-attention computation can be represented as follows:
 
$$SelfAttention = \text{softmax}(QK/d + M)V \quad (1)$$

### 4.3 Fine Tuning

Many studies have demonstrated the advantages of implementing data augmentation techniques to enhance the performance of deep learning models [37,40]. Data augmentation techniques are typically applied to the training set to increase the diversity of input. Commonly-used augmentation techniques include *variable renaming*, *dead store*, and *constant replacement*.

Despite the effectiveness of data augmentations, several studies [25,36] showed that deep learning models are vulnerable to adversarial examples, i.e., minor changes to code can result in significant performance degradation. To address this issue, as many studies have suggested [21,33,40] leveraging a curriculum learning (CL) technique that involves the gradual introduction of data augmentation techniques. Specifically, we applied *variable renaming* as a data augmentation technique. We followed a gradual approach, starting with the original data without any augmentation during the first epoch. In each subsequent epoch, we augmented the original data by progressively increasing the proportion of renamed variables. Specifically, during the second epoch, we renamed 10% of the variables for each sample in the original training set. In the third epoch, we continued this approach and renamed 20% of the variables per sample. As we progressed to the fourth epoch, we further increased the update ratio, renaming 30% of the variables per sample. Finally, starting from the fifth epoch and throughout the remaining epochs, we consistently maintained an update ratio of 40%, resulting in the renaming of 40% of the variables per sample in each subsequent epoch.

## 5 Experimental Results

To evaluate the effectiveness of our proposed model, OMPIFY, we conducted several experiments to answer our research questions. All experiments were con-

<sup>7</sup> <https://github.com/tree-sitter/tree-sitter>

ducted on an NVIDIA A100 GPU. Furthermore, for the sake of consistency, we utilized the original implementations of the models as presented in their respective papers.

## 5.1 Dataset & Preprocessing

In our work, we developed a novel dataset, named *Open-OMP-Plus*, comprising more than 54,000 code snippets from C and C++ for OpenMP analysis (Table 1). The dataset was collected from `github.com` using the *github-clone-all*<sup>8</sup> script, which enables searching for repositories that satisfy specific criteria. We used this tool to locate all repositories that include C or C++ files and also feature the term “OpenMP” in their title, description, or README.

| Description    | C      | C++    | Clauses          | Amount | # Lines | Amount |
|----------------|--------|--------|------------------|--------|---------|--------|
| With OpenMP    | 14,906 | 8,241  | <i>private</i>   | 6,758  | < 15    | 40,745 |
| Without OpenMP | 17,193 | 14,323 | <i>reduction</i> | 3,267  | 16-50   | 10,607 |
| Total          | 32,099 | 22,564 | Total            | 10,025 | > 50    | 3,311  |

(a) Number of loops paralleled with OpenMP for each programming language (C and C++).

(b) Number of common OpenMP shared memory attributes.

(c) Code snippet length in *Open-OMP-Plus*.

Table 1: The distribution of each class for each programming language.

To minimize the noise in our dataset, we employed inclusion and exclusion criteria inspired by Harel et al. [22]. Specifically, we included only C/C++ files that contained OpenMP *pragmas* in their code. This criteria operates on the assumption that the developers were aware of OpenMP parallelization and that any non-parallelized loops intentionally have not been parallelized. We excluded duplicates, empty loops, and loops that used *barrier*, *critical*, or *atomic pragmas*, which can be bottlenecks on code execution and are not optimal samples.

Once we identified files that contained OpenMP pragmas, we parsed them using *pycparser* [12] parser, which converts the code into an AST format. Each sample in the dataset comprises several fields, such as the plain *for-loop* code, its corresponding *pragma* (if any), the AST of the *for-loop*, the AST of the functions called within the *for-loop*, the declaration of each variable used in the *for-loop*, all the assignment instructions from the context of the loop that involves each of the variables used in the loop, and the DFG of the *for-loop* and its extended scope. By analyzing the AST, we can extract code structures such as loops and identify the relevant functions and variables from its outer scope.

To evaluate the performance of our model, we divided the dataset into three sets: train, validation, and test, using standard 80-10-10 split. Additionally, we collected three benchmarks that were known to use OpenMP correctly, namely

<sup>8</sup> <http://github.com/rhysd/github-clone-all>

NAS [11], PolyBench [7], and SPEC [8], and used them to further test our model. To avoid fair evaluation, we removed from the training set any samples that could be found in the benchmarks.

## 5.2 Results

We now present the results of our experiments to answer the research questions.

**RQ1: Code modalities.** To compare the various code modalities, we utilized three distinct models. CODEBERT was pre-trained on natural language (NL) and programming language (PL) and SPT-CODE was pre-trained on PL and AST, and GRAPHCODEBERT was pre-trained on PL and DFG. To apply these models to *pragma classification* task, we added a fully connected layer of size two

and a SoftMax layer at the end of these models and fine-tuned them using our corpus, *Open-OMP-Plus*. We also included PRAGFORMER, which was fine-tuned on our corpus, for comparison with these models.

Based on the results presented in Table 2, it can be inferred that the use of multimodal models, which combine code representations such as AST or DFG with the original code, has a positive impact on the performance of the model in the *pragma classification* task. However, despite being trained on the same dataset, i.e., *CodeSearchNet*, CODEBERT, SPT-CODE, and GRAPHCODEBERT achieved significantly different performance, with GRAPHCODEBERT outperforming others. This indicates that the DFG representation and pretraining tasks proposed by Guo et al. [20] were more beneficial this task. It is worth noting that while the use of AST in SPT-CODE significantly improves performance compared to using only the code, the DFG in GRAPHCODEBERT has only a minimal impact on performance. This could be due to the DFG representation’s inability to effectively capture the relationship between arrays and their indexing. As shown in Figure 3, there is no direct connection between the array and the index variable  $i^{11}$ . In scientific codes, the relationship between the array and the index is often critical as it determines the feasibility of parallelization.

**RQ2: Extended scope.** In this experiment, we aimed to investigate the impact of extended scope on the performance of OMPIFY in solving the CLPP task. To achieve this, we trained OMPIFY on two distinct corpora: one comprising only the *for-loop* structured block, and the other consisting of the extended version that includes the surrounding scope of the *for-loop*, incorporating assignments

| Model Name                | Metrics      |              |              |
|---------------------------|--------------|--------------|--------------|
|                           | P            | R            | Acc          |
| PragFormer                | 0.826        | 0.780        | 0.830        |
| CodeBERT                  | <b>0.848</b> | 0.813        | 0.852        |
| SPT-Code                  | 0.812        | 0.784        | 0.831        |
| SPT-Code (code only)      | 0.792        | 0.786        | 0.820        |
| GraphCodeBERT             | 0.836        | <b>0.835</b> | <b>0.862</b> |
| GraphCodeBERT (code only) | 0.834        | 0.833        | 0.861        |

Table 2: Effect of different code modalities on the task of *pragma classification*. (P=Precision, R=Recall, Acc=Accuracy)

to variables used within the loop (Figure 5). The results, as presented in Table 4, demonstrate the effect of including the outside scope in determining the necessity of OpenMP *pragma*. The observed increase in recall indicates that the model exhibits improved identification of *for-loops* requiring OpenMP *pragma*, resulting in fewer false negatives. This finding suggests that considering the outside scope provides valuable information for accurately identifying the need for *pragma* in *for-loops*.

| Data Type  | P            | R            | Acc          |
|------------|--------------|--------------|--------------|
| No Scope   | <b>0.833</b> | 0.831        | 0.860        |
| With Scope | 0.829        | <b>0.844</b> | <b>0.863</b> |

Fig. 4: Effect of *context*. (P=Precision, R=Recall, Acc=Accuracy)

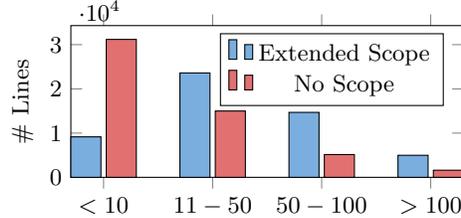


Fig. 5: Code length comparison.

### RQ3: Data augmentation.

Through evaluating the impact of data augmentation techniques on performance, we investigated the effectiveness of PRAGFORMER and GRAPHCODEBERT in the binary classification task of OpenMP *pragma* classification. The results of this experiment are presented in Table 3. We employed the *variable renaming* augmentation, where each variable was replaced with *var*, concatenated with a random index number. This augmentation is referred as *replaced* in the table. The results reveal the vulnerability of these models to adversarial examples created by fully replacing variable names, leading to degraded performance compared to the unmodified variables. However, by gradually introducing code augmentations using the curriculum learning method (referred as *curriculum* in the table), we observed improved accuracy.

| Model         | Augmentation | Metrics      |              |              |
|---------------|--------------|--------------|--------------|--------------|
|               |              | P            | R            | Acc          |
| PragFormer    | original     | 0.793        | <b>0.847</b> | 0.841        |
|               | curriculum   | 0.825        | 0.815        | 0.848        |
|               | replaced     | 0.727        | 0.826        | 0.794        |
| GraphCodeBERT | original     | <b>0.851</b> | 0.841        | 0.870        |
|               | curriculum   | 0.849        | 0.846        | <b>0.872</b> |
|               | replaced     | 0.838        | 0.781        | 0.843        |

Table 3: Effect of data augmentation techniques. (P=Precision, R=Recall, Acc=Accuracy)

### RQ4: Multi-label classification.

Table 4 shows the results of PRAGFORMER, GRAPHCODEBERT, and OMPIFY when applied to all the three tasks of *pragma* classification, classification of *private* clause, and *reduction* clause.

| Model                     | Task             | Metrics      |              |              |
|---------------------------|------------------|--------------|--------------|--------------|
|                           |                  | P            | R            | Acc          |
| PragFormer                | <i>pragma</i>    | 0.793        | 0.847        | 0.841        |
|                           | <i>private</i>   | 0.716        | 0.663        | 0.924        |
|                           | <i>reduction</i> | 0.632        | 0.598        | 0.953        |
| GraphCodeBERT (Separated) | <i>pragma</i>    | <b>0.850</b> | 0.841        | 0.870        |
|                           | <i>private</i>   | <b>0.768</b> | 0.684        | 0.937        |
|                           | <i>reduction</i> | <b>0.690</b> | 0.688        | 0.963        |
| OMPify                    | <i>pragma</i>    | 0.849        | <b>0.848</b> | <b>0.872</b> |
|                           | <i>private</i>   | 0.755        | <b>0.689</b> | <b>0.938</b> |
|                           | <i>reduction</i> | <b>0.690</b> | <b>0.700</b> | <b>0.966</b> |

Table 4: Effect of multi-label classification problem formulation. (P=Precision, R=Recall, Acc=Accuracy)

Note that OMPIFY approaches all three tasks together as a multi-label classification problem, while PRAGFORMER, GRAPHCODEBERT approach each task independently. In the table, we present the result of OMPIFY for the combined task but split the results according to labels.

The results convey that OMPIFY achieves significantly better performance compared to PRAGFORMER, underscoring the hypothesis that these three tasks are not independent. In addition, our model slightly outperforms the GRAPHCODEBERT model. The results show a significant improvement in recall for OMPIFY for all three tasks, with a major decrease in the number of false negative predictions. In our context, a false negative prediction means that a sample is incorrectly classified as not requiring *pragma*, *private*, or *reduction*. Therefore, the unified prediction strategy of OMPIFY can better identify samples that require *pragma*, *private*, or *reduction*. This suggests that the understanding of each task contributes to the overall prediction. For instance, if OMPIFY predicts the need for *pragma*, it will also influence the prediction of shared-memory attributes, which may also appear in the *pragma*.

**Real-world benchmarks.** In order to test the performance of OMPIFY on real-world programs, we obtained C/C++ programs that were using OpenMP *pragma* from three scientific code benchmarks, namely, NAS, SPEC, and PolyBench. Table 5

| Benchmark | With OMP | Without OMP | <i>private</i> | <i>reduction</i> |
|-----------|----------|-------------|----------------|------------------|
| NAS       | 166      | 146         | 12             | 2                |
| PolyBench | 63       | 85          | 36             | 0                |
| SPEC      | 157      | 1,000       | 1              | 0                |

Table 5: Benchmark statistics

shows the statistics of the collected programs. These benchmarks are manually-written as parallel programs using OpenMP, so they serve as a good test case for OMPIFY. As a comparison, we applied PRAGFORMER to the same test. Our model exhibited a significant increase in performance when compared to PRAGFORMER (Table 6).

Moreover, given the recent popularity of ChatGPT [10] in programming-related tasks, we decided to evaluate it on our CLPP task. For this evaluation, we randomly sampled 2500 test inputs from our test dataset. We then fed those test programs to ChatGPT one by one and then used the prompt “*Generate the optimal OpenMP pragma if possible*” to check if ChatGPT’s response matches with the expected label for the test program. Although ChatGPT performs well on various NLP tasks, it performed

| Benchmark     | Model      | Metrics      |              |              |
|---------------|------------|--------------|--------------|--------------|
|               |            | P            | R            | Acc          |
| SPEC          | PRAGFORMER | 0.445        | 0.802        | 0.837        |
|               | OMPIFY     | <b>0.572</b> | <b>0.854</b> | <b>0.894</b> |
| PolyBench     | PRAGFORMER | 0.703        | 0.301        | 0.648        |
|               | OMPIFY     | <b>0.836</b> | <b>0.810</b> | <b>0.851</b> |
| NAS           | PRAGFORMER | 0.635        | 0.734        | 0.634        |
|               | OMPIFY     | <b>0.731</b> | <b>0.886</b> | <b>0.766</b> |
| 2500 examples | ChatGPT    | 0.401        | <b>0.913</b> | 0.401        |
|               | PRAGFORMER | 0.815        | 0.721        | 0.817        |
|               | OMPIFY     | <b>0.839</b> | 0.818        | <b>0.860</b> |

Table 6: Comparison on different benchmarks. (P=Precision, R=Recall, Acc=Accuracy)

Although ChatGPT performs well on various NLP tasks, it performed

poorly in our specific task, often suggesting the use of OpenMP *pragma* even when it was not applicable.

## 6 Conclusions & Future Work

This paper aims to investigate the potential of multimodal models in accurately predicting the need for shared-memory parallelization in code. Our research discovered that incorporating additional code representations, such as ASTs and DFGs, significantly improves their performance compared to models that rely solely on the original code. Building upon this knowledge, we introduced a novel model called OMPIFY, based on GRAPHCODEBERT. OMPIFY takes advantage of the inter-dependencies between the task of predicting the need for parallelization and the prediction of shared-memory attributes, such as *private* and *reduction* variables. By leveraging these relationships, OMPIFY demonstrates enhanced accuracy and robustness in determining the need for shared-memory parallelization. In addition to developing the OMPIFY model, we also constructed a comprehensive database called *Open-OMP-Plus*. This database includes the *for-loop* itself and extends its scope to include assignment statements of variables found within the *for-loop*. By incorporating this extended scope, we demonstrate that OMPIFY can effectively utilize this additional information to improve its predictions further.

For future research, we aim to address several areas of improvement. Firstly, since the multimodal models analyzed in RQ1 were not pre-trained on C/C++ programming languages, there is a potential for enhancing their performance by pretraining them on datasets that include C/C++ code. This approach can contribute to better code understanding and comprehension. Additionally, in RQ4, we observed improvements in multi-label prediction. To further enhance this aspect, we intend to explore the conversion of the multi-label prediction problem into *pragma* generation. By generating *pragmas* directly, we can achieve more precise and fine-grained control over parallelization tasks. Furthermore, an important question arises regarding the correctness of the generated *pragmas*. To address this concern, we plan to investigate techniques and approaches for evaluating the accuracy and correctness of the generated *pragmas*.

**Acknowledgments:** This research was supported by the Israeli Council for Higher Education (CHE) via the Data Science Research Center, Ben-Gurion University of the Negev, Israel; Intel Corporation (oneAPI CoE program); and the Lynn and William Frankel Center for Computer Science. Computational support was provided by the NegevHPC project [6] and Intel Developer Cloud [27]. The authors thank Re'em Harel, Israel Hen, and Gabi Dadush for their help and support.

## References

1. Automatic Parallelism and Data Dependency. <https://web.archive.org/web/20140714111836/http://blitzprog.org/posts/automatic-parallelism-and-data-dependency>, [Online]
2. Compute-optimized machine family. <https://cloud.google.com/compute/docs/compute-optimized-machines>, [Online]
3. GitHub Copilot. <https://github.com/features/copilot>, [Online]
4. High Performance Computing as a Service Market Forecast. <https://www.alliedmarketresearch.com/high-performance-computing-as-a-service-market>, [Online]
5. ML-enhanced code completion improves developer productivity. <https://ai.googleblog.com/2022/07/ml-enhanced-code-completion-improves.html>, [Online]
6. NegevHPC Project. <https://www.negevhpc.com>, [Online]
7. PolyBench Benchmarks. <https://web.cse.ohio-state.edu/~pouchet.2/software/polybench/>, [Online]
8. SPEC-OMP2012 website. <https://www.spec.org/omp2012/>, [Online]
9. The next wave of Google Cloud infrastructure innovation: New C3 VM and Hyperdisk. <https://cloud.google.com/blog/products/compute/introducing-c3-machines-with-googles-custom-intel-ipu>, [Online]
10. ChatGPT. <https://chat.openai.com/> (2023), [Online]
11. Bailey, D.H., Barszcz, E., Barton, J.T., Browning, D.S., Carter, R.L., Dagum, L., Fatoohi, R.A., Frederickson, P.O., Lasinski, T.A., Schreiber, R.S., et al.: The nas parallel benchmarks. *The International Journal of Supercomputing Applications* **5**(3), 63–73 (1991)
12. Bendersky, E., et al.: *Pycparser* (2010)
13. Chen, M., Tworek, J., Jun, H., Yuan, Q., Pinto, H.P.d.O., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., et al.: Evaluating large language models trained on code. arXiv preprint arXiv:2107.03374 (2021)
14. Creusillet, B., Keryell, R., Even, S., Guelton, S., Irigoin, F.: *Par4all: Auto-parallelizing c and fortran for the cuda architecture* (2009)
15. Dave, C., Bae, H., Min, S.J., Lee, S., Eigenmann, R., Midkiff, S.: *Cetus: A source-to-source compiler infrastructure for multicores*. *Computer* **42**(12) (2009)
16. Dever, M.: *AutoPar: automating the parallelization of functional programs*. Ph.D. thesis, Dublin City University (2015)
17. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: BERT: Pre-training of deep bidirectional transformers for language understanding. In: *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. pp. 4171–4186. Association for Computational Linguistics, Minneapolis, Minnesota (Jun 2019). <https://doi.org/10.18653/v1/N19-1423>, <https://aclanthology.org/N19-1423>
18. Fagin, R., Vardi, M.Y.: *The theory of data dependencies: a survey*. IBM Thomas J. Watson Research Division (1984)
19. Feng, Z., Guo, D., Tang, D., Duan, N., Feng, X., Gong, M., Shou, L., Qin, B., Liu, T., Jiang, D., et al.: *Codebert: A pre-trained model for programming and natural languages*. arXiv preprint arXiv:2002.08155 (2020)
20. Guo, D., Ren, S., Lu, S., Feng, Z., Tang, D., Liu, S., Zhou, L., Duan, N., Svyatkovskiy, A., Fu, S., et al.: *Graphcodebert: Pre-training code representations with data flow*. arXiv preprint arXiv:2009.08366 (2020)

21. Guo, S., Huang, W., Zhang, H., Zhuang, C., Dong, D., Scott, M.R., Huang, D.: Curriculumnet: Weakly supervised learning from large-scale web images. In: Proceedings of the European conference on computer vision (ECCV). pp. 135–150 (2018)
22. Harel, R., Pinter, Y., Oren, G.: Learning to parallelize in a shared-memory environment with transformers. In: Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming. pp. 450–452 (2023)
23. Harel, R., Mosseri, I., Levin, H., Alon, L.o., Rusanovsky, M., Oren, G.: Source-to-source parallelization compilers for scientific shared-memory multi-core and accelerated multiprocessing: analysis, pitfalls, enhancement and potential. *International Journal of Parallel Programming* **48**(1), 1–31 (2020)
24. Hasabnis, N., Gottschlich, J.: Controlflag: A self-supervised idiosyncratic pattern detection system for software control structures. In: Proceedings of the 5th ACM SIGPLAN International Symposium on Machine Programming. p. 32–42. MAPS 2021, Association for Computing Machinery, New York, NY, USA (2021). <https://doi.org/10.1145/3460945.3464954>, <https://doi.org/10.1145/3460945.3464954>
25. Henke, J., Ramakrishnan, G., Wang, Z., Albarghouth, A., Jha, S., Reps, T.: Semantic robustness of models of source code. In: 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER). pp. 526–537. IEEE (2022)
26. Husain, H., Wu, H.H., Gazit, T., Allamanis, M., Brockschmidt, M.: Codesearchnet challenge: Evaluating the state of semantic code search. arXiv preprint arXiv:1909.09436 (2019)
27. Intel: Intel Developer Cloud. <https://www.intel.com/content/www/us/en/developer/tools/devcloud/overview.html> (2023), [Online]
28. Mosseri, I., Alon, L., Harel, R., Oren, G.: Compar: Optimized multi-compiler for automatic openmp S2S parallelization. In: Milfeld, K.F., de Supinski, B.R., Koesterke, L., Klinkenberg, J. (eds.) *OpenMP: Portable Multi-Level Parallelism on Modern Systems - 16th International Workshop on OpenMP, IWOMP 2020*, Austin, TX, USA, September 22-24, 2020, Proceedings. *Lecture Notes in Computer Science*, vol. 12295, pp. 247–262. Springer (2020). [https://doi.org/10.1007/978-3-030-58144-2\\_16](https://doi.org/10.1007/978-3-030-58144-2_16), [https://doi.org/10.1007/978-3-030-58144-2\\_16](https://doi.org/10.1007/978-3-030-58144-2_16)
29. Neamtiu, I., Foster, J.S., Hicks, M.: Understanding source code evolution using abstract syntax tree matching. *ACM SIGSOFT Software Engineering Notes* **30**(4), 1–5 (2005)
30. Niu, C., Li, C., Ng, V., Chen, D., Ge, J., Luo, B.: An empirical comparison of pre-trained models of source code. arXiv preprint arXiv:2302.04026 (2023)
31. Niu, C., Li, C., Ng, V., Ge, J., Huang, L., Luo, B.: SPT-Code: Sequence-to-sequence pre-training for learning the representation of source code. arXiv preprint arXiv:2201.01549 (2022)
32. Parr, T.: *The definitive ANTLR 4 reference*. Pragmatic Bookshelf (2013)
33. Platanios, E.A., Stretcu, O., Neubig, G., Poczos, B., Mitchell, T.M.: Competence-based curriculum learning for neural machine translation. arXiv preprint arXiv:1903.09848 (2019)
34. Prema, S., Jehadeesan, R., Panigrahi, B.: Identifying pitfalls in automatic parallelization of nas parallel benchmarks. In: *Parallel Computing Technologies (PAR-COMPTECH)*, 2017 National Conference on. pp. 1–6. IEEE (2017)

35. Prema, S., Nasre, R., Jehadeesan, R., Panigrahi, B.: A study on popular auto-parallelization frameworks. *Concurrency and Computation: Practice and Experience* **31**(17), e5168 (2019)
36. Quiring, E., Maier, A., Rieck, K., et al.: Misleading authorship attribution of source code using adversarial learning. In: *USENIX Security Symposium*. pp. 479–496 (2019)
37. Rebuffi, S.A., Goyal, S., Calian, D.A., Stimberg, F., Wiles, O., Mann, T.A.: Data augmentation can improve robustness. *Advances in Neural Information Processing Systems* **34**, 29935–29948 (2021)
38. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
39. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *CoRR* **abs/1706.03762** (2017), <http://arxiv.org/abs/1706.03762>
40. Wang, D., Jia, Z., Li, S., Yu, Y., Xiong, Y., Dong, W., Liao, X.: Bridging pre-trained models and downstream tasks for source code understanding. In: *Proceedings of the 44th International Conference on Software Engineering*. pp. 287–298 (2022)
41. Yang, G., Zhou, Y., Yu, C., Chen, X.: Deepsc: Source code classification based on fine-tuned roberta. *CoRR* **abs/2110.00914** (2021), <https://arxiv.org/abs/2110.00914>