

Deep Language Models for Software Testing and Optimisation

Foivos Tsimpourlas



Doctor of Philosophy

Laboratory for Foundations of Computer Science

CDT Pervasive Parallelism

School of Informatics

The University of Edinburgh

2023

Abstract

Developing software is difficult. A challenging part of production development is ensuring programs are correct and fast, two properties satisfied with software testing and optimisation. While both tasks still rely on manual effort and expertise, the recent surge in software applications has led them to become tedious and time-consuming. Under this fast-pace environment, manual testing and optimisation hinders productivity significantly and leads to error-prone or sub-optimal programs that waste energy and lead users to frustration. In this thesis, we propose three novel approaches to automate software testing and optimisation with modern language models based on deep learning. In contrast to our methods, existing few techniques in these two domains have limited scalability and struggle when they face real-world applications.

Our first contribution lies in the field of software testing and aims to automate the test oracle problem, which is the procedure of determining the correctness of test executions. The test oracle is still largely manual, relying on human experts. Automating the oracle is a non-trivial task that requires software specifications or derived information that are often too difficult to extract. We present the first application of deep language models over program execution traces to predict runtime correctness. Our technique classifies test executions of large-scale codebases used in production as “pass” or “fail”. Our proposed approach reduces by 86% the amount of test inputs an expert has to label by training only on 14% and classifying the rest automatically.

Our next two contributions improve the effectiveness of compiler optimisation. Compilers optimise programs by applying heuristic-based transformations constructed by compiler engineers. Selecting the right transformations requires extensive knowledge of the compiler, the subject program and the target architecture. Predictive models have been successfully used to automate heuristics construction but their performance is hindered by a shortage of training benchmarks in quantity and feature diversity. Our next contributions address the scarcity of compiler benchmarks by generating human-likely synthetic programs to improve the performance of predictive models.

Our second contribution is BENCHPRESS, the first steerable deep learning synthesizer for executable compiler benchmarks. BENCHPRESS produces human-like programs that compile at a rate of 87%. It targets parts of the feature space previously unreachable by other synthesizers, addressing the scarcity of high-quality training data for compilers. BENCHPRESS improves the performance of a device mapping predictive model by 50% when it introduces synthetic benchmarks into its training data.

BENCHPRESS is restricted by a feature-agnostic synthesizer that requires thousands of random inferences to select a few that target the desired features. Our third contribution addresses this inefficiency. We develop BENCHDIRECT, a directed language model for compiler benchmark generation. BENCHDIRECT synthesizes programs by jointly observing the source code context and the compiler features that are targeted. This enables efficient steerable generation on large scale tasks. Compared to BENCHPRESS, BENCHDIRECT matches successfully $1.8\times$ more Rodinia target benchmarks, while it is up to 36% more accurate and up to 72% faster in targeting three different feature spaces for compilers.

All three contributions demonstrate the exciting potential of deep learning and language models to simplify the testing of programs and the construction of better optimization heuristics for compilers. The outcomes of this thesis provides developers with tools to keep up with the rapidly evolving landscape of software engineering.

Acknowledgements

First of all, I would like to express my deep gratitude to my three advisors Ajitha Rajan, Hugh Leather and Pavlos Petoumenos for providing me with the best life, research, professional advice that I could ever wish for. Their mentorship throughout this whole journey was crucial for its successful completion. They were always there for me and they influenced me, shaped me and improved me as a researcher, professional and person. I also want to thank my collaborator Min Xu for supporting me and providing me with great advice when I needed it.

I want to thank all the great friends that I made during my PhD which made the journey so much more fun. Andrey for the fun times we had in Edinburgh and London, Anuraag for being a great colleague at ARM, Stefanos for all the fun rides and discussions and Riyasat, Marco and Ansong for all the laughs we had at the office. Their friendship has been invaluable to me.

Next, I want to express my gratitude to all the special people in my life. My parents, Dimitris and Lalita, for supporting me unconditionally throughout my whole life and wanting the best for me. Dimitra, for being always there for me and supporting me unconditionally to all my steps for so many years. I owe to her and my parents a lot for all the support, advice and love that I have received.

Finally, I want to thank all my Greek life-long friends for all the memories and experiences that we have shared together. Each in their own way, they have influenced me to be a better person throughout all these years. I feel extremely lucky to spend quality time with them no matter how many years go by.

Declaration

I declare that this thesis was composed by myself, that the work contained herein is my own except where explicitly stated otherwise in the text, and that this work has not been submitted for any other degree or professional qualification except as specified.

(Foivos Tsimpourlas)

To my late friend.

Table of Contents

1	Introduction	1
1.1	Software Testing and Machine Learning	2
1.2	Compiler optimisation and Machine Learning	4
1.3	Contributions	5
1.4	Publications	7
1.5	Structure	8
2	Background	11
2.1	Software Testing	11
2.1.1	Terminology	11
2.1.2	Software Testing Workflow	13
2.1.3	Testing Practices	14
2.1.4	Program Instrumentation	15
2.1.5	Test Oracle	15
2.2	Compiler Infrastructure	16
2.3	Machine Learning	17
2.3.1	Neural Networks	19
2.3.2	Language Modeling	20
2.3.3	Active Learning	24
2.4	Summary	25
3	Related Work	27
3.1	Introduction	27
3.2	The Test Oracle	27
3.3	Program Representation	31
3.4	Language Modeling for Program Synthesis	34

4	Supervised learning over test executions as a test oracle	39
4.1	Introduction	39
4.1.1	Extended Contributions	42
4.2	Approach	42
4.2.1	Instrument and Gather Traces	43
4.2.2	Training Set	45
4.2.3	Preprocessing	45
4.2.4	Neural Network Model	46
4.3	Experiment	49
4.3.1	Labelling Traces	50
4.3.2	Subject Programs	51
4.3.3	Performance Measurement	54
4.3.4	Hierarchical Clustering	54
4.3.5	Results	55
4.3.6	Q1. Precision, Recall and Specificity	55
4.3.7	Q2. Size of training set	62
4.3.8	Q3. Comparison against state of art	65
4.3.9	Q4. Generalisation	65
4.3.10	Threats to Validity	67
4.4	Summary	68
5	BenchPress: A Deep Active Benchmark Generator	71
5.1	Introduction	71
5.2	Motivation	72
5.3	Approach	73
5.3.1	Learning Corpus	74
5.3.2	Language Modeling	76
5.3.3	Benchmark Generation	78
5.3.4	Feature Space Search	79
5.4	Experimental Setup	80
5.4.1	Platforms	80
5.4.2	Language Modeling for source code	81
5.4.3	Feature Spaces	81
5.4.4	Analysis of BENCHPRESS and CLgen language models	83
5.4.5	Targeted Benchmark Generation	83

5.4.6	Active Learning for Feature Selection	84
5.5	Results	85
5.5.1	Analysis of BENCHPRESS and CLGEN language models	85
5.5.2	Targeted Benchmark Generation	90
5.5.3	Active Learning for Feature Selection	96
5.6	Summary	98
6	Deep Directed Language Modeling for Compiler Features	99
6.1	Introduction	99
6.2	Approach	100
6.3	Directed Language Modeling	101
6.4	Experimental Setup	103
6.4.1	Platforms	103
6.4.2	Language Modeling for source code	103
6.4.3	Feature Spaces	104
6.4.4	Targeted Benchmark Generation	104
6.5	Results And Analysis	105
6.5.1	Targeted Benchmark Generation	105
6.6	Summary	115
7	Conclusion	117
7.1	Contributions	117
7.1.1	Automate the Test Oracle	117
7.1.2	Steerable Program Generation of Compiler Benchmarks	118
7.1.3	Directed Language Modeling for Compiler Benchmark Generation	118
7.2	Critical Analysis	119
7.2.1	Using Machine Learning as a Test Oracle	119
7.2.2	Generative Modeling for Compiler Benchmarks	119
7.2.3	Directed Program Synthesis	120
7.3	Future Work	121
7.3.1	Efficient Directed Program Synthesis	121
7.3.2	Autonomous Predictive Models	121
7.4	Concluding Remarks	122
	Bibliography	125

List of Figures

1.1	Test oracles compute the test inputs' expected outputs. They identify program correctness by comparing the actual output of a test execution with the expected output.	3
1.2	Training pipeline of a predictive model.	4
1.3	Three contributions presented in applying ML to automate software testing and compiler optimisation: (a) We propose an ML oracle, which allows developer to label only 14% of test cases with the remaining being classified automatically. (b) We present BENCHPRESS, a steerable program generator that uses active learning to improve the training data of predictive models for compilers. (c) We develop BENCHDIRECT, the first directed language model for targeted compiler benchmark generation. This contribution targets benchmarks written by compiler experts up to 36% more accurately and up to 72% faster compared to BENCHPRESS.	6
2.1	An overview of the software testing process.	12
2.2	The three-phase pipeline of a compiler.	17
2.3	A feed-forward neural network with four input features, two hidden layers and one output neuron.	19
2.4	Recurrent Neural Network architecture	21
2.5	Long Short-Term Memory cell architecture	22
2.6	The architecture of a Transformer Encoder-Decoder.	23
4.1	Key idea in our approach.	40

4.2	Gathering traces, encoding them, and using NNs to classify them. ENCODER 1 constructs a fixed vector representation per trace line. A second LSTM Encoder receives all trace line representations as a sequence and outputs a vector that summarises the execution trace. Both models are jointly trained with the MLP such that a low error is achieved in predicting the trace’s label.	43
4.3	ENCODER 1 representing a single line in a trace as a vector containing function caller, callee names, arguments and return values.	46
4.4	ENCODER 2 representing a sequence of trace lines as a single vector. .	46
4.5	Labelling test executions by matching actual and expected behavior. .	51
4.6	Precision and recall achieved by classification model over each PUT. .	63
4.7	Precision and recall achieved by classification model over each PUT. .	64
4.8	Precision-Recall curve for ETHEREUM-CD.	66
4.9	Biff trained model - Precision and recall for unseen fsms.	67
4.10	WHOIS trained model - Precision and recall for unseen fsms.	67
5.1	# Memory operations and # computational instructions for (a) Rodinia benchmarks in purple diamonds and (b) CLGEN’s samples in red dots. Generating samples with missing features is vital for predictive modeling’s performance.	73
5.2	BENCHPRESS’s high-level approach. We highlight the corpus collection and processing in green, the language modeling for source code in red and the feature space search for benchmark generation in orange. .	75
5.3	When a [HOLE] is inserted to a kernel at a random index, it hides a random number of tokens, unknown to BENCHPRESS. On this example, BENCHPRESS learns to predict the first hidden token, P	77
5.4	During sampling, BENCHPRESS receives an input and predicts iteratively the fitting tokens. BENCHPRESS predicts [ENDHOLE] to indicate a [HOLE] is complete.	78
5.5	Probability distribution of (a) token length and (b) LLVM-IR Instruction count among BENCHPRESS’s and CLGEN’s generated benchmarks. BENCHPRESS’s benchmarks presented here are generated at a single inference step without iteratively directing program synthesis.	88

5.6	PCA-2 representation of feature space coverage of BENCHPRESS and CLGEN for (a) Grewe’s et al., (b) InstCount and (c) Autophase feature spaces. In this experiment, BENCHPRESS’s generation is undirected and no iterative space search is performed.	89
5.7	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over Grewe’s et al. feature space. Relative proximity is 1 minus the distance of the two kernels in the feature space relative to the distance of the Rodinia benchmark from the axes origin. 100% means an exact match in features and is highlighted with a white asterisk (*). A score towards 0% indicates the closest match is closer to the axes origin than the benchmark, i.e., a very small or empty kernel.	91
5.8	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over InstCount feature space.	92
5.9	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over Autophase feature space.	93
5.10	# Memory operations and # computational instructions for (a) Rodinia benchmarks in purple diamonds, (b) CLGEN’s samples in red dots and BENCHPRESS’s benchmarks in green crosses after performing directed search for all Rodinia benchmarks.	95
5.11	BENCHPRESS’s performance enhancement of Grewe et al. heuristic model when using active learning compared to passively targeting random parts of the feature space over the course of 10 sampling epochs. The y-axis shows the performance enhancement as a percentage for each sampling epoch (0 to 9) shown in x-axis.	98
6.1	BENCHDIRECT’s directed language model design.	102

6.2	Pareto fronts of the average relative proximity versus total inferences for BENCHDIRECT and BENCHPRESS in targeting Rodinia benchmarks over three feature spaces ((a) Grewe’s et al., (b) InstCount and (c) Autophase). Higher relative proximity and fewer inferences are better, therefore optimal points, i.e., Pareto-dominant, are those towards the top left. We annotate the workload size configuration per Pareto point.	106
6.3	BENCHDIRECT’s acquired execution time speedup and relative proximity improvement over BENCHPRESS per workload size configuration for (a) Grewe’s et al., (b) InstCount and (c) Autophase feature spaces.	107
6.4	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on Grewe’s et al. feature space. We show the best match for BENCHDIRECT and BENCHPRESS.	109
6.5	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on InstCount feature space. We show the best match for BENCHDIRECT and BENCHPRESS.	110
6.6	Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on Autophase feature space. We show the best match for BENCHDIRECT and BENCHPRESS.	111
6.7	A comparative visualization of the	112
6.8	A comparative visualization of the	113
6.9	A comparative visualization of the	114

List of Tables

1.1	The growth of GitHub since it was founded in 2008, depicted by the total number of registered developers and total number of repositories, including both public and private repositories. As of 2022, almost 30 million out of 200 million repositories were public.	2
2.1	Software testing terms and definitions used throughout this thesis. . .	12
2.2	An simple test suite for the provided binary search example on Listing 2.1.	14
4.1	Precision, Recall and True Negative rate (TNR) using our approach and hierarchical clustering.	56
4.2	Precision (P), Recall (R) and Specificity (TNR) for each PUT omitting certain trace information.	58
4.3	Precision (P), Recall (R) and Specificity (TNR) for each PUT omitting certain trace information.	59
5.1	Throughput comparison between BENCHPRESS and CLGEN on generated OpenCL benchmarks when BENCHPRESS does not use feature-directed program generation. The column acronyms are as follows: (a) number of unique benchmarks, (b) number of compiling benchmarks, (c) rate of compilation, (d) largest compiling sample in tokens, (e) largest compiling sample in LLVM-IR instructions (-O1) (f) and inference time per sample in ms.	86
5.2	Grewe et al. heuristic model's performance, precision, recall, and specificity when trained on each technique. Speedup is the geometrical mean of speedups over all benchmarks relative to the optimal static decision, i.e. running on the GPU. Precision, recall, and specificity treat GPU labels as positive and CPU labels as negative.	97

Chapter 1

Introduction

Software has become a crucial part of modern life. The exponential growth of software in technological applications has scaled the complexity of programs significantly [165, 1]. Taking GitHub as an example, the number of developers and repositories have been increasing exponentially since 2008. In just two years (2020 to 2022), the number of developers has risen from 40 to 139 million and the number of repositories has increased from 139 to 200 million (Table 1.1).

The success of software-based products relies on fulfilling their promise to make our every day lives more productive, safer, more convenient. Whether that may be a self driving car, a robot vacuum cleaner or a streaming service on the world wide web, they all share the same success factors: Being safe and fast [5, 4]. Software correctness is a program's ability to perform the exact tasks as defined by their specification and is required for it to be established and trusted by millions of users. In privacy or real-time applications (e.g., an autopilot) it is critical indeed; a minor bug can be the cause of a devastating tragedy. Execution efficiency is equally important. An under-optimised application leads to user frustration and can cause failure on real time systems. In both these domains, compiler engineers and software testing experts have put an enormous effort in producing software that is fast and correct.

However, when the complexity of software increases, so does the effort of testing and optimising it. This is exacerbated by the recent advance of hardware with GPUs, FPGAs and heterogeneous platforms enabling the amount of applications for software to explode. Existing approaches to testing and optimising software require enormous effort and time, while also relying on compiler and low-level systems expertise. Such manual effort is no longer sustainable, leading to under-optimised and buggy software being pushed to production. To keep up with the pace of change in software appli-

	2008	2010	2012	2014	2016	2018	2020	2022
# Git Devs	100K	1M	3M	9M	14M	31M	40M	94M
# Git Repos	6.2K	1M	5M	10M	29M	100M	139M	200M

Table 1.1: The growth of GitHub since it was founded in 2008, depicted by the total number of registered developers and total number of repositories, including both public and private repositories. As of 2022, almost 30 million out of 200 million repositories were public.

cations development, researchers must develop novel techniques that enable the automation of software testing and optimisation processes. This thesis proposes methods that reduce the overhead and improve performance for software testing and compiler optimisation, two important parts of development cycle. In the next two sections, the relevant problems are stated and elaborated, existing techniques are discussed and our proposed approaches are motivated.

1.1 Software Testing and Machine Learning

Software testing is a critical part of software’s development cycle as it ensures reliability, security and high performance which further results in time saving, cost effectiveness and customer satisfaction. Testing is divided in many different categories, such as unit, functional and regression testing. All these types heavily rely on human experts developing test cases, namely programs that will exercise a PUT’s (*Program Under Test*) features.

As the scale of software increases, the number of tests needed for effective validation becomes extremely large, ultimately making software development challenging and costly [15]. To achieve cheaper and faster testing, as much of the process as possible needs to be automated. With respect to test input generation, researchers have made remarkable progress in generating effective test inputs [96, 24, 34]. Automated test input generation tools, however, generate substantially more tests than manual approaches. This becomes an issue when determining the correctness of test executions, a procedure referred to as the *test oracle*. In Figure 1.1, we show how a test oracle compares the actual and the expected output of a test execution to determine the subject program’s correctness.

The test oracle is still largely manual, relying on developer expertise. Recent sur-

veys on the test oracle problem [20, 118, 96] show that automated oracles based on formal specifications, metamorphic relations [105] and independent program versions are not widely applicable and are difficult to use in practice. The recent success of artificial intelligence in many applied tasks in computer science has motivated researchers, including the author of this thesis, to explore machine learning as a tool to address the test oracle problem. For codebases with thousands of test inputs that are automatically generated, a machine learning classifier that can automatically predict their outcome when executed would be especially useful.

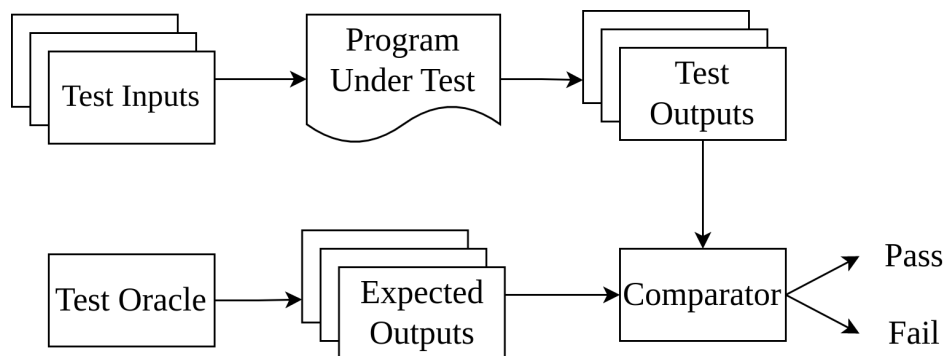


Figure 1.1: Test oracles compute the test inputs' expected outputs. They identify program correctness by comparing the actual output of a test execution with the expected output.

Previous works exploring the use of machine learning for test oracles have been in a restrictive context - applied to very small programs with primitive data types, and only considering their inputs and outputs [148, 87]. Information in execution traces has not been considered by existing ML-based approaches. Other bodies of work in program analysis have used neural networks to predict method or variable names and detect name-based bug patterns [12, 127] relying on static program information, namely, embeddings of the Abstract Syntax Tree (AST) or source code.

One of this thesis objectives is to propose machine learning based methods to identify the runtime behavior of test executions. Automating the test oracle reduces the cost of software construction and maintenance and improves overall functionality and user experience. The proposed methodology for a machine learning oracle is designed to be widely applicable to real, large-scale codebases and provide testing experts with intuitive abstractions to incorporate it into their testing pipeline.

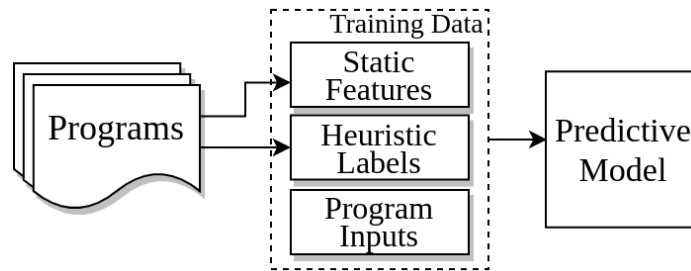


Figure 1.2: Training pipeline of a predictive model.

1.2 Compiler optimisation and Machine Learning

Program optimisation is the second area of focus in this thesis. The most essential software tool in producing efficient programs is the compiler. Compilers optimise programs through a set of transformations, from an input code to machine instructions that best utilise the resources of a target architecture. The selection of these transformations, or their parameters, is based on *heuristics* applied by compiler engineers. Each architecture supported by a compiler requires extensive manual tuning by experts to achieve great performance. As knowledge of the whole compiler is required to optimise a single heuristic, this become increasingly difficult when software complexity increases, leading to high development cost, complexity [159] and slow adaptation to the rapidly changing hardware landscape [38, 114, 89]. When compilers cannot keep up with the pace of change, the result is sub-optimal executables that consume more time and energy.

To aid the labour-intensive process of constructing optimisation heuristics, machine learning has been successfully used in several approaches [111, 157, 108, 39, 153, 158, 154, 51, 37, 99, 120, 142] in the form of predictive modeling. Predictive models predict outcomes by analyzing patterns in a given set of *features* extracted from input data as shown in Figure 1.2. For example, instead of engineers expertly crafting the loop unrolling heuristic through intuition and experimentation, a predictive model can be trained on empirical data of the performance of loops under multiple configurations. It can then be used to predict the best loop unrolling decision for any program on any hardware. Unlike manual-driven techniques, a predictive model can be easily adapted to new architectures and compilers simply by repeating the data collection for re-training. Estimating compiler optimisation heuristics through predictive modeling has been shown to outperform human experts and reduce development time [40, 42].

However, designing effective predictive models requires extensive and diverse train-

ing data to help learn accurate optimisation heuristics. In the field of compilers there is an acute shortage of benchmarks, both in quantity and diversity of features [155, 42, 40]. The average number of benchmarks used in performance tuning papers was 17 in 2017 [40, 155, 33, 69, 163, 18, 36, 55], while other areas of machine learning rely on orders of magnitude more data [45]. A shortage of benchmarks in training leads to poor feature space coverage that degrades the performance of predictive models [42, 59]. Because of this, their potential for success depends on data augmentation techniques.

The most common approach to generating programs is fuzzing [164, 103]. Fuzzers generate programs by inserting random statements and expressions that conform to a target language's standard. In the field of compilers, fuzzers [164, 103] are commonly used to produce compiler benchmarks in C and OpenCL respectively. However, fuzzed programs differ significantly from code that has been written by humans so much that predictive models perform worse when trained on their extracted features [40]. Prior research [40] has shown generative models based on deep learning are an alternative solution by generating an unbounded number of benchmarks. Synthesized benchmarks resemble programs written by humans but a recent survey [59] shows they are short, repetitive and do not extend the diversity of features of existing training benchmarks, they are therefore ineffective. There are specific areas of the feature space that are missing from our existing datasets, but no research work covers these missing features with new programs.

To improve compiler benchmarks for predictive models, we need a generative approach that identifies and targets those features that are missing from existing datasets. This thesis aims to bridge the gap between the achieved and potential performance of predictive modeling for compiler heuristics by proposing novel techniques to generate high quality compiler benchmarks at scale.

1.3 Contributions

Many software testing and compiler optimisation tasks that are vital for the software industry's success still remain largely manual. Inspired by the ever increasing performance of machine learning applications in many tasks, we claim they have the potential to be faster and better in performing tasks that have been long considered achievable only by experts. This thesis tackles this issue and proposes three novel contributions that reduce development cost and take away the need for domain expertise, human

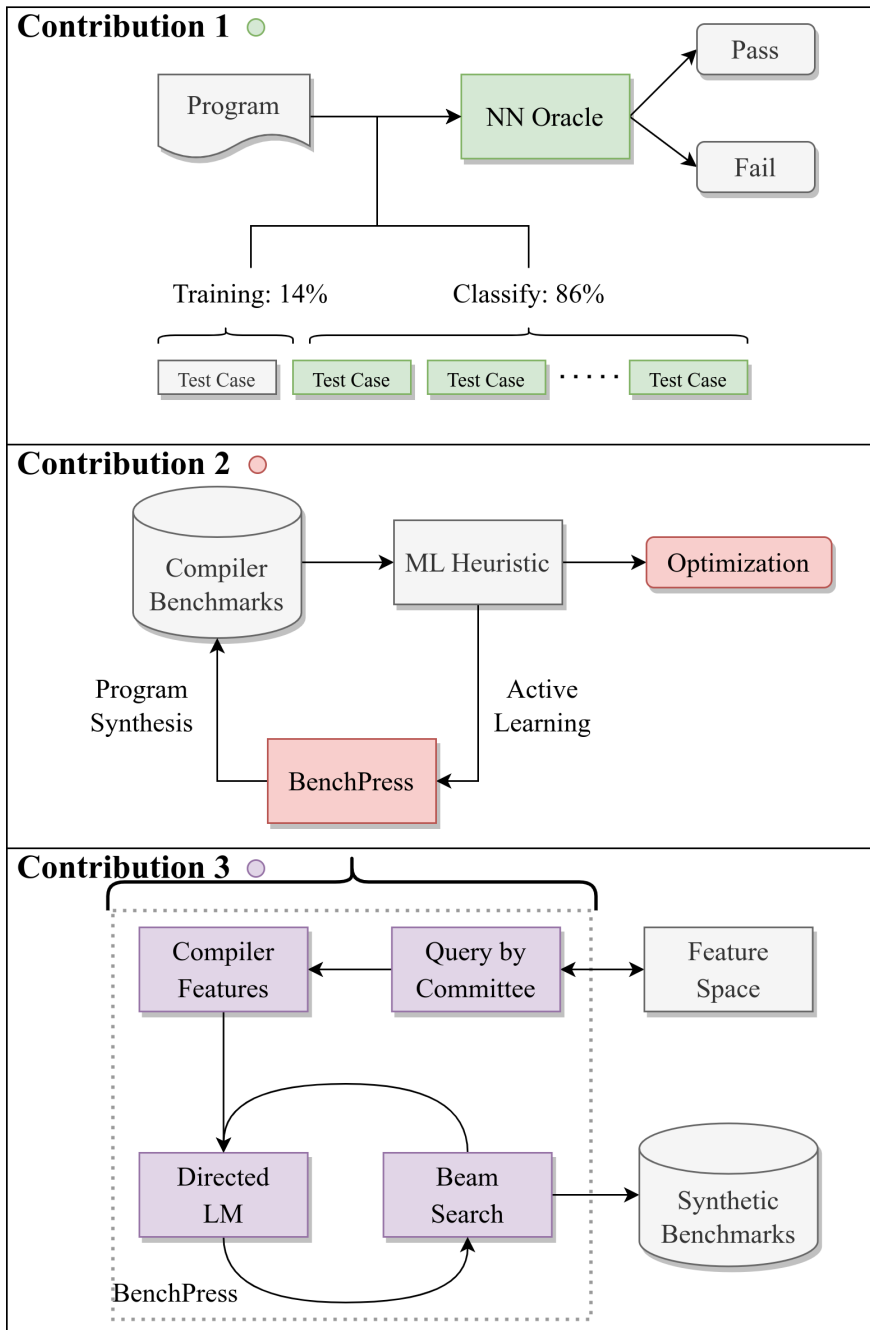


Figure 1.3: Three contributions presented in applying ML to automate software testing and compiler optimisation: (a) We propose an ML oracle, which allows developer to label only 14% of test cases with the remaining being classified automatically. (b) We present BENCHPRESS, a steerable program generator that uses active learning to improve the training data of predictive models for compilers. (c) We develop BENCHDIRECT, the first directed language model for targeted compiler benchmark generation. This contribution targets benchmarks written by compiler experts up to 36% more accurately and up to 72% faster compared to BENCHPRESS.

intuition and manual effort. One contribution lies in the field of software testing and two in compiler optimisation. The overview of this thesis’s contributions are shown in Figure 1.3.

The key contributions of this thesis are the following:

- The first application of deep learning over program execution traces to predict runtime correctness. We develop the first NN-based test oracle that classify test executions as “pass” or “fail”. Our model achieves a near maximum classification accuracy on 15 real, large-scale codebases by training only on 14% of the initial labelled data. This addresses the manual effort of labelling manually the expected output of test executions for large scale production code.
- BENCHPRESS, the first steerable deep learning program synthesizer to generate compilable, executable benchmarks for compilers. BENCHPRESS produces human-like programs that compile at a rate of 87% in contrast to 2.33% achieved by the current state of the art generator. BENCHPRESS synthesizes benchmarks that target parts of the feature space previously unreachable by human-written code from GITHUB. This contribution addresses the scarcity of high-quality training data for compiler predictive models.
- BENCHDIRECT, the first directed language model for compilers. We develop a language representation model that synthesizes compiler benchmarks by jointly conditioning on source code context and the desired compiler features in any feature space. BENCHDIRECT outperforms BENCHPRESS in the task of directing program generation towards the features of human-written benchmarks by targeting them up to 36% more accurately and up to 72% faster. This contribution addresses the inefficiency of existing directed program generation techniques, enabling large scale tasks.

1.4 Publications

This thesis consists of four publications describing our research ideas and results.

Chapter 4 elaborates on our approach to classify program executions, published in:

- “*Supervised learning over test executions as a test oracle*”,
F. Tsimpourlas, M. Allamanis, A. Rajan, SACSE 2021 [145]).

- “*Embedding and classifying test execution traces using neural networks*”, F. Tsimpourlas, G. Rooijackers, A. Rajan, M. Allamanis, IET Software 2022 [146].

Chapter 5 presents BENCHPRESS, the first guided program generator for compiler benchmarks, based on language modeling and active learning:

- “*BenchPress: A Deep Active Benchmark Generator*”, F. Tsimpourlas, P. Petoumenos, M. Xu, C. Cummins, K. Hazelwood, A. Rajan, H. Leather, PACT 2022 [144].

Chapter 6 discusses BENCHDIRECT, a directed language model for efficient and accurate steerable compiler benchmark generation:

- “*BenchDirect: A Directed Language Model for Compiler Benchmarks*”, F. Tsimpourlas, P. Petoumenos, M. Xu, C. Cummins, K. Hazelwood, A. Rajan, H. Leather, Submitted to TACO, March 2023 [143].

The source code and experimental data for all three contributions are publicly available on GitHub for other researchers to use on the following repositories:

1. <https://github.com/fivosts/Learning-over-test-executions>
2. <https://github.com/fivosts/BenchPress>

1.5 Structure

This thesis is organized as follows:

Chapter 2 provides background. The relevant terminology is defined and all techniques used in this work are described.

Chapter 3 offers a review of existing literature, divided into three categories: program testing, program optimisation and language modeling for program generation.

Chapter 4 presents a novel approach for solving the test oracle problem. A qualitative evaluation shows the effectiveness of machine learning identifying incorrect program executions with near maximum accuracy.

Chapter 5 introduces BENCHPRESS, a novel directed program synthesizer that finds important program features with active learning and generates programs with such features. BENCHPRESS outperforms the state of the art in multiple program generation tasks and improves the performance of predictive models for compilers.

Chapter 6 presents **BENCHDIRECT**, an optimised steerable program synthesizer based on a novel, feature-conditioned language model. **BENCHDIRECT** outperforms **BENCHPRESS** in targeting compiler features in all three accuracy, speed and code quality measured as its human-likeness.

Chapter 7 summarizes the overall findings of the thesis, provides a review of current limitations and outlines potential future directions.

Chapter 2

Background

In this Chapter, we provide necessary background information on various aspects and concepts of software testing, compiler optimisation and machine learning that are relevant to the scope of this thesis.

Section 2.1 describes the software testing workflow and concepts involved in the testing process. Section 2.2 illustrates the fundamentals of compilers and software optimisation techniques. Section 2.3 illustrates modern machine learning methodologies that are used in all thesis contributions in software testing and compiler optimisation.

2.1 Software Testing

Software testing is the procedure that examines the behavior of subject programs, known as *Programs Under Test* (PUT). The primary purpose of software testing is to expose failures so they can be discovered and corrected before programs are released to users [92]. Testing's scope includes the examination as well as the execution of code in various conditions to determine correctness. To validate programs through execution, a set of test inputs are needed, also known as *test cases*. An overview of the software testing process is shown in Figure 2.1.

2.1.1 Terminology

First, in Table 2.1 all the definitions and terms used in this thesis are provided to avoid confusion and be consistent with the current literature.

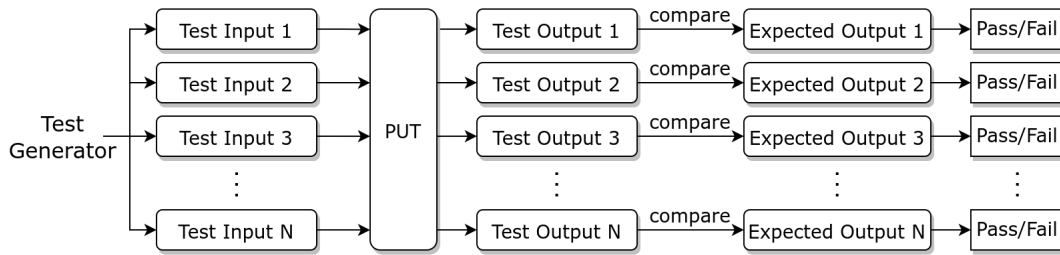


Figure 2.1: An overview of the software testing process.

Definition	Description
System/Program Under Test (SUT/PUT)	Software or System that is tested.
Test Oracle	a mechanism that determines whether a test execution is correct or not.
Software Specification	A description of a PUT's intended behaviour and the foundation of the test oracle.
Test Input	The input used to execute a PUT.
Test Output	The real output data collected from a PUT's execution given a specific test input.
Expected Output	The expected output data from the execution of a specific test input.
Test Case	The test input with its expected output.
Test Suite	A collection of test cases.
Test Failure	A state where a PUT raises an exception or an undefined behavior leading to not producing a test output.
Test Result	The status indicating whether a test has passed or failed.

Table 2.1: Software testing terms and definitions used throughout this thesis.

2.1.2 Software Testing Workflow

We illustrate the process of software testing with an example. Listing 2.1 presents a binary search function implemented in Python. The function has four inputs: a sorted array of integers *array*, two numbers *low* and *high* indicating the bounds of the array, and a number *target*, for which to search in the array. It performs a binary search in the array and returns the first index at which the target is found.

```
1 def bin_search(array, low, high, target):
2     if high >= low:
3         mid = low + (high - low) // 2
4         if array[mid] == target:
5             return mid
6         elif array[mid] > target:
7             return bin_search(array, low, mid - 1, target)
8         else:
9             return bin_search(array, mid + 1, high, target)
10    else:
11        return -1
```

Listing 2.1: An example Python program performing binary search on sorted integer array

The functional specification for this binary search is the following:

1. It accepts a sorted array of integer values, two integers to index the array and one integer to search the array for.
2. It returns the first index at which the target is present in the array. If the number is not present, it returns -1.
3. If the array is empty, it returns -1.

Table 2.2 shows a test suite for this program based on this specification which contains four tests. Listing 2.2 presents an implementation for these tests. In this case the PUT is the `bin_search()` function. The test inputs are declared and initialized, the PUT is executed with the test inputs and the test output is recorded. The test output is compared to the expected output using assertions.

Test ID	<i>array</i>	<i>low</i>	<i>high</i>	<i>target</i>	Expected Output
1	{}	0	0	2	-1
2	{1, 2, 2, 4}	0	3	2	1
3	{1, 6, 7}	0	2	7	2
4	{1, 6, 7}	0	2	8	-1

Table 2.2: An simple test suite for the provided binary search example on Listing 2.1.

```

1 def test1():
2     assert bin_search([], 0, 0, 2) == -1
3
4 def test2():
5     assert bin_search([1, 2, 2, 4], 0, 3, 2) == -1
6
7 def test3():
8     assert bin_search([1, 6, 7], 0, 2, 7) == 2
9
10 def test4():
11     assert bin_search([1, 6, 7], 0, 2, 8) == -1

```

Listing 2.2: A simple implementation for the example test suite given

2.1.3 Testing Practices

The purpose of software testing is to provide confidence that a system's behaviour conforms to specification. Ideally, test suites will exhaustively sample the entire input space of a system, but in practice this is infeasible because it would require billions of tests even for trivial programs. However, the goal of software testing is to uncover faults in the system which can be fixed, therefore sufficient testing should eventually provide enough confidence that no critical faults remain.

Testing can be applied at any level of system granularity - individual functions, modules and entire systems. *Unit testing* checks the behaviour of the smallest functional units of a program. *Integration testing* checks the correctness of the interactions between units and modules. *System testing* focuses on the behaviour of the system as a whole. *Regression testing* aims to ensure that changes to the existing codebase do not introduce faults in the system. *Mutation testing* is a type of testing that injects

errors by modifying the source code of the PUT with random changes. The PUT is executed with a test suite and it is checked whether the test cases are able to expose these injected alterations [85].

2.1.4 Program Instrumentation

Program instrumentation is used to measure performance, discover errors and collect runtime information in the form of execution traces by inserting extra code to a program under test [109]. The inserted code must not alter the original functionality of the program and must also be used to monitor certain metrics from a program's execution.

The most common use cases of program instrumentation include:

1. **Profiling** [138] measures dynamic program behaviour through execution. Profiling information helps developers analyse dynamic information that cannot be measured with static analysis.
2. **Performance estimation** [156] uses timers to code segments that are computationally extensive. Timing information is used to estimate overhead and reveal the time consuming sections of a program.
3. **Logging execution information** [80] involves recording events related to a program's runtime including crashes, code coverage or control flow. Execution information is used to measure the achieved test effectiveness of test executions.

There is a plethora of instrumentation tools available, written in different languages [141, 31, 81]. In this thesis, we use LibTooling [106] to apply static analysis on C, C++ and OpenCL programs for feature extraction. More details on this approach will be given in Chapter 5. For runtime instrumentation, we use the LLVM [97] framework to record the control and data flow of test executions and collect data on a program's runtime behaviour. The details of this technique is discussed in Chapter 4.

2.1.5 Test Oracle

Test oracles [82] belong in the family of black-box testing techniques. Black-box testing examines the functionality of an application without peering into its internal structures [57, 84] and can be applied to any type of testing, unit, functional or system. Test oracles determine the correctness of a program by comparing the test output of a PUT for a given test input with its expected output. Determining the correct output given a

test input is known as *the oracle problem*, a difficult challenge that involves working with the controllability and observability of a system. This process still remains largely manual.

There exist several categories of test oracles found on the literature. The most common forms are the following:

1. **Specified Oracle** judges a program under test's correctness based on formal specifications.
2. **Implicit Oracle** relies on implied information and assumptions such as conclusions based on program crashes.
3. **Derived Oracle** uses code documentation or system executions when specified oracles are unavailable.
4. **Human Oracle**. When no other oracle can be used human experts estimate a program's expected behaviour.
5. **ML Oracle**. Statistical methods based on machine learning that estimate a program's test execution correctness.

This thesis focuses in machine learning oracles. The first contribution described is a machine learning-based technique to automate the test oracle by classifying test executions as “pass” or “fail”. This approach is presented in Chapter 4.

2.2 Compiler Infrastructure

The compiler is a programming tool that translates programs from a given source language to a lower-level target language. Throughout the compilation process, program semantics are preserved. Compilers are expected to produce a good quality representation of programs in the target language, being optimal with respect to objective functions. An important objective function is execution time, i.e., the optimisation goal is to produce a target language representation of the program that will execute as fast as possible.

Compilers are usually presented as three stage architecture, as shown in Figure 2.2. These stages are the front end, the middle end (or optimiser) and the back end. The front end parses and validates the source code, making sure it conforms to the

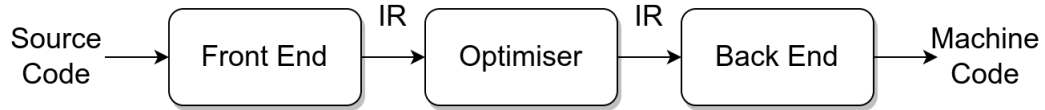


Figure 2.2: The three-phase pipeline of a compiler.

source language’s grammar. This parsed code is translated into an *Intermediate Representation (IR)*. The compiler’s optimiser performs a broad variety of transformations to the IR, called *optimisations*. The compiler’s optimisations are usually independent of the source language or the target hardware and their goal is to improve the code’s performance. The back end is the compiler’s *code generator* and translates the IR to the target language. Hardware-related optimisations can also take place during code generation to exploit the architecture’s supported features.

One of the most commonly used compilers is the LLVM [97]. LLVM’s intermediate representation, *LLVM-IR*, is a human-readable, assembly-like representation that is used by the compiler’s middle-end to apply optimisations. LLVM-IR’s most notable feature is *Static Single Assignment (SSA)* form. In SSA form, each variable is assigned to only once, and every reference to a variable is a reference to its single assignment. SSA form has a number of useful properties that make it useful for optimising code. For example, it makes it easy to determine the live ranges of variables, which can be used for register allocation. It also makes it easy to detect when variables are used before they are defined, which can help expose errors in the code.

2.3 Machine Learning

Machine learning is a family of statistical methods and algorithms. They can be predictive to make predictions in the future, or descriptive to gain knowledge from data. ML algorithms are broadly classified into four major categories, depending on the nature of the learning response available to a learning system:

1. **Supervised Learning** models learn by labelled input/output examples.
2. **Unsupervised Learning** models learn by similarities on unlabelled/unstructured data.
3. **Reinforcement Learning** models learn by trial and error to maximize their accumulated reward based on a reward function.
4. **Semi-supervised Learning** models learn on partly-labelled or incomplete data.

Predictive models make predictions by correlating their input variables, known as *features*, with their outputs, or *labels*. The n -dimensional space described by n input features is called a *feature space*. A *feature vector* is the set of features describing a single point in the feature space. The deduction of features from raw data is called *feature selection* and it is an important process in machine learning architecture design, which significantly determines the model's performance. Instead of manually selecting features, ML architectures can also be used by learning a numerical vector representation of raw data in the latent space. This process is commonly referred to as *summarization*. The learnt features they extract from raw data are called *embeddings* [29]. Trained embeddings have been shown to significantly outperform manual feature selection in their quality to represent input data [112].

For classification tasks, there are multiple ways to measure a machine learner's performance in predicting labels. All metrics make use of the following scores:

- **True Positives (TP):** The number of positive class items correctly labelled.
- **False Positives (FP):** The number of negative class items labelled as belonging to the positive class.
- **True Negatives (TN):** The number of negative class items correctly labelled.
- **False Negatives (FN):** The number of positive class items labelled as belonging to the negative class.

To evaluate machine learning models in this thesis, we make use of *precision*, *recall* and *specificity*. These three metrics are defined as follows:

$$Precision = \frac{TP}{TP + FP} \quad (2.1)$$

$$Recall = \frac{TP}{TP + FN} \quad (2.2)$$

$$Specificity = \frac{TN}{TN + FP} \quad (2.3)$$

2.3.1 Neural Networks

Artificial Neural Networks (ANNs) are the most successful and commonly used type of machine learning architecture and are divided into many different categories. The simplest one is *feed-forward neural networks* (FNN) [21]. FNNs consist of a sequence of neuron layers, where every neuron of a previous layers forms a weighted connection with all neurons of the next layer. This way, the flow of information is unidirectional, mapping the input variables to the neural network's outputs, as shown in Figure 2.3.

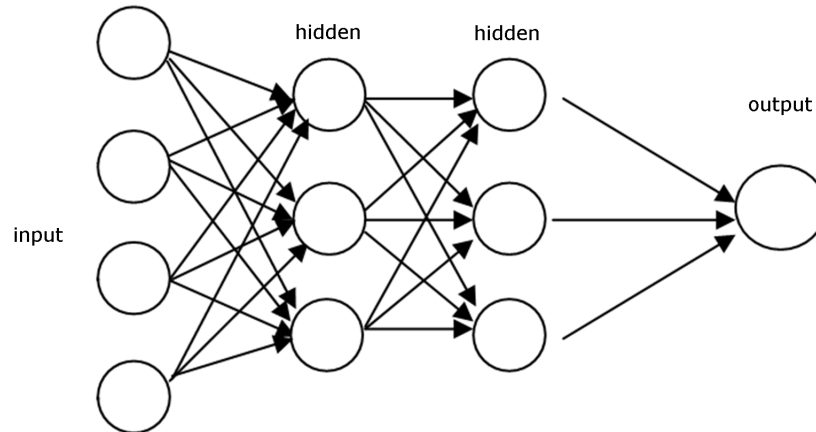


Figure 2.3: A feed-forward neural network with four input features, two hidden layers and one output neuron.

ANNs consist of the *input layer* that holds as many neurons as the feature space and the *output layer* that consists of one neuron per predicted label. For regression, only one neuron is needed. For classification, there is one neuron per label in the label space. All intermediate layers for which there are no ground truth values are known as *hidden layers*. Feed-forward neural networks are powerful in approximating functions. In theory, they can learn any bounded continuous function to arbitrary precision given the right amount of neurons [115].

The prevailing method of training a neural network is back-propagation. A *batch* of B observations is propagated through the network in a *forward* pass, from the inputs to the outputs. The final output vector of the network \hat{y} are compared against the labels y and the distance between them is computed, $L(\hat{y}, y)$. The error metric, or *loss function*, depends on the task. A common loss function used for the task of classification is *categorical cross entropy* [116]. Neural networks are prone to *over-fitting*, when the parameters of the model lose the ability to generalise to unseen data by becoming too specialized on the training data. Many *regularisation* techniques have been adopted

to mitigate the risk of over-fitting. One of the most successful is *Dropout*, in which a parameter in the range $[0, 1]$ is used to determine a proportion of artificial neurons to be removed. This helps training by preventing complex co-adaptations on training values [137].

2.3.2 Language Modeling

Language models are a descriptive type of neural networks that gain knowledge from natural or programming languages used by humans to communicate. The contributions of this thesis heavily rely on the use of language models to represent, embed and generate human-readable programming languages including C, C++ and OpenCL.

Modeling languages is a challenging task that requires large amount of data and large models that are computationally expensive. The high-level steps for designing a NN-based language model architecture for source code are the following:

1. **Tokenize:** Separate source code into words.
2. **Encode:** Convert string-formatted words into numerical vectors.
3. **Train:** Apply language modeling task on encoded corpus.

For LM tasks, there are certain architectures that have been shown to capture more effectively than FNNs the relationship of words in a sentence. The most common of them is the *Recurrent Neural Network* (RNN) [133]. The RNN is a type of ANN where the connections between neurons form a feedback loop from successive layers back to predecesing ones. This enables processing arbitrarily large sequences of tokens while learning the correlations between neighbouring words in the sentence. The RNN cell's hidden state at each time step t is calculated based on the input and previous hidden state h_{t-1} , using the weight matrices W and bias vector b to combine the inputs and apply the activation function:

$$h_t = f(W_{xh}x_t + W_{hh}h_{t-1} + b_h) \quad (2.4)$$

Where, h_t represents the hidden state at time step t , x_t represents the input at time step t , and f represents the activation function. W_{xh} and W_{hh} are the weight matrices that connect the input and hidden states, and b_h is the bias vector.

RNNs can be trained with ordinary back-propagation by unfolding the computation graph over time as shown on Figure 2.4. Back-propagating through time allows the

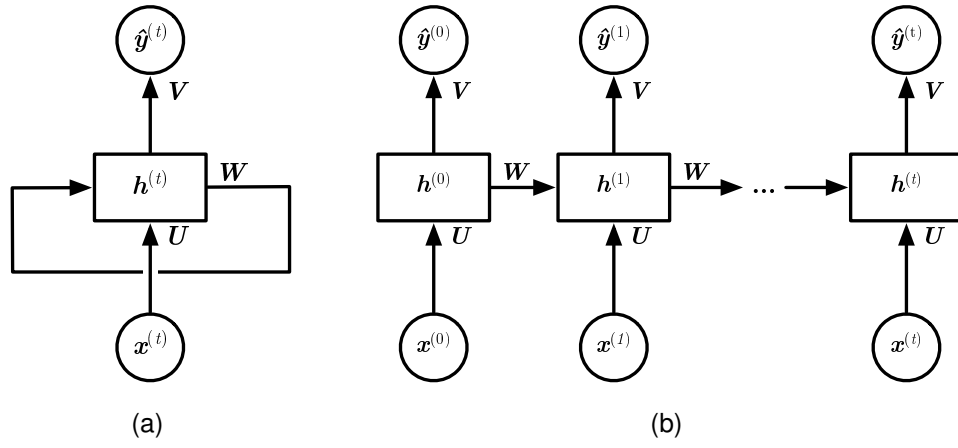


Figure 2.4: The computational graph of a RNN, shown as a recurrence relation in a, and unfolded in b. $\mathbf{x}^{(t)}$ is the input, $\mathbf{h}^{(t)}$ is the hidden state, and $\hat{\mathbf{y}}^{(t)}$ is the output. The network comprises of three weight matrices: inputs-to-hidden weights \mathbf{U} , hidden-to-hidden weights \mathbf{W} , and hidden-to-output weights \mathbf{V} .

propagation of error in the temporal domain in the same manner as through layers. RNN's accuracy suffers when they are tasked to learn correlations over long sequences. This is caused by the exponential diminishing or increase of gradients as they are propagated through the activation functions by the recurrence relation. This issue is also known as the *vanishing gradients problem* [76].

A bolstered, RNN-based architecture addresses this issue, namely the *Long Short-Term Memory* (LSTM) [77]. The LSTM inherits the RNN design with the addition of a cell that stores information and three gates which control the flow of it into and out of the cell, shown in Figure 2.5. For a time-step t , an input at this time-step x_t , the hidden-state h_t , the memory's cell-state c_t and input, forget and output gates i_t , f_t , o_t respectively, the output prediction from the LSTM for the next time-step is as follows:

$$i_t = \sigma(W_{xi}x_t + W_{hi}h_{t-1} + W_{ci}c_{t-1} + b_i) \quad (2.5)$$

$$f_t = \sigma(W_{xf}x_t + W_{hf}h_{t-1} + W_{cf}c_{t-1} + b_f) \quad (2.6)$$

$$c_t = f_t \cdot c_{t-1} + i_t \cdot \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \quad (2.7)$$

$$o_t = \sigma(W_{xo}x_t + W_{ho}h_{t-1} + W_{co}c_t + b_o) \quad (2.8)$$

$$h_t = o_t \cdot \tanh(c_t) \quad (2.9)$$

Where:

σ is the sigmoid activation function [66] and W , b are the weight matrices and bias vectors, respectively, that are learnt during training. This formula represents the

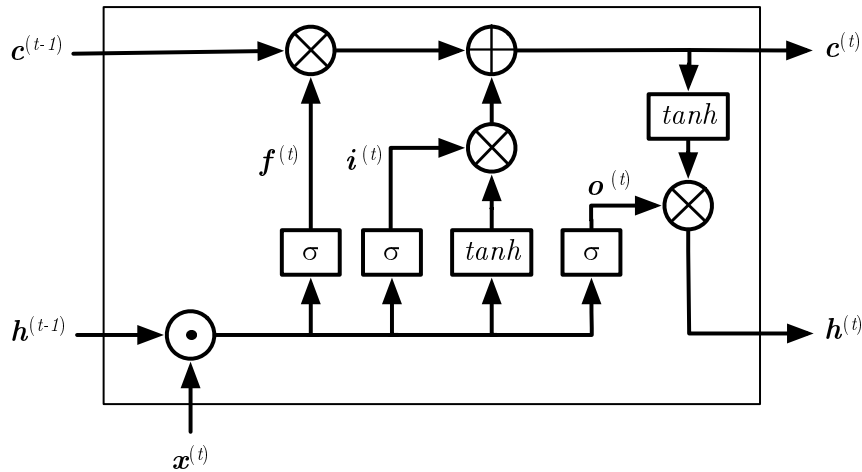


Figure 2.5: A Long Short-Term Memory cell. Input $x^{(t)}$ is concatenated with prior hidden state $h^{(t-1)}$ and used with prior cell state $c^{(t-1)}$ to compute the next step's hidden state $h^{(t)}$ and cell state $c^{(t)}$.

computations performed by one LSTM unit at each time step. In an LSTM network, multiple LSTM units are typically stacked together, with the hidden state of one unit being used as the input to the next unit. The final hidden state summarizes the state of the whole time series and is used as an input to the downstream task, whether that be classification or regression.

The LSTM has been firmly established as the state of the art in sequence modeling and transduction problems. However, its approach to processing words within a sentence in a sequential way precludes parallelization within a training example, meaning words cannot be computed in parallel because of their temporal dependency. This constraint poses a severe performance bottleneck when it comes to designing large language models with billions of training data.

The *Transformer* [149] addresses this issue and enables the creation of language models with hundreds of billions of parameters. The Transformer is based on the *attention* mechanism [56, 136, 32]. Its input is first embedded into a continuous representation and then processed by multiple layers of self-attention. The self-attention layers allow the model to attend to different positions of the input sequence simultaneously and weigh their importance when computing the output:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V \quad (2.10)$$

Where, Q , K , and V are the query, key, and value matrices, respectively, and d_k is

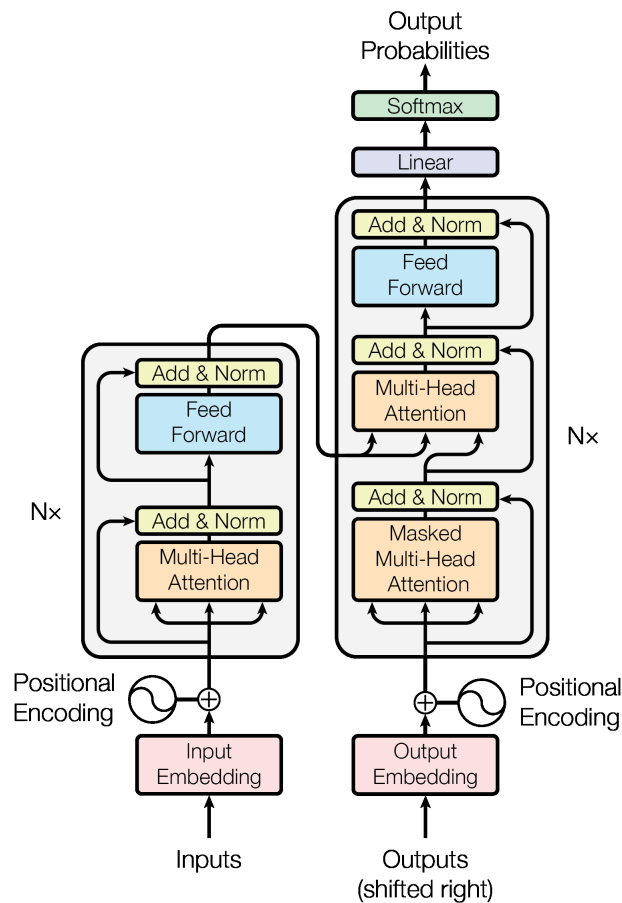


Figure 2.6: The architecture of a Transformer Encoder-Decoder.

the dimension of the keys. The output of the self-attention function is a weighted sum of the values, where the weights are computed based on the dot product between the query and the keys. Then this output is processed in parallel by a feed-forward layer followed by a ReLU activation [2]. The Transformer's architecture is shown on Figure 2.6.

The Transformer is applied to numerous tasks in language modeling, generation and predictive modeling. It also has been used as a foundation for many derived architectures. The most notable is *BERT* (Bidirectional Encoder Representations from Transformers) [46]. BERT learns deep bidirectional representations by jointly conditioning on both left and right context of an input. One of BERT's key features is its ability to preserve the context of the words in a sentence, rather than just processing individual words in isolation. This allows capturing the relationship between words in a more accurate way compared to traditional NLP models, leading to excellent performance on a range of NLP tasks [46]. BERT has been widely adopted as the state-of-the-art in language translation, text summarization, and sentiment analysis

tasks [95, 52, 91]. It has also been used to improve the performance of other machine learning models, such as those used for image classification and object detection [101]. The second and third contribution of this thesis extends BERT from a masked language model to a pre-trained generative model for programming languages.

2.3.3 Active Learning

This thesis applies active learning to search feature spaces applied to the field of compilers for areas with few or no representative samples. Active learning is a type of machine learning that involves learning from actively selected, informative examples rather than passively accepting a pre-determined dataset. In active learning, the learning algorithm can interactively query the user (or some other information source) to obtain the desired outputs, and then use these to improve the accuracy of the model.

Active learning is useful in situations where there is a limited amount of labeled data available, or where the cost of labeling data is high. By allowing the learning algorithm to select which examples to label, active learning can improve the efficiency of the learning process and result in better performance compared to using a fixed dataset. There are several strategies for selecting the most informative examples in active learning. In this thesis, we explore two:

- **Query by Committee**, an approach to selective sampling in which disagreement among a committee of oracles is used to select data for labeling [135].
- **Expected Error Reduction**, a technique that selects a datapoint based on its estimated impact on the classifier's future error over all other datapoints [131].

Query by Committee (QbC) is a simple yet powerful algorithm. During training, a committee of students is trained on the same set of data. The next selected datapoint for labeling (query) is chosen according to the principle of *maximal disagreement*. This disagreement is typically measured with *entropy*, a statistical metric which indicates the amount of uncertainty in a set of samples. Given an input datapoint x that is candidate for labelling, the entropy among N committee members' predictions is defined as:

$$H(x) = - \sum_{i=1}^N p(o_i^x) \log_2 p(o_i^x), \quad o_l \in L := \{o_0, o_1, \dots, o_L\} \quad (2.11)$$

Where l is one possible label that belongs in the label space L and $p(o_i^l)$ is the likelihood of the predicted label l by member i . Query by Committee has shown it

is an effective active learning algorithm for relatively simple tasks where the label space has only a few dimensions. For more complicated tasks, QbC suffers from poor accuracy with noisy predictions (high entropy) among members' predictions. Also, members of the same architecture (e.g. NNs) tend to produce the same predictions regardless of their number of parameters, initialization of weights and the worthiness of a query. This also hurts the committee's overall accuracy.

Expected Error Reduction (EER) is a more advanced technique that tackles QbC's shortfall in complex tasks. In EER, a single predictive model is directly prompted to query a datapoint among a provided sample set. The point that is estimated to cause the lowest possible aggregated error on the training dataset is selected. This is measured by re-training the predictive model on one pair (\mathbf{x}, \mathbf{y}) at a time, where \mathbf{x} is a datapoint in the sample set and \mathbf{y} is a label in the label space, and then measuring the aggregated error over the training set with the categorical cross entropy error function. This process is repeated for all sample datapoints \mathbf{x} combined with each possible label \mathbf{y} in the label space. EER has shown to learn faster than QbC the classification task, requiring less labelled data. However, to query a single datapoint in the feature space, it is required to re-train the predictive model $N \times L$ times, where N is the number of datapoints in the sample set and L is the size of the label space. This is a significant overhead which is prohibitive for large predictive models or multi-dimensional feature spaces.

2.4 Summary

In this Chapter, all necessary background related to the problem and the contributions presented in this thesis is presented. The reader is introduced to the process of software testing and its categories in Section 2.1. Section 2.2 elaborates on modern compiler architectures and how they are used to optimize programs. Finally, Section 2.3 provides an overview of all machine learning tools and approaches that were used in this thesis. The following Chapter surveys the existing related literature.

Chapter 3

Related Work

3.1 Introduction

This Chapter surveys the existing literature relevant to the contributions of this thesis. Section 3.2 presents the existing literature in the field of software testing and runtime analysis, focused on the test oracle. Section 3.3 reviews current research in program embedding and representation relevant to our test oracle approach. In Section 3.4, a literature review of benchmark generation for compilers is presented, to provide intuition about our second and third contributions.

3.2 The Test Oracle

Test oracles are an important part of the software testing process as they are used as a baseline to compare the actual output of a test input to the expected output and determine whether the test execution has passed or failed.

The importance of oracles as an integral part of the testing process has been a key topic of research for over three decades. We distinguish three different kinds of test oracles, based on the survey by [20]. The most common form of test oracle is a *specified oracle*, one that judges behavioural aspects of the system under test with respect to formal specifications. Although formal specifications are effective in identifying failures, defining and maintaining such specifications is expensive and also relatively rare in practice. *Implicit* test oracles require no domain knowledge and are easy to obtain at no cost. However, they are limited in their scope as they are only able to reveal particular anomalies like buffer overflows, segmentation faults, deadlocks. *Derived* test oracles use documentations or system executions, to judge a system's behaviour,

when specified test oracles are unavailable. However, derived test oracles, like metamorphic relations and inferring invariants, is either not automated or it are inaccurate and irrelevant making it challenging to use.

For many systems and much of testing as currently practised in industry, the tester does not have the luxury of formal specifications or assertions or even automated partial oracles [70, 71]. Predicting the behavior of a large scale program is a tedious task; for this reason researchers have tried several methods to automate this process. Statistical analysis and machine learning techniques provide a useful alternative for understanding software behaviour using data gathered from a large set of test executions. Briand et al. [27] present a comprehensive overview of existing techniques that apply machine learning for addressing testing challenges. Among these, the closest related work is that of Bowring et al. [25]. They propose an active learning approach to build a classifier of program behaviours using a frequency profile of single events in the execution trace. Evaluation of their approach was conducted over one small program whose specific structure was well suited to their technique.

More recently, Rigger et al. [130] introduce *Intramorphic Testing*, a white-box methodology for the test oracle that comes in contrast to common black-box approaches such as differential and metamorphic testing. Rigger's et al. approach includes changing one component of the SUT in a known way so that the SUT's original output can be related to the changed output. They illustrate how Intramorphic Testing can expose bugs with three example programs: (a) AST printing, (b) Monte Carlo simulations and (c) the Knapsack problem [93]. Intramorphic Testing is a novel approach to the test oracle problem, however it poses several limitations. First, a developer must not only know the whole codebase and how its components interact, but also how a small alteration in one component changes the program's output. This is very difficult on industry-scale codebases, where the tester is often not the developer. Second, constructing relations between program versions for simple operations might be relatively easy but such relations in complex software is near-impossible. To make things worse, many mutations may have no effect to the SUT's output thus leading to no intramorphic relations. These two factors increase exponentially the amount of effort and time needed to construct high-quality relations in large-scale software and make use of this technique.

Machine learning techniques have also been used in fault detection. Brun and Ernst [28], explored the use of support vector machines and decision trees to rank program properties, provided by the user, that are likely to indicate errors in the program.

Podgurski et al. [126] use clustering over function call profiles to determine which failure reports are likely to be manifestations of an underlying error. A training step determines which features are of interest by evaluating those that enable a model to distinguish failures from non-failures. The technique does not consider passing runs. In their experiments, most clusters contain failures resulting from a single error.

Almaghairbe et al. [11] propose an unsupervised learning technique to classify unlabelled execution traces of simple programs. They gather two kinds of execution traces, one with only inputs and outputs, and another that includes the sequence of method entry and exit points, with only method names. Arguments and return values are not used. They use agglomerative hierarchical clustering algorithms to build an automated test oracle, assuming passing traces are grouped into large, dense clusters and failing traces into many small clusters. They evaluate their technique on 3 programs from the SIR repository [49]. The proposed approach has several limitations. They only support programs with strings as inputs. They do not consider correct classification of passing traces. The accuracy achieved by the technique is not high, classifying approximately 60% of the failures. Additionally, fraction of outputs that need to be examined by the developer is around 40% of the total tests.

Almaghairbe et al. [11] assumed that all the passing traces present the same behaviour, leading to them being organized in one, large cluster. On the other hand, according to them, failing traces tend to present non-uniform, wide-ranging patterns, which results in failing traces being spread among many small clusters. According to their methodology, the clusters that are sized below the average of the set, are considered to contain failing traces and clusters sized above the average are considered to contain passing traces. For their evaluation, they used different clustering techniques and experimented on multiple cluster sizes. Almaghairbe et al. used ‘Daikon’ for instrumentation, an invariant detector that uses machine learning to observe program values and summarize them into a set of formulas.

Existing work using execution traces for bug detection has primarily been based on clustering techniques. Neural networks, especially with deep learning, have been successful for complex classification problems in other domains like natural language processing, speech recognition, computer vision. There is limited work exploring their benefits for software testing problems.

NNs were first used by Vanmali et al. [148] to simulate behaviour of simple programs using their previous version, and applied this model to regression testing of unchanged functionalities. Aggarwal et al. [3] and Jin et al. [87] apply the same ap-

proach to test a triangle classification program, that computes the relationship among three edge inputs to determine the type of triangle. Mytkowicz [48] et al. propose TOGA, a Transformer-based approach to infer exceptional and assertion bug finding test oracles. Based on the observation that oracles in developer-written unit tests typically follow a small number of common patterns, Mytkowicz et al. define a grammar that expresses a taxonomy of these patterns. Using this grammar, a two-step neural ranking procedure scores candidate oracles. They evaluate TOGA on test oracle inference, reporting an accuracy improvement of 11% and 33% over related work in two oracle inference datasets. Their tool is also evaluated in bug finding, exposing 57 real world bugs in Java benchmark, DEFECTS4J [90]. TOGA outperforms the comparing oracles in exposing more real bugs and developing a lower *False Positive Rate* of 25%.

The few existing approaches using NNs have been applied to simple programs having small I/O domains. The following challenges have not been addressed in existing work,

1. Training with test execution data and their vector representation – Existing work only considers program inputs and outputs that are of primitive data types (integers, doubles, characters). Test data for real programs often use complex data structures and data types defined in libraries. There is a need for techniques that encode such data. In addition, existing work has not attempted to use program execution information in NNs to classify tests. Achieving this will require novel techniques for encoding execution traces and designing a NN that can learn from them.
2. Test oracles for industrial case studies - Realistic programs with complex behaviours and input data structures has not been previously explored.
3. Effort for generating labelled training data - Training data in existing work has been over simple programs, like the triangle classification program, where labelling the tests was straightforward. Availability of labelled data that includes failing tests has not been previously discussed. Additionally, the proportion of labelled data needed for training and its effect on model prediction accuracy has not been systematically explored.

The performance of neural networks as classifiers was boosted with the birth of deep learning in 2006 [74]. Deep learning methods have *not* been explored extensively for software testing, and in particular for the test oracle problem. Recently, a few techniques have been proposed for automatic pattern-based bug detection. For example, Pradel et al. [127] propose *DeepBugs*, a deep learning-based static analysis tool for automatic name-based bug detection. Allamanis et al. [8] use graph-based neural static

analysis to detect variable misuse bugs. In addition to these techniques, several other deep learning methods for statically representing code have been developed [13, 10]. We do not discuss these further since we are interested in execution trace classification and in NNs that use dynamic trace information rather than a static view of the code.

Our first contribution is a deep-learning based approach to embed and classify program executions. In the next section, we provide the reader with the current research in program embeddings.

3.3 Program Representation

Creating representations is a vital process in transforming input data into a format that is readable by a machine learning model. Good representations highlight the data's key qualities and optimize the machine learner's overall performance. Source code representation is an active research field with a large impact on language models' ability to understand and generate synthetic programs.

Natural language representation techniques are usually also relevant to source code representation. Natural and programming languages share some structure similarities, a property known as *naturalness*. The *naturalness hypothesis* [7, 73, 72] states that software is a form of human communication and software corpora have similar statistical properties to natural language corpora. Using these properties is critical to build better software engineering tools. The naturalness hypothesis, inspires the goal to apply machine learning approaches to learn how developers naturally write and use code. These models can be used to augment existing tools with statistical information and enable new machine learning-based software engineering tools, such as recommender systems and program analyzers.

Representation of words as continuous vectors is a field with a long research history [50, 75, 132]. The Neural Network Language Model (NNLM) [88] is a popular architecture for word representation, where a feed-forward NN with a linear projection layer and a non-linear hidden layer learn jointly word vector representations. Following the NNLM, neural probabilistic models were presented [113, 23] where word vectors are first learnt using neural network with a single hidden layer. They are then used to train the NNLM. The benefit of this approach compared to the NNLM, is that word vectors can be obtained even without fully constructing the NNLM.

The N-gram model [26] is a simple, yet effective language representation model. N-gram models are statistical models and sentences as sequences of atomic units, each

statically representing a single word. N-gram models are easy to implement for simpler tasks but their performance is limited on more complex ones. In tasks such as machine translation where there is an abundance of training data, more advanced techniques are needed to achieve good performance.

Following the N-gram's limitations, a more advanced word embedding algorithm is *Word2Vec* [112]. This technique is based on neural networks that are trained to extract word embeddings from a corpora of natural language data efficiently. *Word2Vec* constructs vector representations of words that are dispersed on the feature space in such way so words that have a similar meaning or context in humans' language have vector representations that are also close to the feature space of the embedded vectors. Compared to non neural-based embedding algorithms, *Word2Vec* produces significantly more accurate word representations at a much lower computational cost.

These representation models have been widely used by language models to embed software. Such example is *code2vec* proposed by [14]. However, as source code has the feature of being represented as a graph, such representation models are naturally unequipped to capture this kind of structure which may be sub-optimal. This issue is addressed with graph representation models, most notably *Graph2Vec* [117]. *Graph2Vec* is a neural embedding framework that learns data-driven distributed representations of arbitrary sized graphs. *Graph2Vec*'s embeddings are learnt in an unsupervised manner and can be used for a range of downstream tasks, including graph classification apart from programming language modeling.

Embedded representations for programs are not restricted to encoding text-level source code. They can also be used to represent dynamic information, such as program executions, or *execution traces*. An execution trace is a recorded sequence of instructions executed for a given program, as well as any other data that have been accessed or modified throughout the execution. Execution traces capture the state of a program at different runtime states and can provide critical insight to software engineers about the behavior of a program. Our first contribution focuses in embedding execution traces to automate the classification of program runtime behaviour.

Wang et al. [151] propose an approach to embed program runtime flow. They use execution traces captured as a sequence of variable values at different program points. A program point is when a variable gets updated. Their approach uses RNNs to summarise the information in the execution trace. The execution trace embeddings are given as an input to a program repair tool. This embedding technique has several limitations - 1. Capturing execution traces as sequences of updates to every variable in

the program has an extremely high overhead and will not scale to large programs. The paper does not describe how the execution traces are captured, they simply assume they have them. 2. The approach does not discuss how variables of complex data types such as structs, arrays, pointers, objects are encoded. It is not clear if the traces only capture updates to user-defined variables, or if system variables are also taken into account. 3. The evaluation uses three simple, small programs (eg. counting parentheses in a string) from students in an introductory programming course. The complexity and scale of real programs is not assessed in their experiments. Their technique for capturing and directly embedding traces as sequences of updates to every variable is infeasible in real programs.

A novel blended approach to learn program representations with execution traces is presented by Wang et al. [152]. In this work, they collect a set of symbolic traces from programs, one for each execution path. They also obtain concrete traces from program executions, one for each test input. They create a blended trace by merging one symbolic trace with all concrete traces that exercise this corresponding execution path. They develop a NN architecture called *LiGer*, an attention-based RNN. Their model consists of a vocabulary embedding layer, a fusion layer and a programs embedding layer. The first encodes words to embedding vectors. In the fusion layer, one RNN embeds statements and a second RNN embeds all program states of that statement within the same time step. Attention vectors are calculated and concatenated using these embeddings as input. Finally, all attention vectors are fed into an RNN sequentially and all time outputs are pooled. *LiGer*'s embedding quality is evaluated with COSET [150] benchmark. Wang et al. [150] also extend their model into an encoder-decoder architecture and evaluate their model for the purpose of method name prediction. *LiGer* outperforms three relevant code embedding approaches across a set of benchmarks. However, their execution trace processing technique implies significant complexity. They only evaluate it on small functions with simple contexts. It is unlikely whether this technique can be scaled across multiple of functions of a real codebase. On the other hand, we show that our approach scales effectively over real and complex programs from different domains.

Recent works in the field of program representation includes GNN-based approaches that take advantage of the intermediate representations' graph structure. Guo et al. [63] develop a GNN-based approach to embed program binaries and fuse them with the semantics of control flow, data flow and call graphs into one model. They abstract programs into multiple graphs for multi-layer analysis. Their approach aims to enable

program analysis and compilation-related tasks and it achieves an accuracy of 83% in binary similarity detection and dead store prediction.

In GRAPHCODE2VEC, Ma et al. [107] propose a generic approach for task-agnostic, generic embeddings that capture syntax and semantics. GRAPHCODE2VEC’s embeddings are evaluated on a range of downstream tasks, including method name prediction, mutation testing classification and is compared against other state of the art embedding techniques such as CODEBERT [52] and CODE2SEQ, developing high performance in both generic and task-specific baselines. CODEBERT is a relevant publication to our second and third contributions and is discussed further in the next section.

Xu et al. [161] present M3V, a multi-modal, multi-view program embedding for repair operator prediction. Their proposed approach attempts to capture the context of a faulty location in a program using two models, (a) a tree-LSTM receiving text that captures the fault’s signature in natural language and (b) a Graph Neural Network that encapsulates its structure in two views, data and control dependencies. Xu et al. evaluate M3V against state of the art context embedding approaches in repairing two common types of bugs in Java, null pointer exceptions and index out of bounds. They improve the prediction of repair operators in repairing null pointer exceptions by 11% to 41% using their context embeddings and 9% to 30% in index out of bounds.

3.4 Language Modeling for Program Synthesis

In the previous two sections, we survey the current literature of test oracles and their intersection with machine learning in the form of language modeling and program representation. This section discusses the use of language modeling and program representation for compiler benchmark synthesis that aims in optimizing compiler heuristics.

In their 2017 survey, Allamanis et al. [7] describe the fast-moving field of deep language models for source code [7]. Wong et al. [160] develop *AutoComment* that mines StackOverflow to automatically generate code comments. Allamanis et al. [6] develop *Naturalize* which employs techniques developed in the natural language processing domain to model coding conventions. *JSNice* [129] leverages probabilistic graphical models to predict program properties such Javascript identifier names. Allamanis et al. [9] use attention-based neural networks to generate summaries of source code. *Nero* [43] uses an encoder-decoder architecture to predict method names in stripped binaries. This technique receives an input sequence of call sites from the execution of

a binary as an input and produces a predicted method name.

A recent work that has unlocked numerous applications in code generation and classification is BERT, a natural language representation model by Devlin et al.[46]. Contrary to previous language modeling tools [125, 128], BERT is designed to learn on unlabeled text data by jointly conditioning on both left and right context in all layers. This is achieved by learning to predict single words hidden behind [MASK] tokens. BERT achieves state of the art results in 11 natural language tasks and enables multiple applications of this architecture to a wide variety of difficult machine learning tasks, such as machine translation, question answering etc.

BERT has found many applications in programming languages. In CuBERT [91], Kanade et al. apply BERT over Python programs and evaluate it on the identification of typical mutation faults such as variable misuse localization, swapped operands, function-docstring mismatch and exception type checking. In CodeBERT [52], Feng et al. fine-tune BERT to perform NL-PL and PL-NL transformations. First, they synthesize functions from a set of natural language specifications. They also attempt to generate natural language documentation, provided a source code function. In both settings, they measure BERT’s prediction accuracy with respect to the ground truth, i.e. programs and documentation strings.

In the field of compiler benchmarks, there is limited work coming from generative modeling techniques. In 2017, Cummins et al. [40] develop CLGEN, a deep learning generator based on LSTM [78] for OpenCL programs. CLGEN learns source code representations with program fragments collected from GITHUB’s open source repositories. CLGEN is developed to enhance existing benchmark suites with synthetic ones in order to tackle the compiler benchmarks shortage. Mitigating this shortage compiler predictive models are improved by training on more datapoints with new, unseen features. They generate synthetic benchmarks to enhance existing datasets and use the enhanced datasets to train the Grewe et al.[61] heuristic model that predicts whether an OpenCL kernel should execute on the CPU or the GPU for optimal performance. Training the predictive model on training data enhanced with synthetic benchmarks improved its performance by $1.27\times$.

However, CLGEN poses several limitations, illustrated by Goens et al. recent case study [59]. First, after reproducing CLGEN, they show the Grewe et al. [61] predictive model performs better when trained on any of standard benchmarks or GITHUB code, while its performance gets worse when synthetic benchmarks are included in the training set. They also analyze the AST depth distribution of CLGEN’s samples

and compare it to standard benchmarks and code from GITHUB. They prove synthetic benchmarks are significantly smaller, $3\times$ on average. They provide similar results on the feature space coverage comparison, showing CLGEN covers a narrow space of features compared to human-written code, which already exists. This is evidence that CLGEN’s synthetic benchmarks do not improve the feature diversity of training data for compilers. CLGEN generates millions of different examples but the probability of generating a kernel with syntactic errors increases exponentially as new tokens are appended. This is due to its LSTM’s sequential architecture which is unable to regress to earlier parts of the generated sequence and repair them. As a result, only a tiny fraction of synthetic benchmarks compile and those that do are always small in size.

More recent techniques include SketchAdapt by Nye et al. [119]. SketchAdapt is a synthesizer that generates code from specifications. Their model learns program sketches by training a generator-synthesizer [19, 47] on input-output pairs of functions. Then they generate new program sketches that match a given I/O specification. The RNN generator samples a range of possible generic sketches that match the specification of an input/output pair. Sketches contain `<HOLE<` tokens, which the symbolic synthesizer fills sequentially with statements. They evaluate SketchAdapt on 2 types of input specifications: 1) A list of integers as I/O pairs and 2) a natural language description of source code. Their model performs better compared to the generator-only [19] and synthesizer-only architectures [47] separately.

SketchAdapt samples a pre-defined pool of operations (in the form of lambda functions) that may match a program’s behavior. This restricts the amount of different operations it can generate to the size of its operation pool. Also, the functions they infer are required to have inputs and outputs as input specifications, which further restricts its diversity. Finally, they train and evaluate on short sketches of up to 4 operations. We are inspired to use the `<HOLE>` token as a means to generate statements within the middle of an existing function. However, we do not face SketchAdapt’s restrictions.

Bruen et al. [44], propose a Tree2Tree approach for source code generation using Variational Autoencoders. They use GloVe [124] to embed representations of programs’ AST nodes. They encode and decode AST representations using Tree-LSTMs as defined by Tai et al. [140]. They learn latent code representations by minimizing the reconstruction loss of the variational autoencoder. Bruen et al. train their model on two million C++ compiling functions gathered from coding competitions. They evaluate their Tree2Tree model against a simple VAE with a simple LSTM encoder-decoder architecture (Seq2Seq) and measure BLEU [122] and compilation rate scores between

these two architectures. They also experiment with using both architectures as generative models. In this case, they sample 1000 random latent vector representations and feed them to the decoder to collect a new sample. Their Seq2Seq model achieves a compilation rate of up to 67% with greedy search, however the authors argue this happens because the model greedily selects always the most probable labels, leading to repetitive samples that compile. When sampling with temperature, their Seq2Seq model achieves a compilation rate of 40%. Their Tree2Tree architecture is able to generate a wider variety of unique samples, but only achieves a compilation rate of up to 22%, which translates to 200 functions out of these 1000 random latent vectors. This shows bad generalization on the trained dataset which accounts for 2 million C++ functions.

Researchers have recently begun to explore the field of infilling language models for more large scale tasks. Infilling is defined as the bi-directional synthesis process that produces valid tokens by observing the left and right context of an incomplete input sentence. Li et al. [102] develop *AlphaCode*, a coding problem-solving tool based on deep language modeling and reinforcement learning. AlphaCode uses deep natural language reasoning to translate a coding problem description into a valid solution. Among 5,000 participants, it achieves a ranking of top 54.3%. Fried et al. [54] develop *InCoder*, a Transformer-based language model that is trained on 28 programming languages to perform left-to-right generation or bi-directional generation (infilling) to provide executable functions. Guo et al. [62] develop *GRAMMFORMER* which guides sketch generation by the programming language grammar. This model places “holes” where the model is uncertain during generation, which will be filled at a later edit step.

There have also been code generating approaches coming from the field of program repair / program reconstruction from input/output specifications. Gupta et al. [64] develop a program generator - program repair framework named SED. SED is a two-stage code generator. First, a synthesizer receives the input/output specifications that must be met and generates a set of program candidates that are likely to satisfy them. Second, a neural debugger evaluates each candidate program and performs program repair to reform generated candidates into a function that will match the ground truth. Their model is trained and tested on Karel, an educational programming language. Gupta et al. evaluate the performance of three different synthesizer architectures with multiple configurations and measure (a) how well do generated programs generalize across the test cases they are expected to pass and (b) the accuracy of their debugger

to repair synthesized programs across a different count of mutations, 1 to 5. SED is a meaningful research work in the intersection between program generation and program repair. However, Karel is an unrealistic language and SED's generative performance on a large scale, complex programming is not evaluated.

Faustino et al. develop Anghabench [42], a collection of real world C programs mined from GITHUB. They argue generative models cannot be easily employed in general purpose compilers because benchmarks usually target very specific aspects of target hardware or programming language. In their approach, they mine C code from open-source repositories to tackle the benchmark shortage [40, 155]. The main challenge of their work is to automatically compile collected code. To do so, they use Psyche-C [110] type inference engine for C to apply type reconstruction and resolve missing dependencies. Structs, unions and other custom data types are all omitted or re-declared with primitive types, if possible. They collect around 1.5 million C compiling functions. Anghabench does not support custom data types and user-defined functions and this is a serious limitation which alters original programs' semantics and leads to more simplistic data type relations. Furthermore, their benchmarks are compiling, but cannot be executed.

Chapter 4

Supervised learning over test executions as a test oracle

4.1 Introduction

The research contributions presented in this Chapter are the design and implementation of a NN-based test oracle for automatic program runtime correctness classification. This Chapter describes the underlying approach, the implementation and the evaluation of this technique using fifteen industry-standard codebases.

We explore supervised machine learning to infer a test oracle from labelled execution traces of a given system. In particular, we use neural networks (NNs), well suited to learning complex functions and classifying patterns, to design the test oracles. Our technique is widely applicable and easy to use, as it only requires execution traces gathered from running tests through the program under test (PUT) to design the oracle. This is shown in Figure 4.1 where a small fraction of the gathered execution traces labelled with pass/fail (shown in light gray) is used to train the NN model which is then used to automatically classify the remaining unseen execution traces (colored dark gray).

Previous work exploring the use of NNs for test oracles has been in a restricted context – applied to very small programs with primitive data types, and only considering their inputs and outputs [148, 87]. Information in execution traces which we believe is useful for test oracles has not been considered by existing NN-based approaches. Other bodies of work in program analysis have used NNs to predict method or variable names and detecting name-based bug patterns [12, 127] relying on static program information, namely, embeddings of the Abstract Syntax Tree (AST) or source code. Our

proposed approach is the first attempt at using *dynamic execution trace information in NN models for classifying test executions*.

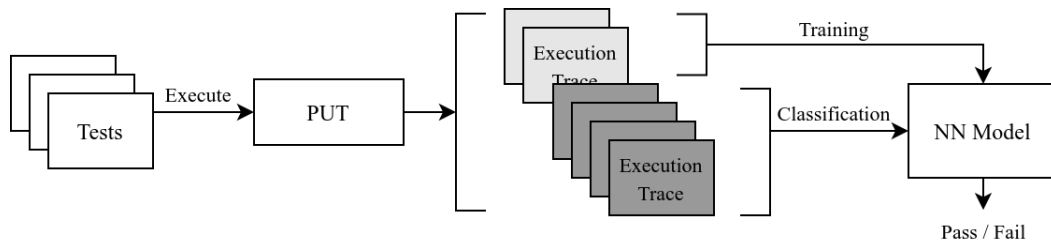


Figure 4.1: Key idea in our approach.

Our approach for inferring a test oracle has the following steps,

1. Instrument a program to gather execution traces as sequences of method invocations.
2. Label a small fraction of the traces with their classification decision.
3. Design a NN component that embeds the execution traces to fixed length vectors.
4. Design a NN component that uses the line-by-line trace information to classify traces as pass or fail.
5. Train a NN model that combines the above components and evaluate it on unseen execution traces for that program.

The novel contributions in this technique are in Steps 3, 4 and 5. Execution traces from a program vary widely in their length and information. We propose a technique to encode and summarise the information in a trace to a fixed length vector that can be handled by a NN. We then design and train a NN to serve as a test oracle.

Labelled execution traces. Given a PUT and a test suite, we gather execution traces corresponding to each of the test inputs in the test suite with our instrumentation framework. Effectively learning a NN classifier for a PUT that distinguishes correct from incorrect executions requires labelled data with both passing and failing examples of traces. We require a small fraction of the overall execution traces to be labelled, which is likely to be a manual process. As a result, our proposed approach for test oracle is *not* fully automated. We hypothesize that the time invested in labelling a small proportion of the traces is justified with respect to the benefit gained in automatically classifying the remaining majority of traces. In contrast, with no classifier, the developer would have had to specify expected output for all the tests, which is clearly more time consuming than the small proportion of tests we need labelled.

NN Architecture. An execution trace in our approach comprises multiple lines, with each line containing information on a method invocation. Our architecture for encoding and classifying an execution trace uses multiple components: (1) Value encoder for encoding values within the trace line to a distributed vector representation, (2) Trace encoder encoding trace lines within a variable-length trace to a single vector, and (3) Trace Classifier that accepts the trace representation and classifies the trace. The components in our architecture are made up of LSTMs, one-hot encoders, and a multi-layer perceptron. We select LSTMs to represent execution traces because they are fast to train, require significantly less training data than other architectures (e.g. Transformers) and show high accuracy for the task of predicting failing execution traces.

Case Studies. We evaluate our approach using 4 subject programs and tests from different application domains - a single module from Ethereum project [30], a module from Pytorch [123], one component within Microsoft SEAL encryption library [134] and a Linux stream editor [104]. One of the 4 subject programs were accompanied by both passing and failing tests that we could directly use in our experiment. The remaining three programs were only accompanied by passing tests. We treated these programs as reference programs. We then generated PUTs by artificially seeding faults into them. We generated traces through the PUTs using the existing tests, labelling the traces as passing or failing based on comparisons with traces from the reference program. We trained a NN model for each PUT using a fraction of the labelled traces. We found our approach for designing a NN classification model was effective for programs from different domains. We achieved high accuracies in detecting both failing and passing traces, with an average precision of 89% and recall of 88%. Only a small fraction of the overall traces (average 15%) needed to be labelled for training the classification models.

In summary, the Chapter makes the following contributions:

- Given a PUT and its test inputs, we provide a framework that instruments the PUT and gathers test execution traces as sequences of method invocations.
- A NN component for encoding variable-sized execution traces into fixed length vectors.
- A NN for classifying the execution traces as pass or fail.
- We provide empirical evidence that this approach yields effective test oracles for programs and tests from different application domains.

4.1.1 Extended Contributions

In addition to our initial publication [145] of the contributions presented in the previous Subsection, we published extended contributions in [146]. The extended contributions are summarised as follows,

1. Support for Java programs. Our work in [145] provided tool support in the LLVM [97] framework to instrument the intermediate representation (LLVM-IR) of programs to gather execution traces. LLVM, however, does not provide front-end support for Java programs. In this Chapter, we provide tool support to gather execution traces for Java programs using the Soot framework [147].

2. Extensive empirical evaluation. We augment the experiments in [145] with 10 additional subject programs – 9 network protocols from L7-filter [35] and 1 Java utilities library from Defects4J [90], a database of real faults for open-source Java programs. For these subject programs, we evaluate precision, recall and specificity of our approach in classifying execution traces. We also assess the size of training set needed and compare accuracies against a hierarchical clustering technique for classifying execution traces proposed by Almaghairbe et al. [11].

3. Generalisation. We conduct an initial exploration into the ambitious possibility of using a model, trained using traces from one subject program, to classify traces from other programs in the same application domain. We use FSMs (*Finite State Machines*) from the network protocol domain to evaluate this possibility.

We found our approach for designing a NN classification model was effective for all subject programs. We achieved high accuracies in detecting both failing and passing traces, with an average precision of 93% and recall of 94%. Only a small fraction of the overall traces (average 14%) needed to be labelled for training the classification models. We found generalisation of a classification model from one network protocol to others in the domain was feasible. Generalisation accuracies were not as high as accuracy achieved using separate classification models, around 70%, but we believe there is scope for improvement using fine tuning in the future.

4.2 Approach

Our approach for building an automated test oracle for classifying execution traces has the following steps,

Step 1: Instrument the PUT to gather traces when executing the test inputs.

Step 2: Preprocess the traces to prune unnecessary information.

Step 3: Encode the preprocessed traces into vectors that can be accepted by the neural network.

Step 4: Design a NN model that takes as input an encoded trace, and outputs a verdict of pass or fail for that trace.

Figure 4.2 illustrates the steps in our approach, with the bottom half of the figure depicting steps 3 and 4 for any given preprocessed trace from step 2. We discuss each of the steps in the rest of this Section.

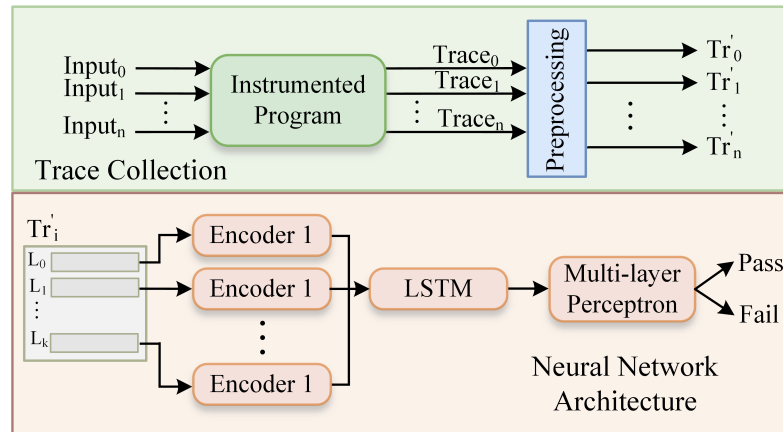


Figure 4.2: Gathering traces, encoding them, and using NNs to classify them. ENCODER 1 constructs a fixed vector representation per trace line. A second LSTM Encoder receives all trace line representations as a sequence and outputs a vector that summarises the execution trace. Both models are jointly trained with the MLP such that a low error is achieved in predicting the trace's label.

4.2.1 Instrument and Gather Traces

For every test input executed through the PUT, we aim to collect an execution trace as a sequence of method invocations, where we capture the name of the method being called, values and data types of parameters, return values and their types, and, finally, the name of the parent method in the call graph. We find gathering further information, eg. updates to local variables within each method, incurs a significant overhead and is difficult to scale to large programs. To gather this information we develop two different tools with support for different programming languages to make our framework widely applicable. Our first instrumentation tool is primarily aimed at C/C++ programs and

uses the middleware of LLVM [97] and instruments the intermediate representation (LLVM-IR) of programs. This allows our implementation to be language-agnostic. LLVM provides front-end support for multiple programming languages in addition to C/C++ like CUDA, Haskell, Swift, Rust among others, along with numerous libraries for optimisation and code generation.

Our second tool is aimed at Java programs and uses Soot [147] to collect execution traces. Soot is a Java optimization framework and provides libraries for users to analyze, instrument and optimize applications. We develop a pass with Soot to compile Java into bytecode and a second pass to convert bytecode to Jimple-IR. Jimple is a typed 3-address intermediate representation suitable for code transformations. Using Soot's API, we develop a second pass to instrument Jimple and finally compile into an executable. Our Soot framework is also compatible with any other programming language that can be compiled into Java bytecode, e.g. Scala.

To perform the instrumentation, we traverse the PUT, visiting each method. Every time a method invocation is identified, code is injected to trace the caller-callee pair, the arguments and the return values. At the end of the program, code is inserted to write the trace information to the output.

Each trace contains a sequence of method invocations. This sequence comprises multiple lines, each line being a tuple (n_p, n_c, r, a) that represents a single method invocation within it having:

- The names of the caller (parent) n_p and called n_c functions.
- Return values r of the call, if any.
- Arguments passed a , if any.

The order of trace lines or method invocations is the order in which the methods complete and return to the calling point.

Our LLVM instrumentation supports all variable types including primitive types (such as `int`, `float`, `char`, `bool`), composite data types (such as structs, classes, arrays) defined by a user or library, and pointers for return and argument values. Structs and classes are associated with a sequence of values for their internal fields. We instrument these data structures in a depth first fashion, until all primitive types are traced. For pointers, we monitor the values they refer to.

Our instrumentation within Soot collects all Java primitive types, strings, primitive wrapper classes, atomic wrapper classes and arrays. We also support custom classes

that are defined within the scope of a subject Java translation unit. Soot allows the instrumentation of public class members only; private methods and variables are not accessed.

4.2.2 Training Set

We execute the instrumented program with each test input in the test suite to gather a set of traces. A subset of the traces is labelled and used in training the classification model. To label the traces as pass or fail, we compare actual outputs through the PUT with expected outputs provided by a reference program or the specifications. Section 4.3.1 discusses how we label traces for the subject programs in our experiment. It is worth noting that in our approach, the developer will only need to provide expected outputs for a *small proportion of tests rather than the whole test suite*. In the absence of expected output in tests, how will tests be labelled is a common question. Answering this question will depend on what is currently being done by the developer or organisation for classifying tests as pass or fail. Our approach will entail applying the same practice to labelling, albeit to a significantly smaller proportion of tests.

To avoid data leakage, i.e. information about the expected outcome of the program existing in the trace, in our experiment in Section 4.3, we ensure that expected output is removed from the traces. We also remove exceptions, assertions and any other information in the program or test code that may act as a test oracle. This is further discussed in Section 4.3.2.

4.2.3 Preprocessing

The execution traces gathered with our approach include information on methods declared in external libraries, called during the linking phase. To keep the length of the traces tractable and relevant, we preprocess the traces to only keep trace lines for methods that are defined within the module, and remove trace lines for declared functions that are not defined, but simply linked to later.

For method invocations within loops, a new trace line is created for each invocation of the same method within the loop. For loops with large numbers of iterations, this can lead to redundancy when the method is invoked with similar arguments and return values. We address this potential redundancy issue by applying average pooling to trace lines with identical caller-callee methods within loops. This helps summarize huge sequences of function calls with identical caller-callee methods and similar

features into compressed representations to save memory and time.

4.2.4 Neural Network Model

In this step, we perform the crucial task of designing a neural network that learns to classify the pre-processed traces as passing or failing. Shape and size of the input traces vary widely, and this presents a challenge when designing a NN that accepts fixed length vectors summarizing the traces. To address this, our network comprises three components that are trained jointly and end-to-end: 1. a VALENC that encodes values (such as the values of arguments and return values) into D_V -dimensional distributed vector representations, shown within ENCODER 1 in Figure 4.3, 2. a TRENC that encodes variable-sized traces into a single D_T -dimensional vector, shown as LSTM in Figure 4.4, and finally, 3. a TRACECLASSIFIER that accepts the trace representation for state and predicts whether the trace is passing or failing. The MULTI-LAYER PERCEPTRON in Figure 4.2 represents the TRACECLASSIFIER. We describe each component in detail in the rest of this Section.

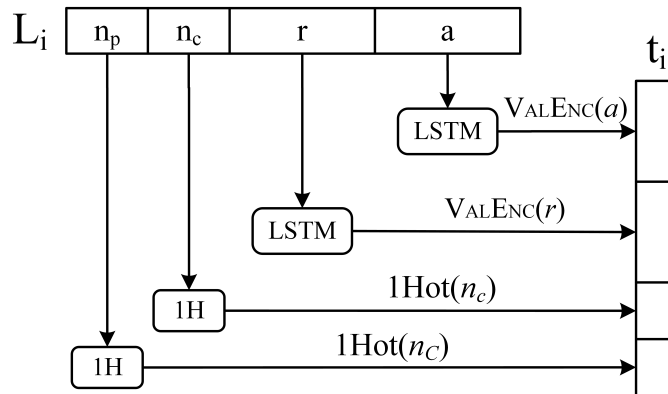


Figure 4.3: ENCODER 1 representing a single line in a trace as a vector containing function caller, callee names, arguments and return values.

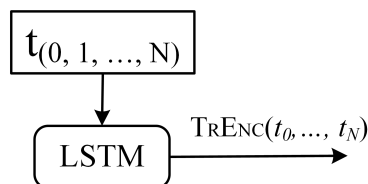


Figure 4.4: ENCODER 2 representing a sequence of trace lines as a single vector.

Encoding Values. Values within the trace provide useful indications about classifying a trace. However, values — such as ints, structs, and floats — vary widely

in shape and size. We, therefore, design models that can summarize variable-sized sequences into fixed-length representations. In the machine learning literature, we predominantly find three kinds of models that can achieve this: recurrent neural networks (RNNs), 1D convolutional neural networks (CNN) and transformers. In this work, we employ LSTMs [79] — a commonly used flavour of RNNs. Testing other models is left as future work. At a high-level RNNs are recurrent functions that accept a vector \mathbf{h}_t of the current state and an input vector \mathbf{x}_t and compute a new state vector $\mathbf{h}_{t+1} = RNN(\mathbf{x}_t, \mathbf{h}_t)$ which “summarizes” the sequence of inputs up to time t . A special initial state \mathbf{h}_0 is used at $t = 0$.

To encode a value v , we decompose it into a sequence of primitives $v = [p_0, p_1, \dots]$ (integers, floats, characters, etc.). Each primitive p_i is then represented as a binary vector $\mathbf{b}_i = e(p_i)$ containing its bit representation padded to the largest primitive data type of the task. For example, if `int64` is the largest primitive then all \mathbf{b}_i s have dimensionality of 64. This allows us to represent all values (integers, floats, strings, structs, pointers, etc.) as a unified sequence of binary vectors. We encode v into a D_V -dimensional vector by computing

$$\text{VALENC}(v) = \text{LSTM}_v(e(p_L)_L, \text{VALENC}([p_0, p_1, \dots, p_{L-1}])),$$

where LSTM_v is the LSTM that sequentially encodes the \mathbf{b}_i s. Note that we use the same `VALENC` for encoding arguments and return values, as seen in Figure 4.3. The intuition behind this approach is that the bits of each primitive can contain valuable information. For example, the bits corresponding to the exponent range of a float can provide information about the order of magnitude of the represented number, which in turn may be able to discriminate between passing and failing traces.

Representing a Single Trace Line. Armed with a neural network component that encodes values, we can now represent a single line (n_p, n_c, r, a) of the trace. To do this, we use `VALENC` to encode the arguments a and the return value r . We concatenate these representations along with one-hot representations of the caller and callee identities, as shown in Figure 4.3. Specifically, the vector encoding \mathbf{t}_i of the i th trace line is the concatenation

$$\mathbf{t}_i = [\text{VALENC}(a), \text{VALENC}(r), \text{1HOT}(n_p), \text{1HOT}(n_c)],$$

where `1HOT` is a function that takes as input the names of the parent or called methods and returns a one-hot vector that uniquely encodes that method name. For methods that are rare (appear fewer than k_{min} times) in our data, `1HOT` collapses them to a single

special Unknown (UNK) name. This is similar to other machine learning and natural language processing models and reduces sparsity often improving generalization. The resulting vector \mathbf{t}_i has size $2D_V + 2k$ where k is the size of each one-hot vector.

Encoding Traces. Now that we have built a neural network component that encodes single lines within a trace, we design TRENC that accepts a sequence of trace line representations $\mathbf{t}_0 \dots \mathbf{t}_N$ and summarizes them into a single D_T -dimensional vector as shown in Figure 4.4. We use an LSTM with a hidden size D_T , and thus

$$\text{TRENC}(\mathbf{t}_0 \dots \mathbf{t}_N) = \text{LSTM}_{tr}(\mathbf{t}_N, \text{TRENC}(\mathbf{t}_0 \dots \mathbf{t}_{N-1})),$$

where $\text{LSTM}_{tr}()$ is an LSTM network that summarizes the trace line representations.

Classifying Traces. With the neural network components described so far we have managed to encode traces into fixed length vector representations. The final step is to use those computed representations to make a classification decision. We treat failing traces as the positive class and passing traces as the negative class since detecting failing runs is of more interest in testing. We compute the probability that a trace is failing as

$$P(\text{fail}) = \text{TRACECLASSIFIER}([\text{TRENC}(\mathbf{t}_0 \dots \mathbf{t}_N)]),$$

where the input of TRACECLASSIFIER is the output vector of TRENC. Our implementation of TRACECLASSIFIER is a multilayer perceptron (MLP) with sigmoid nonlinearities and a single output. The sigmoid enables the model’s raw output to be viewed as a probability of the trace to represent a failing execution. It therefore helps developers measure the model’s prediction certainty. It follows that $P(\text{pass}) = 1 - P(\text{fail})$.

Training and Implementation Details. We train our network end-to-end in a supervised fashion, minimizing the binary cross entropy loss. All network parameters (parameters of LSTM_V and LSTM_{tr} and parameters of the MLP) are initialized with random noise. For all the runs on our network we use $D_V = 128$, $D_T = 256$. The TRACECLASSIFIER is an MLP with 3 hidden layers of size 256, 128 and 64. We use the Adam optimizer [94] with a learning rate of $10e - 5$.

For our subject programs, we find the aforementioned feature values to be optimal for performance and training time, after having experimented with other NN architectures, varying the D_V , D_T sizes, and the hidden layers in the MLP. We explored increasing D_V to 256, 512, D_T to 512, 1024 and size of hidden layers to 512 and 1024.

To handle class imbalance in datasets, we explicitly counteract the imbalance in the loss function by down-weighting the samples within the most popular class such that samples of both class participate equally within this function.

Our implementation of the proposed approach is available at <https://github.com/fivosts/Learning-over-test-executions>.

4.3 Experiment

In our experiment, we evaluate the feasibility and accuracy of the NN architecture proposed in Section 4.2 to classify execution traces for 15 subject programs and their associated test suites. The selection of our subject programs is based on the build system and compiler they use. Our instrumentation tool depends on CMake and LLVM-7, therefore all case studies are required to satisfy this restriction. We also select programs with at least a few hundreds of test cases to enable the NN’s training process. We investigate the following questions regarding feasibility and effectiveness:

Q1. Precision, Recall and Specificity: *What is the precision, recall and specificity achieved over the subject programs?*

To answer this question, we use our tool to instrument the source code to record execution traces as sequences of method invocations, arguments and return values. A small fraction of the execution traces are labelled (*training set*) and fed to our framework to infer a classification model. We then evaluate precision, recall and specificity achieved by the model over unseen execution traces (*test set*) for that program. The test set includes both passing and failing test executions. We use *Monte Carlo cross-validation*, creating random splits of the dataset into training and test data. We created 15 such random splits and averaged precision, recall and specificity computed over them.

Q2. Size of training set: *How does size of the training set affect precision and recall of the classification model?*

For each program, we vary the size of training set from 5% to 30% of the overall execution traces and observe its effect on precision and recall achieved.

Q3. Comparison against state of art: *How does the precision, recall and specificity achieved by our technique compare against agglomerative hierarchical clustering, proposed by Almaghairbe et al. [11] in 2017?*

We choose to compare against the hierarchical clustering work as it is the most relevant and recent in classifying execution traces, even though it is not NN-based like

our approach. Traces used in their work are sequences of method invocations, similar to our approach. Other test oracle work that use NNs is not used in the comparison as they do not work over execution traces, and are limited in their applicability to programs with numerical input and output which is not the case for programs in our experiment.

Q4. Generalisation of classification model: *Can a classification model inferred from execution traces of one program be used to classify test executions over other programs in the same domain?*

For the network protocol domain, we evaluate the accuracy of using a classification model inferred using traces from a single protocol detection finite state machine (FSM) to classify test executions from other protocol FSMs.

4.3.1 Labelling Traces

All our subject programs are open source, and most of them were only accompanied by passing tests. This is not uncommon as most released versions of programs are correct for the given tests. We take these correct programs to be reference implementations. To enable evaluation of our approach that distinguishes correct versus incorrect executions, we need subject programs with bugs. We, therefore, generate PUTs by automatically mutating the reference implementation using common mutation operators [86] listed below,

1. Arithmetic operator replacement applied to $\{+, -, *, /, --, ++\}$.
2. Logical connector replacement applied to $\{\&\&, ||, !\}$.
3. Bitwise operator replacement applied to $\{\&, |, \wedge, \ll, \gg\}$.
4. Assignment operator replacement applied to $\{+ =, - =, * =, / =, \% =, \ll =, \gg =, \& =, | =, \wedge =\}$.

A PUT is generated by seeding a single fault into the reference implementation at a random location using one of the above mutation operators. We used an independent open source mutation tool¹ to generate PUTs from a given reference program. Figure 4.5 shows a PUT generated by seeding a single fault into a reference program. As seen in Figure 4.5, we run each test, T_i , in the test suite, through both the reference program and PUT, and label the trace as *passing* if the expected output, EO_i , from

¹<https://github.com/chao-peng/mutec>

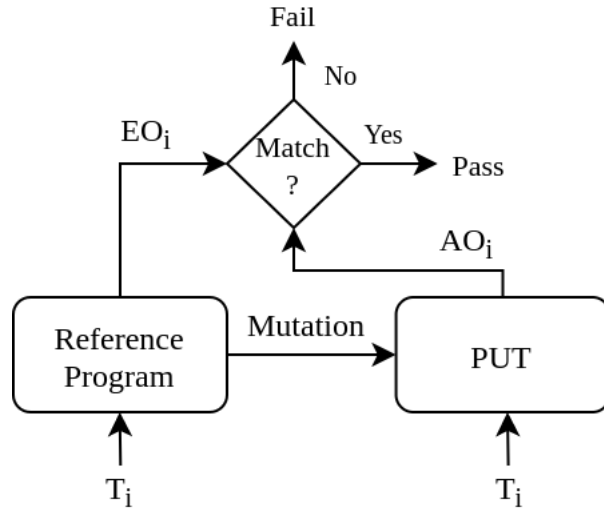


Figure 4.5: Labelling test executions by matching actual and expected behavior.

the reference matches the actual output, AO_i , from the PUT. If they do not match, the trace is labelled as *failing*. We rejected PUTs from mutations that did not result in any failing traces (outputs always match with the reference). This avoids the problem of equivalent mutants. All the PUTs in our experiment had both passing and failing traces.

4.3.2 Subject Programs

We chose subject programs from different domains to assess applicability of our approach, namely from the blockchain, deep learning, encryption and text editing domains. A description of the programs and associated tests is as follows.

- Ethereum** [30] is an open-source platform based on blockchain technology, which supports smart contracts. Within it, we evaluate our approach over the DIFFICULTY module that calculates the mining difficulty of a block, in relation to different versions (eras) of the cryptocurrency (Byzantium, Homestead, Constantinople etc.). The calculation is based on five fields of an ETHEREUM block, specified in the test input.

Tests. We use the default test inputs provided by Ethereum’s master test suite for the DIFFICULTY module. We test this module for the Byzantium era of the cryptocurrency (version 3.0). The test suite contains 2254 *automatically* generated test inputs, using fuzzing. Each test input contains one hex field for the test input of the difficulty formula and another hex field for the expected output of the program. All the test

inputs provided with the module are passing tests with the actual output equal to the expected output. As a result, we use the provided module as a reference implementation. As described in Section 4.3.1, we seed faults into the reference implementation to generate PUTs, each containing a single mutation. For the difficulty module, we generate 2 PUTs – 1. Ethereum-SE with a seeded fault in the core functionality of the difficulty module, and 2. Ethereum-CD with a fault seeded in one of the functions that is external to the core function but appears in the call graph of the module. The balance between passing and failing tests varies between the two PUTs, Ethereum-CD being perfectly balanced and Ethereum-SE being slightly imbalanced (828 failing and 1426 passing tests).

2. Pytorch [123] is an optimized tensor library for deep learning, widely used in research. In our experiment, we evaluate our model over the `intrusive_ptr` class, which implements a pointer type with an embedded reference count. We chose this class because it had a sizeable number of tests (other modules had < 20 published tests).

Tests. Implementation of the class is accompanied by 638 tests, all of which are passing. We, thus, use this as the reference implementation. As with ETHEREUM, we apply mutations to the `intrusive_ptr` implementation to generate a single PUT. Upon comparison with the reference, 318 of the existing tests are labelled passing through the PUT and 320 as failing.

3. Microsoft SEAL [134] is an open-source encryption library. In our experiment, we study one component within Microsoft SEAL, the ENCRYPTOR module, which is accompanied by tests. This component is responsible for performing data encryption.

Tests. The ENCRYPTOR component is accompanied by 133 tests. The provided tests were all passing tests, with matching expected and actual output. As with previous programs, we generate a PUT by mutating the original implementation. On the PUT, 11 tests fail and 122 pass.

4. Sed [104] is a Linux stream editor that performs text transformations on an input stream.

Tests. We use the fifth version of SED available in the SIR repository [49]. This version is accompanied by 370 tests, of which 352 are passing and 18 are failing. The failing tests point to real faults in this version. Since the implementation was accompanied by both passing and failing, we used it as the PUT. We did *not* seed faults to generate the PUT.

5. L7-Filter [35] is a packet identifier for Linux. It uses regular expression match-

ing on the application layer data to determine what protocols are used. It works with unpredictable, non-standard and shared ports. We study the following 9 protocols, implemented as FSMs, separately in our evaluation

1. ARES - P2P filesharing
2. BGP - Border Gateway Protocol
3. BIFF - new mail notification
4. FINGER - User information server
5. FTP - File Transfer Protocol
6. RLOGIN - remote login
7. TEAMSPEAK - VoIP application
8. TELNET - Insecure remote login
9. WHOIS - query/response system (eg. for domain name)

Tests. For each of the network protocol FSMs, we use test suites generated by Yaneva et al.[162] that provide all-transition pair coverage. The test suites for the FSMs have both passing and failing tests.

6. commons-lang [17] is a java library from Apache Commons with utility classes for the `java.lang` API. This is a large codebase and contains Java classes such as `OBJECT` and `CLASS`. We gather this subject program from the Defects4J database that provides several versions of this library and a labelled set of passing and failing test cases for each version.

Tests. Defects4J contains different versions of `COMMONS-LANG`. Most of them have very few, or even no failing tests. These versions do not provide our model with failing examples to learn and predict, therefore we discard them. We use the 34th version of this program, which contains 559 passing and 27 failing tests. These 27 tests are caused by real bugs found in this version of the subject program, therefore we do not seed faults to generate the PUT.

Checks to avoid data leakage. We ensure no test oracle data is leaked into traces. We remove expected outputs, assertions, exceptions, test names and any other information that may act directly or indirectly as a test oracle. For example, Ethereum

uses BOOST testing framework to deploy its unit tests. We remove expected outputs and assertions in the test code that compare actual with the expected output e.g. BOOST_CHECK_EQUAL.

For PUTs generated by seeding faults into the reference implementation, we only use one PUT for each reference implementation except in the case of Ethereum where we generated two PUTs, since faults were seeded in different files. Generating more PUTs for each reference implementation would be easy to do. However, we found our results across PUTs for a given reference program only varied slightly. As a result, we only report results over one to two PUTs for each reference implementation.

4.3.3 Performance Measurement

For each PUT, we evaluate performance of the classification model over unseen execution traces. As mentioned in Section 4.2.4, we use positive labels for failing traces and negative labels for passing. We measure

1. *Precision* as the ratio of number of traces correctly classified as “fail” (TP) to the total number of traces labelled as “fail” by the model (TP + FP).
2. *Recall* as the ratio of failing traces that were correctly identified (TP / (TP + FN)).
3. *Specificity* or true negative rate (TNR) as the ratio of passing traces that were correctly identified (TN / (TN + FP)).

TP, FP, TN, FN represent true positive, false positive, true negative and false negative, respectively.

4.3.4 Hierarchical Clustering

In research question 3 in our experiment, we compare the classification accuracy of our approach against agglomerative hierarchical clustering proposed by Almaghairbe et al. [11]. Their technique also considers execution traces as sequences of method calls, but only encoding callee names. Caller names, return values and arguments are discarded. We attempted to add the discarded information, but found the technique was unable to scale to large number of traces due to both memory limitations and a time complexity of $O(n^3)$ where n is the number of traces. For setting clustering parameters for each subject program, we evaluate different types of linkage (SINGLE, AVERAGE,

COMPLETE) and a range of different cluster counts (as a percentage of the total number of tests): 1, 5, 10, 20 and 25%. We use Euclidean distance as the distance measure for clustering. For each program, we report the best clustering results achieved over all parameter settings.

4.3.5 Results

In this Section, we present and discuss our results in the context of the research questions presented in Section 4.3.

4.3.6 Q1. Precision, Recall and Specificity

Table 4.1 shows the precision, recall and specificity achieved by the classification models in our approach for the different PUTs. Results with the hierarchical clustering approach by Almaghairbe et al. [11] are also presented in Table 4.1 for comparison, but this is discussed in Q3 in Section 4.3.8. The column showing % of traces used in training varies across programs, we show the lowest percentage that is needed to achieve near maximal precision and recall.

The classification models for all 15 PUTs achieve more than 71% precision and 86% recall, with an average of 93% and 94%, respectively. Our technique works particularly well for Pytorch, Sed and all networking protocols, achieving $\geq 94\%$. This implies that the number of false positives in the classification is very low and a large majority of the failing traces are correctly identified.

The classification models for all PUTs also achieve high specificity ($\geq 79\%$, average 96%). This implies that the NN models are able to learn runtime patterns that distinguish not only failing executions, but also passing executions with a high degree of accuracy. These results are unprecedented as we are not aware of any technique in the literature that can classify both passing and failing executions at this level of accuracy.

Analysis. To understand the results in Table 4.1, for each of the PUTs, we inspected and compared passing and failing traces using a combination of longest common subsequence, syntactic diffs, and manual inspection. We also performed *ablation* - systematically removing information (one parameter at a time) from the traces, training new classification models with the modified traces and observing their effect on precision, recall and specificity (TNR). In our experiments, we systematically remove the following parameters from the original traces – function call names, arguments,

PUT	Lines of Code	% Traces for training	Total # Traces	Our Approach			Hierarchical Clustering [11]		
				Precision	Recall	TNR	Precision	Recall	TNR
Ethernum-CD	55927	15	2254	0.80	0.82	0.79	1.0	0.49	1.0
Ethernum-SE	55927	15	2254	0.99	0.82	0.86	1.0	0.25	1.0
Pytorch	21090	10	638	0.99	0.98	0.99	0.48	1.0	0.16
SEAL_Encryptor	25967	30	132	0.75	0.86	0.98	0.16	0.36	0.83
Sed	4492	10	370	0.94	0.94	0.99	0.35	0.63	0.86
commons-lang	49028	40	586	0.71	0.94	0.98	0.07	0.96	0.4
Ares protocol	1261	3	16066	0.97	0.98	0.97	0.94	0.24	0.0
BGP protocol	1025	5	16009	0.99	0.99	0.99	0.18	0.01	0.98
Biff protocol	627	15	1958	0.97	0.99	0.99	0.43	0.22	0.72
Finger protocol	791	10	2775	0.99	0.99	0.99	0.53	0.13	0.92
FTP protocol	995	10	9677	0.99	0.99	0.98	0.07	0.001	0.98
Rlogin protocol	955	10	4121	0.97	0.96	0.99	1.0	0.04	1.0
Teamspeak protocol	3284	10	1945	0.95	0.99	0.96	1.0	0.11	1.0
Telnet protocol	1019	10	319	0.98	0.96	0.95	0.29	0.02	0.87
Whois protocol	784	9	4412	0.98	0.99	0.99	0.49	0.03	0.98

Table 4.1: Precision, Recall and True Negative rate (TNR) using our approach and hierarchical clustering.

and return values. Tables 4.2 and 4.3 show the results from our ablation study. We discuss results for each of the programs in the following paragraphs.

Over SEAL Encryptor, our approach achieves 75% precision, 86% recall and 98% specificity when trained with 30% of the traces. Encryptor requires a higher fraction of traces for training when compared to other PUTs, as the number of failing traces is very small (= 11), unlike other programs. Although we handle imbalance in datasets by weighting samples in the loss function, the NN still needs some representatives of the failing class during training. Using 10% of the traces in training, will only provide one example of failing trace (10% of 11) which is not enough for the NN model to learn to distinguish failing versus passing behaviour. Training using 30% of the traces includes 3 failing traces which allows the NN to achieve 75% precision. High precision with only 3 failing traces is because all the failing traces for this program have the same call sequence, which is sufficiently different from passing traces. Passing traces do not all have the same sequence. However, due to the availability of a larger set of passing traces (training with 30% is 40 passing traces), the NN is able to identify the different method call patterns in passing traces accurately (98% specificity). The ablation study in Tables 4.2 and 4.3 shows that all the parameters contribute to model performance as removing them has a detrimental effect.

PUT	Omitted Info.	P	R	TNR
Ethereum-CD	Function names	0.63	0.64	0.62
	Return values	0.68	0.87	0.60
	Arguments	0.54	0.78	0.35
Ethereum-SE	Function names	0.96	0.84	0.35
	Return values	0.99	0.97	0.93
	Arguments	0.96	0.84	0.33
Pytorch	Function names	0.99	1.0	1.0
	Return values	0.99	0.99	0.99
	Arguments	0.51	0.99	0.04
Seal Encryptor	Function names	0.53	0.87	0.92
	Return values	0.46	0.99	0.90
	Arguments	0.28	0.88	0.76
Sed	Function names	0.19	0.72	0.24
	Return values	0.48	0.52	0.85
	Arguments	0.30	0.40	0.73
commons- lang	Function names	0.58	1.00	0.97
	Return values	0.74	0.88	0.99
	Arguments	0.78	0.95	0.99

Table 4.2: Precision (P), Recall (R) and Specificity (TNR) for each PUT omitting certain trace information.

PUT	Omitted Info.	P	R	TNR
Finger	Function names	0.99	0.95	0.99
	Return values	0.98	0.97	0.99
	Arguments	0.52	0.19	0.88
Telnet	Function names	0.93	1.0	0.76
	Return values	0.82	1.0	0.25
	Arguments	0.76	1.0	0.00
Ares	Function names	0.96	0.98	0.95
	Return values	0.95	0.99	0.75
	Arguments	0.93	0.96	0.68
BGP	Function names	0.99	0.98	0.98
	Return values	0.98	0.99	0.98
	Arguments	0.97	0.97	0.97
Biff	Function names	0.58	0.84	0.41
	Return values	0.56	0.92	0.35
	Arguments	0.51	0.64	0.40
FTP	Function names	0.99	0.99	0.98
	Return values	0.97	0.97	0.98
	Arguments	0.88	0.93	0.84
Rlogin	Function names	0.95	0.96	0.96
	Return values	1.0	0.92	1.0
	Arguments	0.85	0.91	0.94
Teamspeak	Function names	0.91	0.97	0.91
	Return values	0.94	0.98	0.94
	Arguments	0.77	0.86	0.77
Whois	Function names	0.96	0.96	0.96
	Return values	0.96	0.96	0.96
	Arguments	0.72	0.75	0.73

Table 4.3: Precision (P), Recall (R) and Specificity (TNR) for each PUT omitting certain trace information.

For PyTorch, we achieve 99% precision, 98% recall and 99% specificity when trained with 10% of the traces. The dataset for PyTorch PUT is balanced (318 passing and 320 failing). 10% of the traces during training provides sufficient examples from both passing and failing classes for the NN to learn to distinguish them. We find the reason for the superior performance of our model over PyTorch is because all failing traces have significantly fewer trace lines than passing traces. The consistent difference in length of traces between the two classes allows the NN to easily distinguish them. The ablation study in Tables 4.2 and 4.3 show arguments in traces matter for model performance, while method names and return values are irrelevant.

With Sed, our model achieves 94% precision and recall, and 99% specificity using 10% of the traces in training. The dataset for Sed is unbalanced, with only 18 failing and 352 passing. 10% of the traces in training uses 2 failing tests and 35 passing tests. Given the extremely small sample of failing tests, it is surprising that the model classifies and identifies failing traces with such high precision and recall. To understand this, we examined both the passing and failing trace lines. We find the length of passing and failing traces is similar. All failing traces, however, have a call to a function, `getChar`, towards the end of the trace. This function call is absent in passing traces. We believe associating this function call to failing traces may have helped the performance of the NN. The ablation study in Tables 4.2 and 4.3 show all the parameters considered in our traces are important for model performance.

For Ethereum-CD, our model achieves 80% precision, 82% recall and 79% specificity when trained with 15% of the traces - 169 passing and 169 failing. Ethereum-CD was generated from the reference implementation using an arithmetic operator mutation in a function deeply embedded in the call graph for the difficulty module. Differences between failing and passing traces are not apparent, and analysing longest common subsequence, syntactic diff and manual inspections did not reveal any characteristic feature for failing or passing traces. We believe the model performance of around 80% precision, recall and specificity is due to the similarity between passing and failing traces and the esoteric nature of the mutation. Ablation study for this program reveals that all features in the traces slightly impact model performance.

For COMMONS-LANG, our model achieves 71% precision, 94% recall and 98% specificity using 40% of the traces. This subject program only contains 27 failing executions versus 559 passing executions. There is a stark imbalance between passing and failing traces for this program which impacts the precision achieved. We also observe failing execution traces consist of multiple calls to a string conversion func-

tion, `toString` towards the final parts of the sequence. We find this can serve as a distinguishing feature between passing and failing executions. It is worth noting that our classifier’s performance significantly drops when removing function names in the ablation study and it may be because the `toString` function is no longer visible. In contrast, removing arguments or return values does not affect performance visibly.

For Ethereum-SE, our model achieves 99% precision, 82% recall and 86% specificity with 15% traces in training - 214 failing and 124 passing. Unlike Ethereum-CD, mutation to generate Ethereum-SE was in the core functionality. Failing traces when compared to passing traces had differences towards the end of the trace which is easily distinguished by the NN. Curiously, removing return values in the ablation study, increases recall and specificity. This may be because the model was previously overfitting to return values in traces which may not have been relevant to the classification.

For L7-Filter networking protocols, all programs have enough test inputs to help our model learn program features with a small percentage of execution traces. Especially for ARES protocol with 16066 test inputs, our model can achieve 97% precision, 98% recall and 97% specificity labelling only 3% of the total traces for training. For BGP protocol, we train on 5% of the total traces and achieve 99% precision, 99% recall and 99% specificity. In all networking protocols, failing traces correspond to executions that lead to non-accepted states of the protocol’s finite state machine. We observe that the sequence of function invocations is similar in both passing and failing executions. However, state information in return values and arguments is critical in order to determine correctness. The ablation study supports this argument, as removing function names in any networking protocol has no effect in the classifier’s performance. On the other hand, in all protocols except for BGP, removing arguments dramatically decreases the model’s precision, recall and specificity. Removing return values leads to a slight performance reduction. In BIFF and TELNET, return values seem to be as important as arguments for our model’s accuracy.

Summary. Overall, we find NN models for all our PUTs perform well as a test oracle, achieving an average of 93% precision, 94% recall and 96% specificity. The NN models perform exceptionally well for programs whose traces have characteristic distinguishing features between passing and failing executions, such as differences in trace lengths or presence of certain function call patterns. In the absence of such features, NNs can still do well if it has enough training samples, as in Ethereum-CD. We also find our approach can cope effectively with unbalanced datasets – four of the fifteen programs in our experiment have unbalanced passing and failing traces. Even

though we attempt to explain what patterns lead to our model to succeed or fail, it is not always possible to explain the relationship between the input features and the predicted outcome. Explainability in neural networks is still an open problem.

4.3.7 Q2. Size of training set

Figures 4.6 and 4.7 shows precision and recall achieved by our approach with different training set sizes. The fraction of traces needed in training to achieve near maximal performance was 3% to 40% across the PUTs. Excluding SEAL Encryptor and commons-lang, all the other programs only needed to be trained over 15% of the traces to achieve near maximal performance. Both SEAL encryptor and commons-lang had very few failing traces, requiring a larger fraction of traces to get sufficient representation of failing classes during training. As seen in the plots in Figures 4.6 and 4.7, increasing the % of traces used in training does not increase precision and recall for all PUTs. For instance, Pytorch and Sed observe a dramatic increase in precision and recall when going from 5 to 10% traces in training. Performance, however, stagnates after that point with increasing traces. With Ethereum-CD and Ethereum-SE, precision or recall becomes worse after 20% traces. This maybe because the model is overfitting to the training traces.

It is also worth noting that the absolute size of our training set varies across subject programs. We find our approach works with training sets with as few as 3 failing traces to as many as 214. The range of passing tests in training was between 31 and 169.

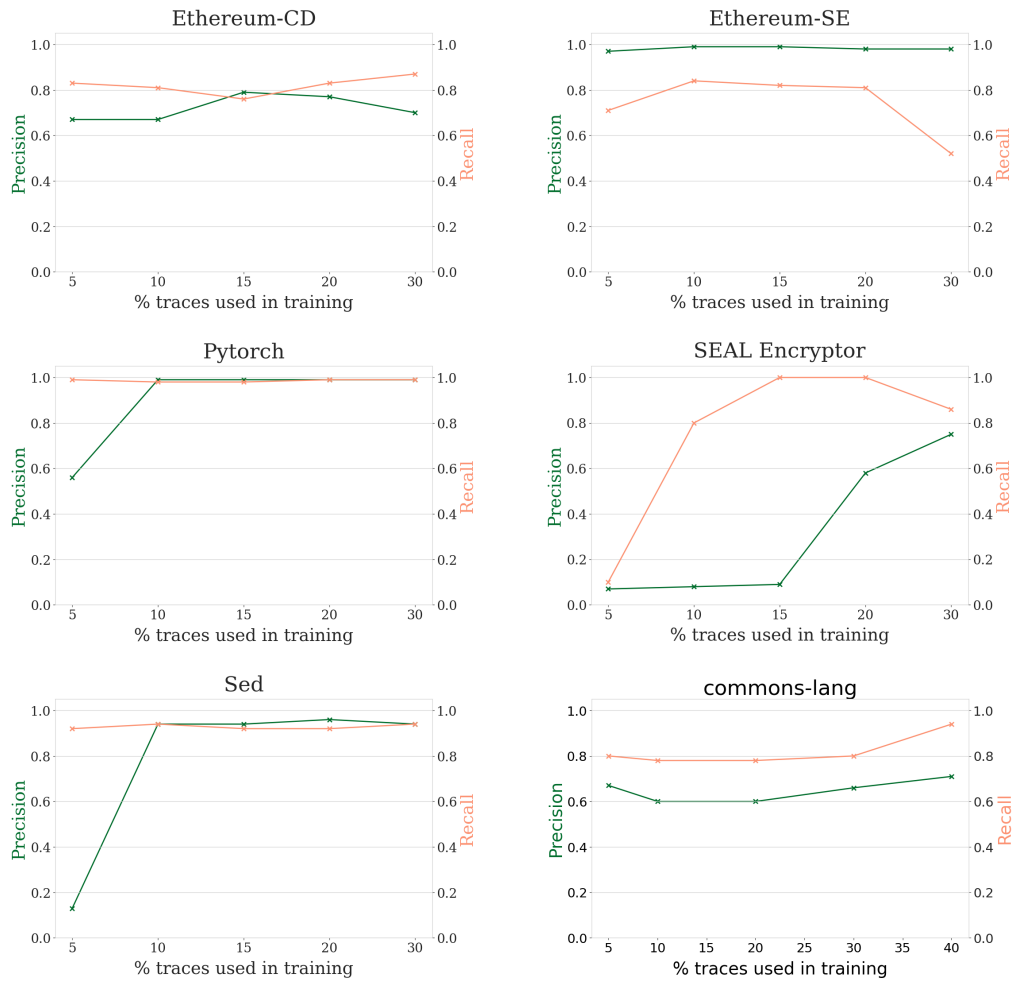


Figure 4.6: Precision and recall achieved by classification model over each PUT.

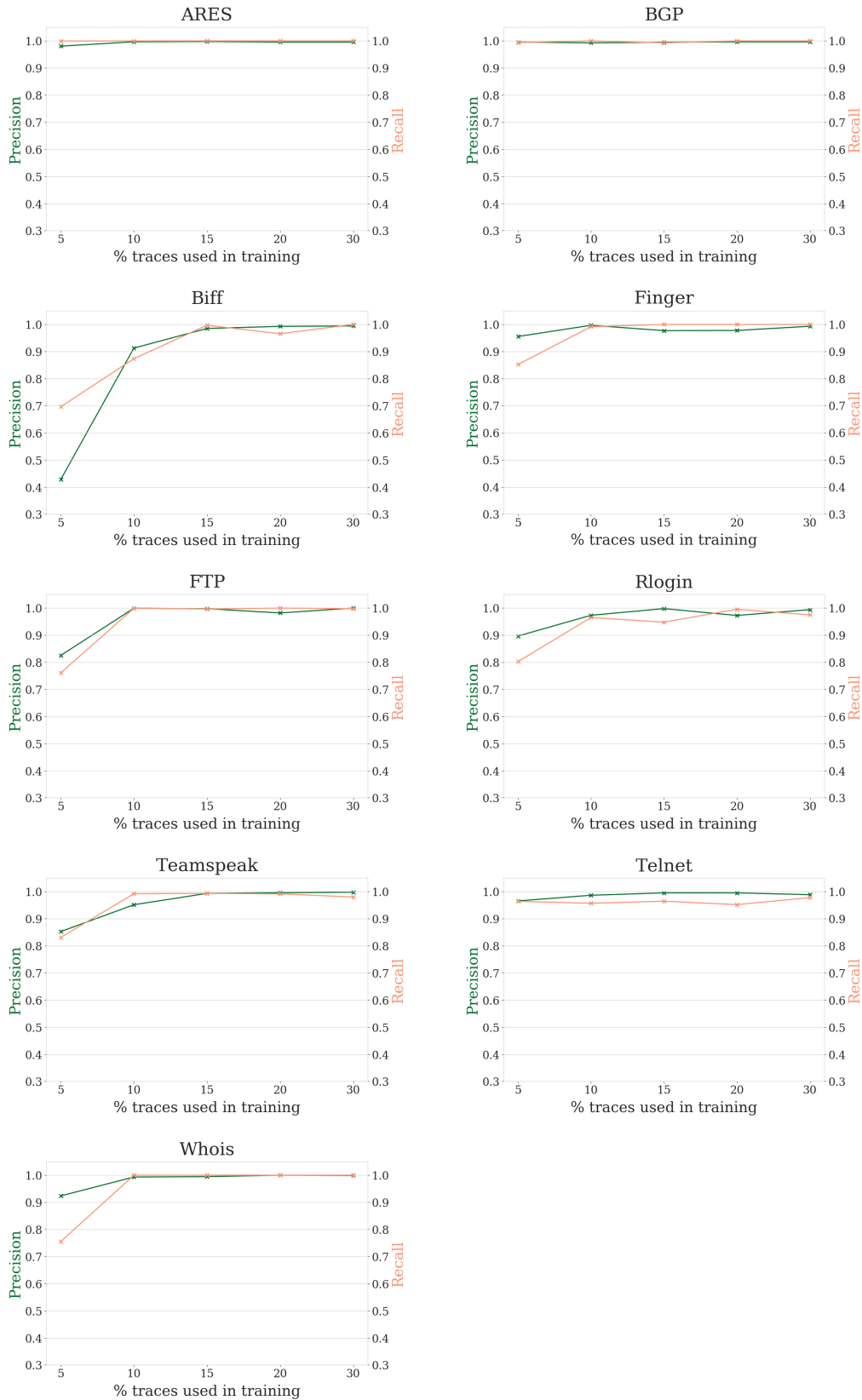


Figure 4.7: Precision and recall achieved by classification model over each PUT.

4.3.8 Q3. Comparison against state of art

Table 4.1 presents precision, recall, and specificity (TNR) achieved by the agglomerative hierarchical clustering proposed by Almaghairbe et al. [11] on each of the PUTs. Comparing the precision, recall and TNR of our approach versus hierarchical clustering, we find our approach clearly outperforms the clustering approach on all but the Ethereum-CD PUT. This is because the hierarchical clustering assumption does not hold for these programs. According to this assumption, passing traces tend to be grouped in a few big clusters and failing traces are grouped into many small clusters. However, for these programs, passing traces tend to be grouped in many small clusters based on their call sequence pattern, making it hard to distinguish them from failing traces by simply comparing cluster sizes.

With Ethereum-CD, the hierarchical clustering approach achieves precision and specificity of 100% and a recall of 49%. This is achieved with complete-linkage clustering, Euclidean distance and a cluster count equal to 10% of total traces. In contrast our approach achieves a precision of 80%, recall of 82% and specificity of 79%. To enable better comparison, we plot the precision-recall curve of the NN model in Figure 4.8 for Ethereum-CD, using 15% of the traces in training.

This curve shows the precision and recall of our trained model with respect to different values of the classification threshold. It is clear from the plot that for the same value of recall (49%), hierarchical clustering performs marginally better than our approach - 100% versus 99%. Hierarchical clustering works well over the Ethereum-CD PUT because the traces are clustered into just one big passing cluster and one failing cluster. Lack of cluster fragmentation improved accuracy of the hierarchical clustering approach. Nevertheless, our model achieves comparable performance for such traces. In addition, our model allows trade off between precision and recall by changing the classification threshold which may be driven by requirements or priorities of the use case. This tradeoff is not possible with the clustering approach.

4.3.9 Q4. Generalisation

In this research question, we conduct an initial exploration into the ambitious possibility of using a model, trained using traces from one subject program, to classify traces from other programs in the same application domain. Figures 4.9 and 4.10 represent precision and recall achieved by models trained using traces from BIFF protocol and WHOIS protocol, respectively, to classify traces produced by other FSMs. We find

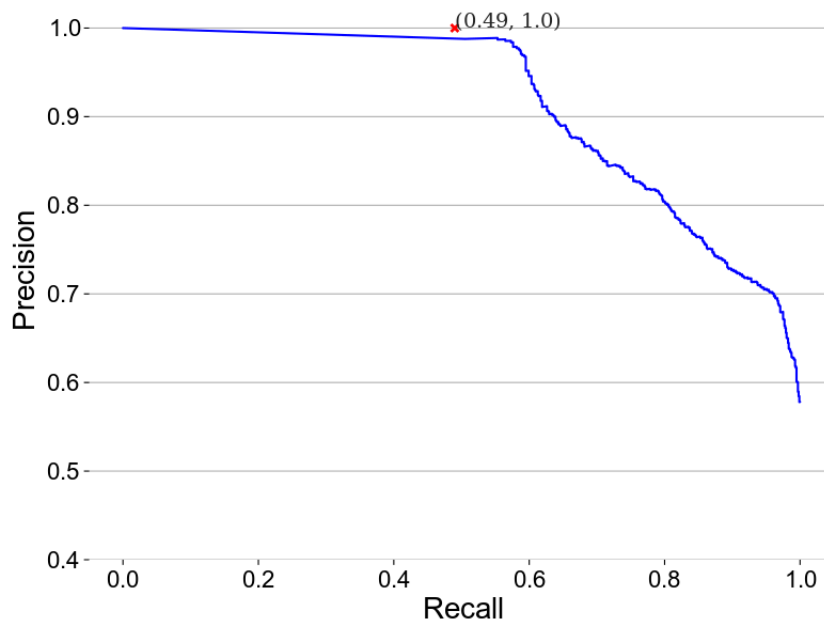


Figure 4.8: Precision-Recall curve for ETHEREUM-CD.

that the model trained using traces from BIFF achieves high accuracy over the ARES protocol with precision and recall close to 1.0, and reasonable precision (> 0.8) for BGP, FTP, RLOGIN, TEAMSPEAK, WHOIS protocols. Lowest precision (0.17) was observed with TELNET. Average precision achieved in classifying traces from unseen FSMs was 0.79. Recall achieved by the model is lower than precision indicating that the model missed identifying failing traces in each of these protocols. Overall, the model trained with BIFF traces was successful in identifying failing traces in other FSMs that have similar patterns to BIFF. Failing traces with differing patterns were missed. We confirmed this observation by checking results from the WHOIS model. Although precision and recall numbers are different from the BIFF model, the reasoning for classification success was the same - extent of similarity in execution patterns between FSMs. With the current approach, we find there is scope to generalise a classification model from a single FSM to multiple FSMs in the networking domain. However, achieving high accuracies with generalisation is a difficult problem and we plan to take small, definitive steps towards addressing this challenge in the future. As a next step, we will explore tuning the classification model from an individual FSM with sample traces from other FSMs to improve generalisation performance.

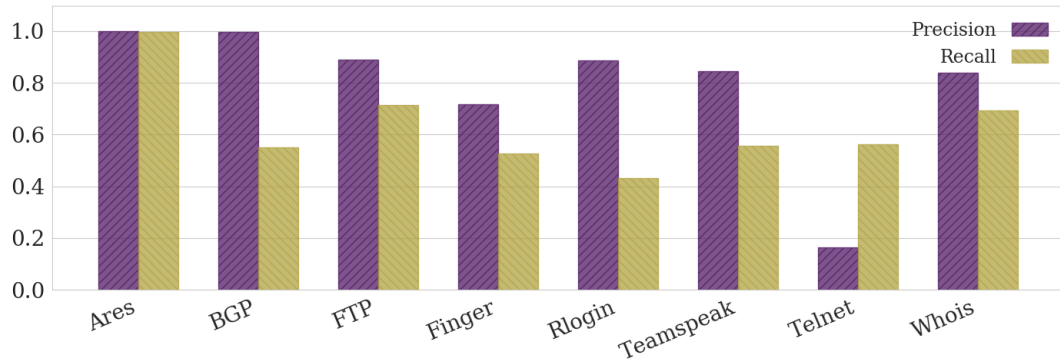


Figure 4.9: Biff trained model - Precision and recall for unseen fsms.

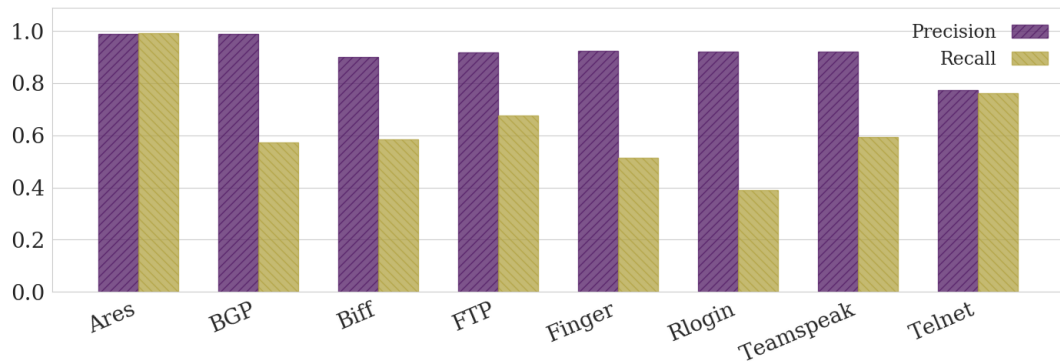


Figure 4.10: WHOIS trained model - Precision and recall for unseen fsms.

4.3.10 Threats to Validity

We see three threats to validity of our experiment based on the selection of subject programs and associated tests.

First, PUTs for 4 out of the 15 subject programs in our experiment were generated by seeding faults into a reference implementation. A reference implementation with only passing tests is not suitable for evaluating our approach. To address this, we generated a faulty implementation and ran the original tests through the PUT to gather both passing and failing traces. It is possible using real faults in place of seeded faults may lead to different results, or introduce faults that are different than those humans make. However, Andrews et al. have shown the use of seeded faults leads to conclusions similar to those obtained using real faults [16, 83]. For one of the subject programs, Sed, we did not artificially seed faults, but instead used the existing implementation as it was accompanied by both passing and failing tests. On a real-life scenario where no real bugs are known, the developer needs to introduce seeded faults into a tested codebase to help the neural network learn the distinctive features between the passing

and failing class. According to [16, 83] such mutations would lead to similar behaviour with real bugs.

Second, the number of tests that accompanied our subject programs was not very large, ranging from 132 to 16066 tests. The NN models in our experiments produced good performance with small to medium sized test suites that may be automatically or manually generated. Our approach is constrained by the amount of training data and not by the size of the test suite. As a result for programs accompanied by large test suites, the NN model will need a larger training set (fraction of traces to be used in training might still be 14%). Depending on the total size of the test suite, the percentage needed to be labelled by an expert could still be large. Nevertheless, the labelling effort (i.e. computing manually the expected value for each test input by understanding a codebase's functionality) for a fraction of the tests in our approach is still less than the current practice of labelling all the tests.

Finally, we conducted our study on subject programs from 5 different application domains which is not representative of all application domains. Different application domains may have ranging programming paradigms, therefore different patterns in errors within execution traces that may be easier or more difficult to predict. For example, a parser program is expected to have many recursive function calls whereas a video game a complex interaction between different classes. It is not guaranteed that a neural network will learn to represent each of them with the same accuracy. Given that our approach has no domain specific constraints, we believe it will be widely applicable.

4.4 Summary

In this Chapter, we describe a novel approach for designing a test oracle as a NN model, learning from execution traces of a given program. We have implemented an end to end framework for automating the steps in our approach

1. Gathering execution traces as sequences of method invocations.
2. Encoding variable length execution traces into a fixed length vector.
3. Designing a NN model that uses the trace information to classify the trace as pass or fail.

We support Java programs in addition to C/C++. In addition, we conduct an extensive evaluation using 15 realistic PUTs and tests. We found the classification model for each PUT was effective in classifying passing and failing executions, achieving an average of 93% precision, 94% recall and 96% specificity while only training with an average 14% of the total traces. We outperform the hierarchical clustering technique proposed in recent literature by a large margin of accuracy for 14 out of the 15 PUTs, and achieved comparable performance for the other PUT. We did an initial experiment with generalising a classification model learned over one protocol FSM to classify executions over other network protocol FSMs. The results for precision and recall over other unseen FSMs was not as high as the individual FSM classification models. In the future, we plan to explore techniques that will improve the generalisation performance of the NN models.

Practical use. Our approach can be applied out of the box for classifying tests for any software that can be compiled to LLVM IR or Jimple IR. We gather execution traces for test inputs automatically, and require a small fraction of the traces to be labelled with their pass or fail outcomes (average 15% in our experiments). The remaining traces will then be classified automatically. Our approach is promising with high accuracy and has clear benefits over current industry practices where developers label *all* the tests. Our future work will focus on methods to improve the classification accuracy while reducing the training data requirement using techniques like transfer learning.

Replacing a PUT with the Test Oracle. Our approach learns to classify program executions as “passing” or “failing” by extracting representations from execution traces. Implicitly, this requires the test oracle to reason about a program’s reference result. Given the test oracle has this information internally it could replace the PUT and act as the actual output provider, given a test input. This could happen with a different NN design, training and labelling process (e.g. regression instead of binary classification) and is out of scope for this work. Our tool does not aim to replace a reference program. Instead we aim to provide the testing expert with a tool that will assist them in labeling test executions. A NN to provide reference outputs for a PUT would be helpful in cases where compiling and executing programs would be slow and acquiring quickly an output with a slight loss of accuracy is acceptable.

Chapter 5

BenchPress: A Deep Active Benchmark Generator

5.1 Introduction

The research contributions presented in this Chapter are the design and implementation of a deep learning, unsupervised generative model for compiler benchmarks that uses active learning to search compiler feature spaces for important features. In the following Sections, the approach, implementation and evaluation of this technique are described.

We develop BENCHPRESS [53], a BERT-based OpenCL benchmark generator [46, 139] that targets and synthesizes benchmarks in desired parts of the feature space. We use active learning to choose parts of the feature space and beam search to steer BENCHPRESS’s generated samples towards the requested features. We train BENCHPRESS with OpenCL code samples that we collect by mining BigQuery [60] and GITHUB directly using its API [58]. We support composite data types and calls to user-defined functions in our dataset and benchmark generation. BENCHPRESS is a bidirectional generative model and learns to generate code in any part of a sequence by jointly considering left and right context. We achieve this with a new learnt token, the [HOLE], which hides a sequence from the input, whose length is unknown to BENCHPRESS during training. BENCHPRESS learns to fill [HOLE] by iteratively predicting an arbitrary number tokens that are likely to lead to a compiling function.

BENCHPRESS outperforms CLGEN in the task of undirected program generation from a fixed input feed, generating $10\times$ more unique OpenCL kernels that are $7.5\times$ longer on average, with a compilation rate of 86% compared to CLGEN’s 2.33%.

BENCHPRESS strongly outperforms benchmark synthesizers CLGEN, CLSMITH [164, 103], and human written code from GITHUB in reaching close to the features of Rodinia benchmarks, developed by compiler experts. Finally, BENCHPRESS uses active learning, specifically query by committee [135], to search the feature space and find missing features to improve Grewe’s et al. [61] CPU vs GPU heuristic. Enhancing the heuristic’s dataset with BENCHPRESS’s benchmarks improves the heuristic’s speedup relative to the optimal static decision by 50%, increasing it from 4% to 6%, when the maximum possible speedup for this task is 12%.

In this Chapter, we present the following contributions:

1. We are the first to develop a feature-space agnostic, steerable code generator towards desired program features.
2. We develop an automated approach to rank the feature space of downstream tasks with active learning.
3. We enable bidirectional source code generation by inserting [HOLE] tokens in any part of a sequence.

5.2 Motivation

Figure 5.1 shows a two-dimensional slice of the Grewe’s et al. [61] feature space: number of computational instructions vs number of memory instructions. Figure 5.1 also shows how the OpenCL benchmarks found in the Rodinia suite map into this plane, represented as purple diamonds. We select to plot the benchmarks’ distribution with respect to these two features because they are accurate in characterizing the type of a subject OpenCL workload. Having enough benchmarks to cover lots of different ratios of computational to memory instructions is desirable for predictive models that generalize well. What we find is that much of this two dimensional space is uncovered. 54 of the 58 Rodinia examples cluster in the lower left corner, the rest of the space having only four examples. Any optimization decision for programs in this area of the space would not be accurate due to lack of representative examples.

CLGEN attempted to address this problem by automatically generating more training examples. However, the generated kernels lacked feature diversity and provided even poorer coverage of the feature space. Figure 5.1 represents their position in the 2D space as red dots. Almost all of them are concentrated in a corner covering a small

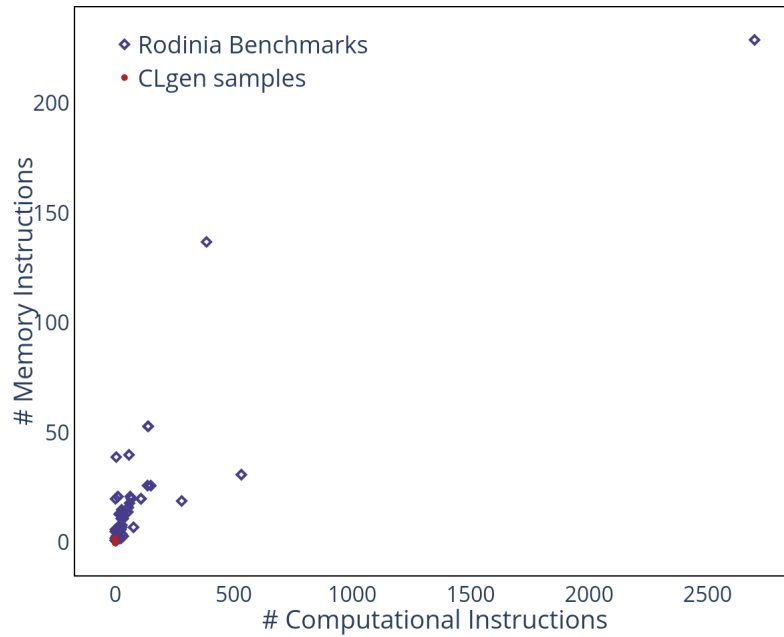


Figure 5.1: # Memory operations and # computational instructions for **(a)** Rodinia benchmarks in purple diamonds and **(b)** CLGEN’s samples in red dots. Generating samples with missing features is vital for predictive modeling’s performance.

percentage of the feature space. While CLGEN can generate hundreds of millions of unique kernels, almost all of them will fail to compile. As the probability of having at least one illegal token in the kernel body increases with the number of tokens, only tiny kernels are valid. In our experiments in Section 4.3.5, the longest compiling CLGEN kernel had 8 lines and 102 tokens. Given the small number of tokens in valid kernels, there is a high degree of repetitiveness in the generated corpus, not only in terms of features but also in terms of structure and functionality. As a result, this approach is not well suited to augmenting the training set with diverse feature benchmarks. There is a compelling need to generate training points for uncovered regions of the feature space and we attempt to address this need with BENCHPRESS. In the following Sections, we discuss our approach and evaluation of BENCHPRESS, comparing it to the existing state-of-the art for feature space coverage.

5.3 Approach

We present BENCHPRESS, a deep learning model for directed compiler benchmark generation. BENCHPRESS is the first steerable synthesizer that can synthesize compiling functions with target features. BENCHPRESS steers its generation with a feature

space-agnostic beam search algorithm to search the space and steer BENCHPRESS’s generation towards the target features. Given a downstream task, BENCHPRESS learns what features to target in order to improve its performance by searching the space with active learning. BENCHPRESS’s language model is based on BERT [46], which we transform into a generative model.

The key feature in BENCHPRESS that enables bidirectional code generation is a new token, namely, the [HOLE] token. Compared to traditional sequence to sequence generative models that can only add tokens at the end of a sequence, bidirectionality is an important feature that allows BENCHPRESS to apply edits on previously generated code and also insert tokens that fit into the general context of an existing function. We train BENCHPRESS to learn and understand how to iteratively fill holes of unknown length that can be found in any part of an input sequence by conditioning it on the left and right context of the [HOLE]. Later, this approach enables us to use beam search and steer benchmark generation into the feature space iteratively as it can regress to previously generated benchmarks with new holes and produce newer samples with better features.

Figure 5.2 illustrates an overview of our approach. BENCHPRESS consists of three main components:

1. Learning corpus collection and processing.
2. Source code language modeling.
3. Feature space search and benchmark generation.

We discuss each step in the following Subsections.

5.3.1 Learning Corpus

Modeling source code accurately requires large amounts of data [100] similarly to other deep learning tasks. We develop a tool to collect data from BigQuery’s GITHUB dataset [60]. We also use GITHUB’s API [58] and mine directly extra repositories that are not included in BigQuery. We choose OpenCL to train and evaluate BENCHPRESS for several reasons. First, many performance critical workloads are still being developed in OpenCL and ensuring optimal performance is important. Second, OpenCL is relatively simple to compile and execute compared to more generic languages (e.g. C/C++) because it has limited dependencies to third-party libraries and header files.

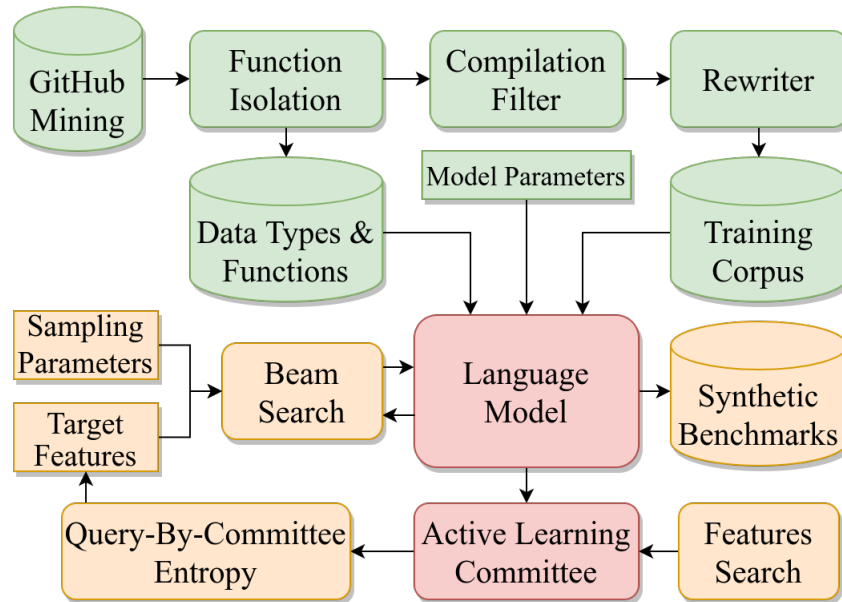


Figure 5.2: BENCHPRESS’s high-level approach. We highlight the corpus collection and processing in green, the language modeling for source code in red and the feature space search for benchmark generation in orange.

Finally, two easily accessible baselines that we can compare against are trained for OpenCL generation. Even though we use OpenCL in this contribution, BENCHPRESS can work without any modifications for any programming language. However, a more appropriate tokenizer for the desired language would make the process more efficient. Also, a driver that satisfies the language’s library dependencies, compiles and executes the generated code has to be developed.

There are a few innovations in how we pre-process the code compared to previous works. First, we inline included header files recursively into source files to resolve type dependencies. Additionally, we automatically extract custom data types (e.g. `STRUCT`, `TYPEDDEF`) and utility functions found in the unprocessed corpus and place them into header files that are accessible throughout BENCHPRESS’s pipeline. This way, we resolve most type dependencies by retaining the functionality and semantics of the original, human-written programs. These two steps enable us to increase significantly the amount of compiling kernels we end up with in our training dataset. Second, we isolate kernels into single instances because BENCHPRESS is trained on complete functions. From the previous steps, the type dependencies of each kernel are known and we automatically provide them to the compiler, retaining their compilability. Finally, we compile all kernels with Clang and reject those that do not compile.

Next, we re-write identifiers by randomly sampling the alphabet, eliminating spuri-

ous naming patterns in the corpus. All kernels are padded to BENCHPRESS's sequence length and kernels that are longer than this are truncated to fit. This helps BENCHPRESS train its later indices' positional embeddings more effectively, for which we have less training information compared to earlier indices. Finally, we derive a tokenizer by parsing the AST of all source code. We reserve tokens for all OpenCL keywords and all intrinsic OpenCL function name identifiers found in the official OpenCL specifications [121]. We analyze the dataset and tokenize by word the most common function names and custom data type identifiers that we have collected. We encode all literals and infrequently used custom types and functions character by character to avoid exploding the size of the vocabulary. We define 5 meta tokens: [START], [END], [PAD], [HOLE], [ENDHOLE]. The derived tokenizer holds in total 2,201 unique tokens.

5.3.2 Language Modeling

BENCHPRESS is based on BERT [46], a Transformer-based model originally designed for natural language modeling. BERT is trained to predict words that have been randomly hidden by [MASK] tokens. This way BERT learns fitting words with respect to their position in a sequence and also the left and right context, i.e., the text sequence before and after the masked token to be predicted. This type of training helps BERT learn what words mean within a given context, improving downstream tasks that rely on that knowledge.

While this is a useful property, it is not enough to turn BERT into a generative model. We also want to be able to extend a kernel by inserting an arbitrary number of tokens in arbitrary positions. We could iteratively add a [MASK] token to get one extra token at a time, until we have a full statement. This would be limiting. Each time the new token would be selected based on its probability of completing forming a compilable kernel. Every intermediate kernel in the iterative process would have to be compilable or similar to a compilable, which is not a general way for augmenting kernels.

Clusters of [MASK] tokens could allow us to insert multiple tokens in each iteration. This is still unsatisfactory. The number of [MASK] tokens in the cluster biases the kind of code that will be generated: if we ask such a generator to produce five tokens, it will give us a five token statement that could be expected to close this gap, not a five token sequence that could be the start of a much longer statement. We could

place the left and right context to the edges of a sequence and fill intermediate positions with [MASK] tokens. BENCHPRESS could predict a vocabulary or a stop token for a [MASK], allowing for arbitrary sequences. We test this configuration and sample a trained model with a fixed input feed. BENCHPRESS is unable to learn the [MASK]s' left and right context conditionally, when many [MASK]s are in a sequence, which leads to zero samples to compile or even resemble reasonable code.

What we do instead is to extend BERT's functionality with a new pair of learnt tokens, the [HOLE] and the [ENDHOLE]. [HOLE] follows the same logic with [MASK], however the number of tokens that have been hidden behind it is unknown to the model during training. The model only learns to predict the first token of an arbitrarily long missing sequence. At inference-time, we iteratively predict the first token of the remaining sequence and re-insert it just before the [HOLE]. This way BENCHPRESS learns to generate arbitrarily large code sequences within any part of a sequence.

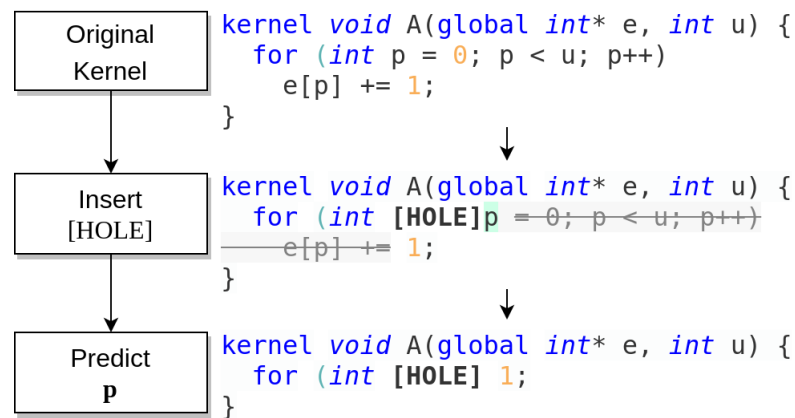


Figure 5.3: When a [HOLE] is inserted to a kernel at a random index, it hides a random number of tokens, unknown to BENCHPRESS. On this example, BENCHPRESS learns to predict the first hidden token, **P**.

Figure 5.3 shows how a [HOLE] is inserted into a function to create a datapoint. A random starting index and a random length are selected. The choice of index and length are only restricted by a potential overlap of the prospective hidden sequence with any of the other meta token or the maximum hole length that is defined as a training parameter for the architecture as a percentage of each function's length. When the specifications of a hole have been settled, the hidden sequence is discarded. Only the first token of it is kept as the target prediction for that hole. A hole can also represent an empty sequence, i.e. hiding 0 tokens. In this case, the target prediction during training is

[ENDHOLE]. The training instances are randomly generated on demand, the entire space of possible instances is too large to be pre-generated. In this work, we only insert 1 hole per training instance for BENCHPRESS to learn. Multiple holes could be used during training, but this is not needed during BENCHPRESS’s current benchmark generation task.

5.3.3 Benchmark Generation

BENCHPRESS’s synthesizer operates as a generative model with the help of [HOLE] / [ENDHOLE] tokens. It receives an input with 1 or more [HOLE] tokens and returns a completed benchmark. For each [HOLE], BENCHPRESS predicts one token that fits in the sequence at the [HOLE]’s index, with respect to its left and right context. If the predicted token is not [ENDHOLE], it moves the [HOLE] and all subsequent tokens one position to the right and inserts the predicted token to the initial target index. This intermediate kernel is iteratively provided as an input for the next token prediction and the process is repeated until BENCHPRESS predicts [ENDHOLE]. This marks a [HOLE] is complete and the final sample is returned, as shown in Figure 5.4.

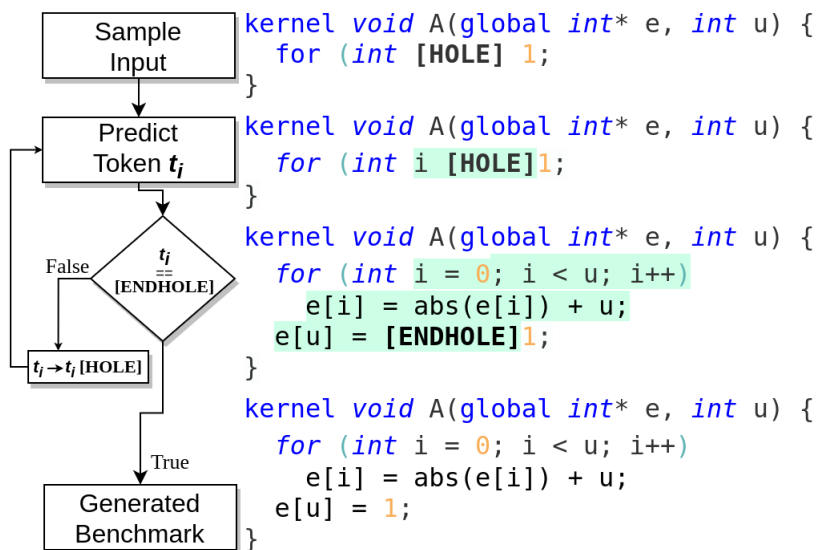


Figure 5.4: During sampling, BENCHPRESS receives an input and predicts iteratively the fitting tokens. BENCHPRESS predicts [ENDHOLE] to indicate a [HOLE] is complete.

On its own, this process only augments kernels. We also make it the first to target desired parts of a feature space by repeatedly generating kernels, selecting the ones closer to the target features, inserting new holes, and generating new augmented kernels. We use beam search to steer generation.

Given a target feature vector, BENCHPRESS samples a starting, fixed input feed `'kernel void [HOLE]'` and yields a collection of starting benchmarks. We reject benchmarks that do not compile and for the remaining we measure the Euclidean distance between their feature vectors and the target features. We select the *top-K* candidates that have the shortest distance from the target and we use them as inputs for the next generation. To improve diversity among promoted benchmarks we introduce randomness in the selection of *top-K* candidates: Each *top-K* sample, has a fixed probability $p = 0.15$ to be replaced by another random candidate of its generation. BENCHPRESS lazily creates multiple different input instances for each selected candidate by placing a random [HOLE] of random length in order to synthesize a new sample. BENCHPRESS generates a successive collection of benchmarks, of which K compiling ones with the shortest distance from the target again are selected with p -randomness and used as inputs. This search continues until a sample achieves a distance of 0 from the target, or until a threshold of generations (i.e. beam search depth) is exhausted. BENCHPRESS returns the closest benchmark to the target's features along with all beam search's intermediate benchmarks that cover the model's traversal of the feature space starting from the origin and ending near the target features. For the benchmark synthesis process, we use categorical sampling with temperature to sample BENCHPRESS's probabilities. The sampling temperature, beam search's width K and depth are defined as sampling parameters.

In the worst case, BENCHPRESS's directed program generation is slow, ranging from a few seconds to one hour, as it typically requires thousands of language model inferences. However, BENCHPRESS is the first program synthesizer that can target a set of desired program features.

5.3.4 Feature Space Search

A steerable synthesizer allows the generation of benchmarks with desired features. However, the automatic selection of those parts of the feature space that are worth targeting is challenging and depends on the downstream task.

BENCHPRESS attempts to solve this by searching the feature space with query by committee [135], a well-known active learning technique. We implement a committee of (a) 7 NN, (b) 7 k-NN and (c) 7 K-means models. We set their initial state by passively training on a small portion of the downstream task's data. We sample the committee with thousands of random points in the space, we collect the predicted

labels and measure the entropy for each sample. The entropy shows the level of uncertainty among the committee about the predicted label of a given point and is defined as:

$$H = - \sum_{l \in L} (p(l) * \log(p(l))) \quad (5.1)$$

where L is the set of all predicted labels and $p(l)$ the probability of label l in the committee’s prediction set for a given input. The highest entropy point is an important feature vector to target and BENCHPRESS steers benchmark generation towards it with the approach explained in 3.3. We collect the labels of generated benchmarks and we train incrementally the committee with them. Then, we sample it to find the next highest entropy point. We continue this process until we saturate the feature space. BENCHPRESS’s committee is agnostic to the downstream task or the feature space and its I/O dimensions are hyper-parameters selected with respect to the task’s feature and prediction dimensions.

Our active learning algorithm determines the useful range of a downstream task’s feature space over multiple iterations. Each part of the feature space that is targeted is likely to improve the task’s overall accuracy. By design, BENCHPRESS does not introduce randomly scattered datapoints to an application dataset’s feature space, therefore no irrelevant benchmarks are added. Once no high-uncertainty areas are found (based on entropy), the active learning process halts.

5.4 Experimental Setup

We describe the configurations used in training BENCHPRESS, and the parameters used in evaluation, namely (1) Feature Spaces - we use three different representations of program features, (2) Target Benchmarks - We use Rodinia benchmarks [33] and their features as the target for synthesis by BENCHPRESS, (3) Comparison to state of the art - we compare BENCHPRESS with code synthesizers and human written code in improving Grewe’s et al. heuristic model.

5.4.1 Platforms

We train BENCHPRESS and conduct all our experiments on two 64-bit systems each having one Intel Xeon E5-2620 16-core CPU, 2x Nvidia GeForce GTX 1080 GPU and 32 Gigabytes of RAM. We use Ubuntu 18.04, PyTorch 1.9.1 [123], CUDA version 11.4

and Nvidia driver version 510.47.03. We use Clang-10 as BENCHPRESS’s compiler and LLVM-10 to compile and execute InstCount and Autophase [65] extracting tools. For compatibility reasons, we are required to use Clang LibTooling from LLVM-6 to execute Grewe’s et al. [61] feature extractor.

5.4.2 Language Modeling for source code

We collect OpenCL code from GITHUB and split it into single function instances. We ensure no kernels that come from benchmarks suites used in the evaluation are included in our corpus by excluding the repositories or their forks that contain them. We do this to avoid the model train on benchmarks that will be requested to target during the evaluation. We pre-process text, re-write variables and reject OpenCL kernels that do not compile. In total we mine 63,918 OpenCL kernels across 12,860 GITHUB repositories and we successfully compile 19,637 of them (31% compilation rate).

We train BENCHPRESS on our OpenCL Corpus for 10M steps with a batch size of 32. For BENCHPRESS’s BERT model parameters, we select 2 hidden layers, 12 attention heads. We set intermediate size, hidden size and max position embeddings (i.e. sequence length) to 768. We set the sequence length to 768 as it is the largest length our hardware can support while keeping the model’s training time to a sustainable level. We set the maximum length of holes to be 90% of a kernel’s token length, i.e. a hole can hide almost all tokens of a training instance. We optimize the model using Adam optimizer with a learning rate that reaches a maximum of 45×10^{-6} after 20,000 warmup steps and decays linearly over the remaining training steps. We train BENCHPRESS’s language model to a final loss value of 0.28. The selection of hyper-parameters and number of training epochs are decided after validating our candidate models in generating kernels and trading-off the synthesis quality against the amount of hardware needed to sample them.

5.4.3 Feature Spaces

Compiler predictive models use static code features to represent programs and learn optimisation heuristics. A vector of independent characteristics represent a single program. Each of them are typically an integer or float value. Features are extracted at the Syntax level by traversing the AST or at the IR level using the compiler’s middle end (e.g. LLVM-IR). A feature space is the collection of all possible program feature vectors.

BENCHPRESS is a generative model that can be steered to generate samples for a desired part of the feature space. We evaluate BENCHPRESS on three source feature representations we find across the literature, (a) Syntax-level Grewe’s et al. features [61], (b) IR-level LLVM-InstCount [98] and (c) IR-level Autophase [65].

Grewe’s et al. features are extracted with Clang’s LibTooling and used to train their predictive model on the CPU vs GPU task for OpenCL kernels. This feature space holds 8 dimensions. 4 dimensions describe the number of 1) computational, 2) relational, 3) atomic and 4) memory access instructions. The feature space also counts the different type of memory instructions, local memory or coalesced. Finally, the computational to memory and coalesced to memory ratios are defined.

InstCount is a standard pass provided by LLVM-IR framework and used in Compiler Gym by Cummins et al. [41]. InstCount holds 70 dimensions: 67 dimensions each counting all 67 LLVM-IR instruction types and total number of 1) instructions, 2) basic blocks and 3) functions. Autophase by Huang et al. [65] holds 56 dimensions. While many of the features used in Autophase are shared with InstCount, they introduce new ones such as number of input arguments to PHI Nodes or total number of memory instructions. On the other hand, they do not include the count of some LLVM instructions that are not considered to contribute to a program’s representation, e.g. CATCHPAD instruction.

In this work, BENCHPRESS targets both static and runtime features representing programs. InstCount and Autophase feature spaces include only static features. Grewe et al. define two runtime features: the local size of a kernel and the amount of transferred bytes during execution. For both types, the targeted benchmark generation for BENCHPRESS does not change. The language model produces candidate benchmarks whose features are extracted. When the feature space depends only on static features, only a compiler pass is typically required. For runtime features, the feature extractor needs to drive the kernel to collect them. This adds an extra time overhead.

BENCHPRESS is also agnostic to the structure of targeted features. The feature spaces we use consist of fixed arrays of numerical features that are easy to navigate to with Euclidean distance. Other features could be in the form of graphs, e.g. control flow or data flow graphs of programs. In such features, an embedding model is needed to convert structured data into numerical representations to enable BENCHPRESS’s steerable generation. Pre-trained ML models are typically used to embed information into vectors. For example, Graph Neural Networks can be used for graph-like data or Transformers receiving a flattened view of graphs as a sequence.

5.4.4 Analysis of BENCHPRESS and CLgen language models

CLGEN [40] is the current state of the art in OpenCL benchmark generation. Its synthetic benchmarks improve the accuracy of Grewe’s et al. predictive model [61] by $1.27\times$. However, Goens et al. [59] perform a case study and show evidence that CLGEN’s synthetic benchmarks do not improve the quality of training data and, consequently, performance of predictive models. They show that a predictive model in fact performs worse with synthetic benchmarks as opposed to human written benchmarks or code from GITHUB.

This study motivates us to perform an analysis of BENCHPRESS’s language model, BERT, with CLGEN in the task of undirected program generation. In this first experiment, we reproduce CLGEN using the authors’ artifacts and we sample it with a fixed input `'kernel void'` to collect a dataset of unique OpenCL kernels. We use BENCHPRESS on the same generative task and sample the model with the same fixed input `'kernel void [HOLE]'` to obtain another dataset of unique benchmarks. In this experiment we focus on the language model’s inference performance. We compare both generative models on their throughput, their ability to create compiling code, feature distribution and code size. In this experiment, we do not direct program generation. BENCHPRESS generates compiling kernels in a single inference step.

5.4.5 Targeted Benchmark Generation

Next, we evaluate BENCHPRESS’s ability to steer towards desired program features. We use well-established compiler benchmarks as our reference and target their features within this space. These benchmarks usually perform intensive operations, such as matrix multiplications or FFT analysis, they contain hundreds of computational and memory instructions and are specifically fine-tuned by experts to exercise compilers from different angles. As a result, we believe features in these benchmarks provide a good target to assess performance of BENCHPRESS’s ability to target complex features.

We choose target benchmarks within the Rodinia suite [33, 22] as it is widely used in the literature [40, 42]. Similar to the training corpus, we collect the suite’s source files, we inline header files and dependent OpenCL libraries into them, we split kernels into single source files and reject those that do not compile. In total, we collect 61 target Rodinia benchmarks out of which 58 compile. For the remaining benchmarks, we collect their features using the feature extractors for Grewe’s et al., Inst-Count and Autophase feature spaces [61, 98, 65]. We target the feature vectors of

these benchmarks and request BENCHPRESS to generate at least one matching benchmark for each. We end up with three collective synthetic benchmark datasets, one for each feature space, that contain code with features matching Rodinia benchmarks. For each Rodinia benchmark’s target feature vector, we measure the minimum Euclidean distance to it achieved between BENCHPRESS, code from GitHub, CLGEN and CLSMITH [164, 103]. For GITHUB’s and CLSMITH’s kernels, we use SRCIROR [67] to apply code mutations exhaustively with beam search.

To make our experiment more intuitive we use two datasets for GITHUB: a) GITHUB consisting of all OpenCL kernels we collected and b) GITHUB-768, a proper subset of GITHUB which contains only the kernels that do not exceed BENCHPRESS’s sequence length of 768 tokens. Since BENCHPRESS benchmarks’ size are restricted to the architecture’s sequence length, we feel it is important to make this distinction in order to present a view of BENCHPRESS’s actual performance on features that may be unreachable within the current sequence length. For example, it may be impossible to generate 2,000 computational instructions within 768 tokens. For such cases, we believe GITHUB-768 with its equally restricted sequence length would allow for a fairer comparison.

For all three feature spaces, we weed out the Rodinia benchmarks that have an exact matching sample (i.e. a Euclidean distance of 0) in GITHUB-768. Since we already have matching samples for them, we do not need to target them with BENCHPRESS or any other generative model. However, we do not skip benchmarks whose features exist only in GITHUB’s full dataset as we wanted to explore the feasibility of using BENCHPRESS to generate a sample with the same features but smaller sequence length. Applying this restriction we end up with 22 Rodinia benchmarks for Grewe’s et al., 52 for InstCount and 36 for Autophase feature spaces.

5.4.6 Active Learning for Feature Selection

BENCHPRESS’s steerable generation is vital for searching the feature space while also finding useful features to target with active learning. In this experiment, we evaluate BENCHPRESS in the downstream task of training the predictive model proposed by Grewe et al. [61], a well-tested problem used by many baseline models.

Grewe et al. train a decision tree model to predict the optimal device to execute a benchmark, choosing between a CPU and a GPU. They measure their model’s performance as speedup achieved with using the predicted device for execution versus

statically executing all benchmarks on the GPU. To train the predictive model, they use OpenCL benchmarks from 7 well-known benchmarks suites [40, 61]. In this experiment, we reproduce Grewe’s et al. heuristic using their artifact and we also retrain it with datasets enriched with executable benchmarks from BENCHPRESS using active learning and passive learning (i.e. targeting random parts of the feature space instead of searching it), CLGEN and GITHUB. We measure the speedup over static mapping for each of them.

To collect our evaluated datasets, we execute OpenCL benchmarks with CLDRIVE [40] by Cummins et al. CLDRIVE automatically generates inputs and drives kernels to the hardware. It measures the execution time per device across thousands of runs and it rejects kernels that produce runtime errors, do not modify any of the inputs (no output) or modify them differently for each run (not deterministic). For (a) the 7 human-written benchmarks suites, (b) BENCHPRESS, (c) CLGEN and (d) GITHUB, we execute their kernel on CLDRIVE using a range of different *local* and *global size* configurations. We label each instance with the fastest measured device (the CPU or the GPU), in the same way Cummins et al. [40] and Grewe et al. [61] performed their evaluation.

5.5 Results

In this Section, we show our experiments’ results and compare BENCHPRESS with state of the art techniques in OpenCL benchmark synthesis. We present case studies of (a) BENCHPRESS’s throughput as a generative model compared to CLGEN, (b) its ability to steer benchmark generation towards desired features and (c) its performance in searching the feature space to enhance a downstream task’s performance.

5.5.1 Analysis of BENCHPRESS and CLGEN language models

We perform an analysis of BENCHPRESS and CLGEN as language models and compare them in generating a collection of benchmarks from a fixed input feed, ‘kernel void [HOLE]’ and ‘kernel void’ respectively. We compare the two approaches measuring (a) the generative models’ throughput and (b) the quality of their generated benchmarks in terms of code size and features. In this experiment, we do not use any directed search or iterative approach for BENCHPRESS’s generation. We perform this evaluation to measure how BERT, BENCHPRESS’s underlying language model,

compares with CLGEN as a generative model. Table 5.1 presents the aggregate measurements for the generated benchmarks using both approaches.

	# unique bench	# compiling bench	comp rate	max tokens	max inst (LLVM)	ms per sample
BENCHPRESS	190,460	142,607	86%	750	161	162
CLGEN	1,564,011	13,035	2.33%	102	32	103

Table 5.1: Throughput comparison between BENCHPRESS and CLGEN on generated OpenCL benchmarks when BENCHPRESS does not use feature-directed program generation. The column acronyms are as follows: (a) number of unique benchmarks, (b) number of compiling benchmarks, (c) rate of compilation, (d) largest compiling sample in tokens, (e) largest compiling sample in LLVM-IR instructions (−01) (f) and inference time per sample in ms.

Compilation rate and code quality. BENCHPRESS generates over $10\times$ more unique compiling benchmarks than CLGEN. This result is observed despite BENCHPRESS generating $8\times$ fewer unique benchmarks than CLGEN. The compilation rate with BENCHPRESS is 86% while CLGEN has an exceedingly small rate of 2.3%. BENCHPRESS’s largest sample is 750 tokens compiling to 161 LLVM-IR instructions. This is a $7.5\times$ and $5\times$ increase in number of tokens and number of LLVM-IR instructions compared to CLGEN’s largest kernel. The only drawback of BENCHPRESS compared to CLGEN is that it is considerably slower in generating candidates. This is because the transformer-based architecture in BENCHPRESS is significantly larger in number of parameters than CLGEN’s LSTM. Additionally, BENCHPRESS tends to generate longer kernels than CLGEN, necessitating more inference steps and longer generation time.

In Figures 5.5a and 5.5b, we show the frequency distribution of the number of tokens and number of LLVM-IR instructions for compiling kernels for both datasets. To visualize our results better, we focus on synthesized kernels with token lengths ≤ 100 and instructions lengths ≤ 25 where the vast majority of benchmarks are found. Most of BENCHPRESS’s benchmarks are found to have 20 to 80 tokens and 3 to 16 LLVM-IR instructions. The majority of CLGEN’s benchmarks are found to have 5 to 45 tokens and only up to 4 LLVM-IR instructions. 94% of CLGEN’s generated benchmarks have only 1 instruction when compiled to LLVM-IR. We analyze the dataset to explain this phenomenon and find CLGEN generates a lot of comments, repeated dead statements

and awkward non-human-like code such as multiple semi-colons. These results agree with the case study by Goens et al. [59] that shows the AST depth distribution of CLGEN’s code is significantly narrower compared to code from GITHUB or standard benchmarks.

Feature space coverage. To further enhance our comparison, we perform an analysis on the feature space coverage of BENCHPRESS’s and CLGEN’s synthesized programs in all three feature spaces. Feature coverage is the most critical metric when evaluating the effectiveness of a benchmark synthesizer for predictive modeling. We use Principal Component Analysis (PCA-2) to represent the feature spaces in an easy to visualize 2-dimensional space. In Figures 5.6a, 5.6b and 5.6c we show the extent of feature space covered by candidates in the two approaches. CLGEN’s samples are clustered around the origin, while there is one outlier for Autophase and two for Grewe’s et al. and InstCount features. Candidates generated by BENCHPRESS are more scattered achieving a much wider coverage of the feature space.

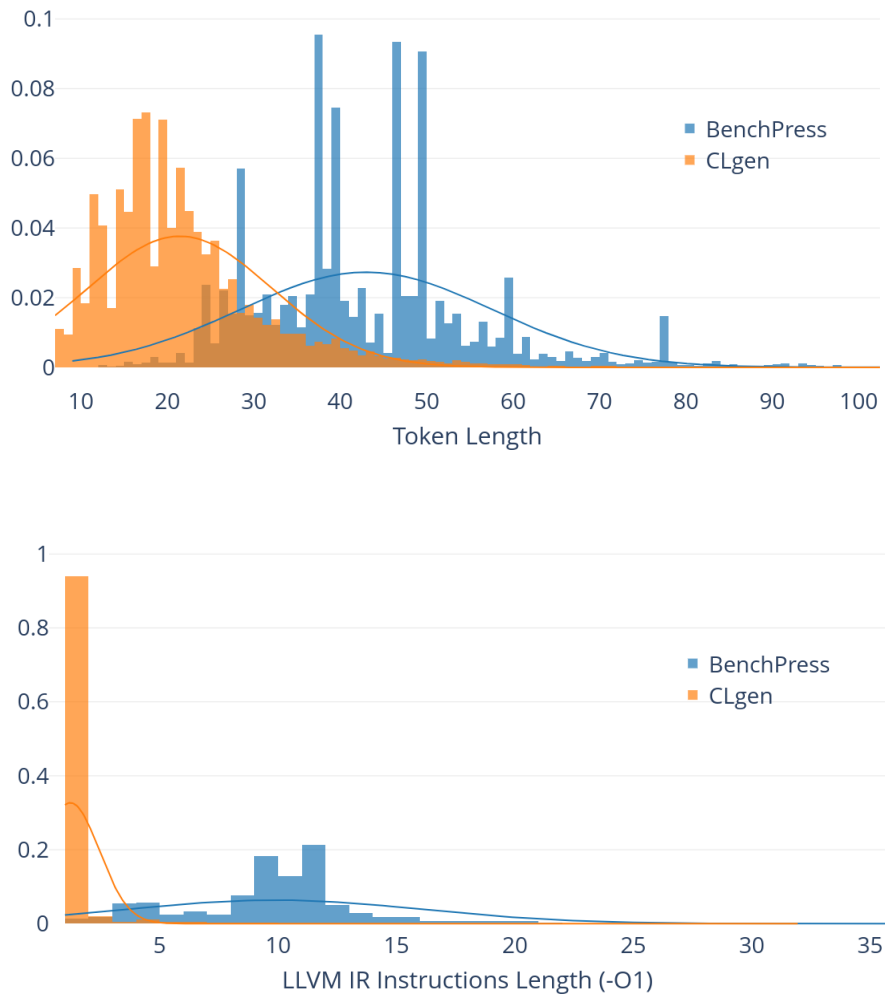


Figure 5.5: Probability distribution of (a) token length and (b) LLVM-IR Instruction count among BENCHPRESS's and CLGEN's generated benchmarks. BENCHPRESS's benchmarks presented here are generated at a single inference step without iteratively directing program synthesis.

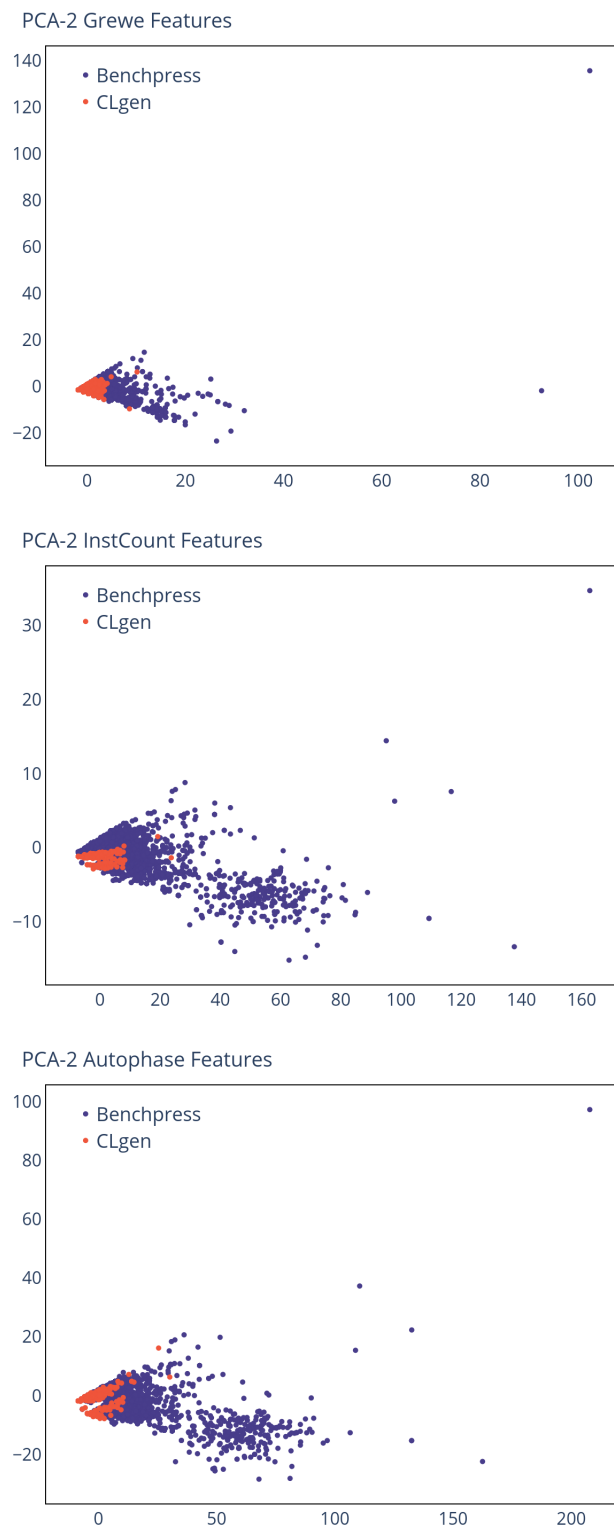


Figure 5.6: PCA-2 representation of feature space coverage of BENCHPRESS and CLGEN for (a) Grewe's et al., (b) InstCount and (c) Autophase feature spaces. In this experiment, BENCHPRESS's generation is undirected and no iterative space search is performed.

5.5.2 Targeted Benchmark Generation

We use beam search to generate samples that target desired parts of the feature space. We compare BENCHPRESS with human-written benchmarks from GITHUB and synthetic benchmarks from CLGEN and CLSMITH in targeting the features of Rodinia benchmarks on three feature spaces. We use SRCIROR code mutator with beam search to collect GITHUB and CLSMITH benchmarks with closer features. For each target benchmark, we gather one OpenCL kernel per evaluated dataset whose features have the minimum available Euclidean distance from the target features. Figures 5.7, 5.8 and 5.9 show the relative proximity of each benchmark to the target. This proximity is the complement of the relative distance of the two kernels, i.e, 1 minus the distance between the two kernels in the feature space relative to the distance of the Rodinia kernel from the axes origin. This allows us to express the quality of the match with an intuitive 0% to 100% scale: 100% means the two kernels have the same features, 0% means the best kernel is as close to the target as an empty kernel. We mark perfect matches with a white asterisk (*).



Figure 5.7: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over Grewe’s et al. feature space. Relative proximity is 1 minus the distance of the two kernels in the feature space relative to the distance of the Rodinia benchmark from the axes origin. 100% means an exact match in features and is highlighted with a white asterisk (*). A score towards 0% indicates the closest match is closer to the axes origin than the benchmark, i.e., a very small or empty kernel.

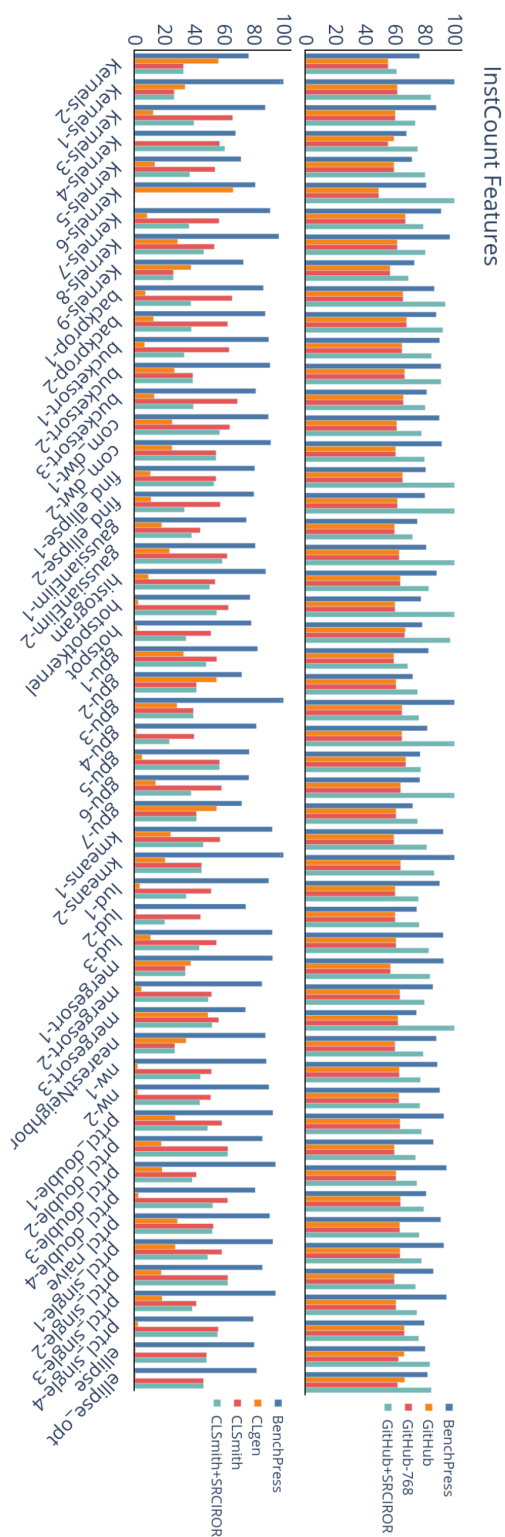


Figure 5.8: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over InstCount feature space.

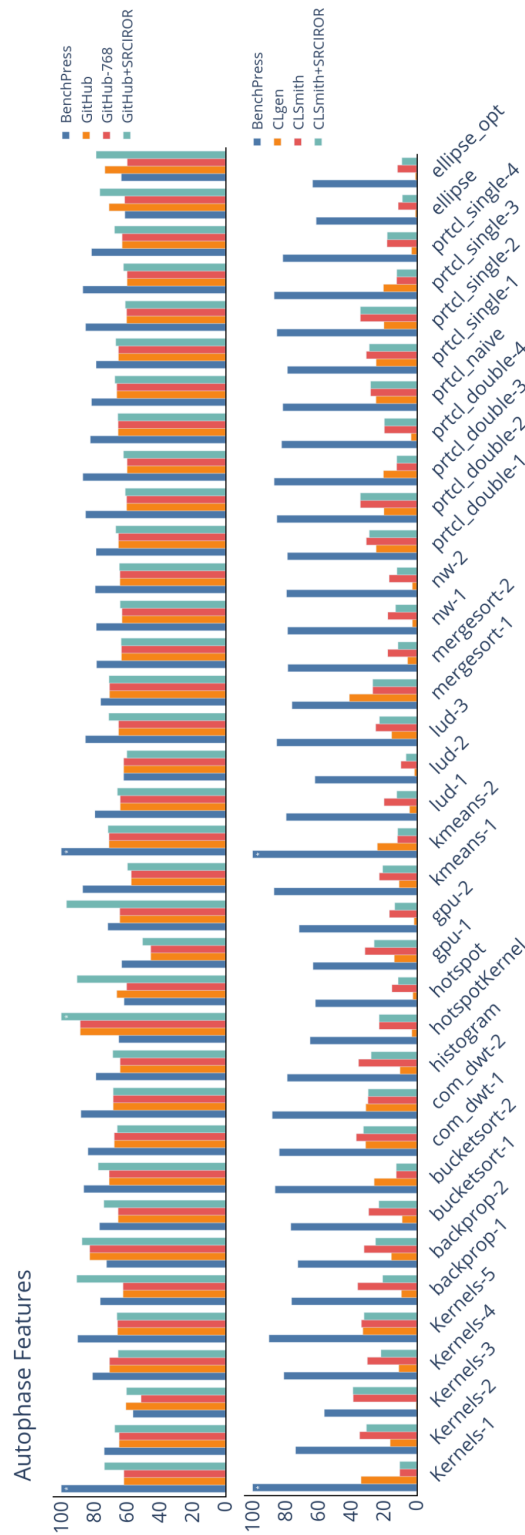


Figure 5.9: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features. We report the best match for seven datasets (BENCHPRESS’s, CLgen’s, GitHub’s and GitHub-768’s datasets also combined with exhaustive mutations with SRCIROR) over Autophase feature space.

Performance on syntactic features. On Grewe’s et al. feature space, BENCHPRESS generates kernels that are the closest ones in features for all 22 Rodinia Benchmarks compared to CLGEN and CLSMITH, and 20 out of 22 compared to GITHUB and GITHUB-768. BENCHPRESS synthesizes an exact match (100% relative proximity) for 14 target benchmarks.

We pick out and discuss a few examples from our results. The absolute distance achieved for ‘nw-1’ and ‘ellipse_opt’, is 1.0. For both targets, almost all features match except for one missing instruction (COALESCED MEM ACCESS and ATOMIC INST respectively). For ‘hotspot’ GITHUB and BENCHPRESS produce a candidate kernel with exact matching features. However, BENCHPRESS generates the matching candidate kernel in 421 tokens, unlike GITHUB’s closest benchmark that has 798 tokens. For the two target benchmarks that BENCHPRESS’s candidates were not closest to, we found only GITHUB contains better samples for `com_dwt-3` and `gpu-1`, while BENCHPRESS does not. We find both benchmarks to be fairly large (901 and 5,200 tokens respectively) and BENCHPRESS cannot reach these features within 768 tokens. For the same reason, GITHUB-768, CLGEN and CLSMITH does worse than BENCHPRESS on these targets.

Performance on LLVM IR features. Autophase and InstCount features are extracted from the LLVM-IR of a program that has been compiled with `-O1` flag to apply basic optimisations such as dead code elimination. BENCHPRESS occasionally generates repeating operations that a compiler will remove or numerical operations that may be reduced to simple assignments. Owing to these optimisations, we find targeting benchmarks on these two feature spaces is more challenging than Grewe’s et al. syntax-level features. With InstCount features, BENCHPRESS generates candidates whose features completely match 2 out of the 52 Rodinia benchmarks. Among the remaining 50, BENCHPRESS outperforms CLGEN, CLSMITH, GITHUB and GITHUB-768 for all target benchmarks, achieving higher proximity. SRCIROR significantly improves GITHUB leading to GITHUB+SRCIROR to achieve better proximity for 18 out of 52 Rodinia benchmarks compared to BENCHPRESS. On Autophase features, BENCHPRESS generates candidates matching the same 2 target benchmarks, while outperforming CLGEN, CLSMITH and GITHUB on 30 out of 36 Rodinia benchmarks in total. GITHUB+SRCIROR performs better than BENCHPRESS for 8 out of 36 target benchmarks and produces an exact match for ‘hotspotKernel’.

We previously explain the importance of having diverse features in compiler benchmarks and we show, in Figure 5.1, how sparse Rodinia benchmarks are on Grewe’s et

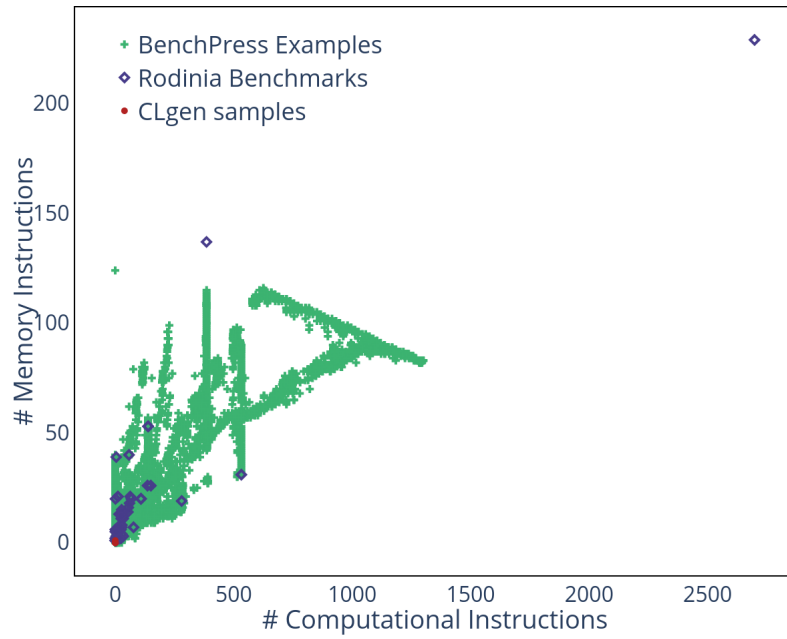


Figure 5.10: # Memory operations and # computational instructions for **(a)** Rodinia benchmarks in purple diamonds, **(b)** CLGEN’s samples in red dots and BENCHPRESS’s benchmarks in green crosses after performing directed search for all Rodinia benchmarks.

al. reduced feature space and how CLGEN fails to provide any additional features. Now we introduce into this 2-dimensional space all BENCHPRESS’s kernels that are generated while performing directed space search to target Rodinia benchmarks and we present them in Figure 5.10. BENCHPRESS densely populates the space around the target benchmarks that are clustered around the lower left corner. We find BENCHPRESS’s samples progressively converge to the target benchmark features with successive generations. For example, BENCHPRESS targets `com_dwt-3` at 385 computational and 137 memory instructions, starting from the axes origin and attempting to reach its features from different directions. One of the directions prevail but does not manage to exactly reach the target. The same happens for the top right point, `gpu-1`. BENCHPRESS’s samples get closer developing a straight line from the origin to 1,000 computational and 100 memory instructions. At this point BENCHPRESS is restricted by its sequence length and cannot augment further its samples. This is depicted by its attempt to reduce the distance by swapping the two instruction types within the same token length, forming a perpendicular line with a negative slope. We argue the area of Grewe’s et al. feature space that BENCHPRESS can cover within 768 tokens to be the area of the triangle formed by the intersections of the axes with the extension of the

negative slope line developed by BENCHPRESS’s samples.

Summary - BENCHPRESS vs GitHub vs CLgen vs CLSmith. 6 of the targeted Rodinia benchmarks exceed BENCHPRESS’s maximum sequence length of 768 tokens. In LLVM-IR feature spaces, care must be taken to generate code that will not be removed by compiler optimisations. This is a difficult challenge for source code generative models. However, our results demonstrate that BENCHPRESS can generate OpenCL kernels that approach target human-written benchmarks compared to GITHUB code and CLGEN candidates. Our experiments also show BENCHPRESS is dramatically better in all cases than CLGEN, the current state of the art in OpenCL synthetic benchmark generation. We further elaborate on BENCHPRESS’s performance in the next Subsections.

5.5.3 Active Learning for Feature Selection

We combine BENCHPRESS’s ability to generate benchmarks targeting desired features with active learning in order to generate benchmarks that improve the training of the Grewe et al. heuristic. We evaluate this against passive training with CLGEN, GITHUB code, and BENCHPRESS with randomly selected target features. All approaches augment the same baseline training set that is taken from [40], containing 7 benchmark suites¹. Table 5.2 shows the effect of each approach on the predictive power of the heuristic. Training only on human written benchmarks improves the heuristic’s performance by 4%, as shown in Table 5.2’s first row. To understand the maximum achievable improvement in the heuristic, we compute the best speedup (= 12%) that is achieved if the model chooses the optimal device as opposed to always picking the GPU. For 71% of the benchmarks, GPU is the optimal device, so no speedup improvement is possible. For the remaining 29% benchmarks, predicting the ‘CPU’ label correctly with Grewe et al. will result in a speedup improvement.

BENCHPRESS using active learning (BENCHPRESS-AL) clearly outperforms all other approaches in terms of average speedup, improving it by 6%. When trained on BENCHPRESS with passive/random feature selection (BENCHPRESS-P), the speedup achieved is only 1%. To our surprise, the same speedup is achieved with GITHUB, which is worse compared with training only on the original benchmark suites. We further analyze the dataset collected from GITHUB code and we find it to be imbalanced with 90% of its training instances are labelled as ‘GPU’. This leads the model

¹The benchmarks have been updated with a wider range of `global` and `local` sizes.

	Speedup %	Precision	Recall	Specificity
BENCHMARKS	+4%	0.81	0.86	0.61
BENCHPRESS-AL	+6%	0.84	0.86	0.64
BENCHPRESS-P	+1%	0.84	0.85	0.48
CLGEN	-1%	0.52	0.86	0.43
GITHUB	+1%	0.85	0.83	0.61

Table 5.2: Grewe et al. heuristic model’s performance, precision, recall, and specificity when trained on each technique. Speedup is the geometrical mean of speedups over all benchmarks relative to the optimal static decision, i.e. running on the GPU. Precision, recall, and specificity treat GPU labels as positive and CPU labels as negative.

having a higher precision of 0.85, i.e. predicting correctly that a kernel should execute on the GPU, but falling short when it comes to correctly predicting the ‘CPU’ label. Training the heuristic with CLGEN actually leads to a slowdown: it is 1% slower to execute kernels on the predicted devices compared to statically executing everything on the GPU, the baseline device. We analyze CLGEN’s dataset and observe the opposite pattern found in GITHUB’s dataset. 63% of its training data execute faster on the CPU than on the GPU. This is a direct consequence of CLGEN generating small benchmarks that are poor in features, as the CPU may be slower than the GPU but the large overhead of transferring data to the GPU makes the CPU a better choice for small workloads. CLGEN containing too many CPU-labeled kernel explains the heuristic’s low precision and specificity, as it becomes biased to select the CPU very often leading to a slowdown.

Our main motivation behind using active learning is that it gives BENCHPRESS the ability to target directly those parts of the feature space that will maximize a downstream task’s performance. To assess the active learner’s performance, we compare the Grewe et al. heuristic’s speedup when trained on BENCHPRESS’s benchmarks that target areas of the feature space selected by the active learner versus benchmarks that target random features. In both cases, we execute BENCHPRESS for the same amount of time, 10 sampling epochs (i.e., performing steered generation for 10 target feature vectors). In Figure 5.11, we show the speedup achieved by the heuristic when trained on the data collected at that step. Using active learning to target features, BENCHPRESS’s dataset improves the heuristic’s speedup by 50% after 5 sampling steps, from 4% to 6%. Targeting random features never leads to a speedup higher than

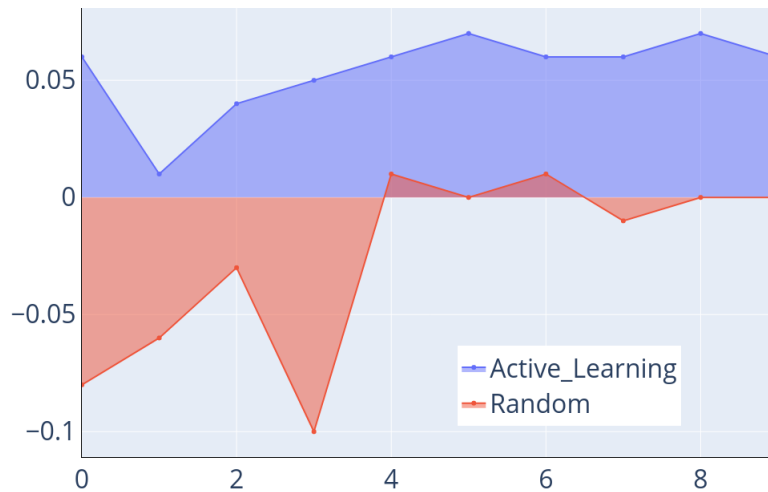


Figure 5.11: BENCHPRESS’s performance enhancement of Grewe et al. heuristic model when using active learning compared to passively targeting random parts of the feature space over the course of 10 sampling epochs. The y-axis shows the performance enhancement as a percentage for each sampling epoch (0 to 9) shown in x-axis.

1%. BENCHPRESS can still develop the same speedup by targeting random features if infinite amount of time was available. Our active learner ensures that missing features are going to be quickly targeted, improving the state of the art within 5 sampling epochs.

5.6 Summary

Predictive models for compilers have been shown to outperform compiler experts but they are restricted by the amount and quality of training data they are exposed to. What is needed is an approach that can synthesize benchmarks and enhance datasets with missing features. In this Chapter we propose BENCHPRESS, a powerful code generator that uses active learning to search the feature space and steers generation towards desired features. BENCHPRESS generates $10\times$ more and $7.5\times$ larger undirected benchmarks with $37\times$ greater compilation rate than CLGEN - a state of the art compiler benchmark generator - from a fixed input feed. BENCHPRESS outperforms CLGEN, CLSMITH, code from GITHUB and applied mutations with SRCIROR in generating OpenCL kernels that target the features of Rodinia benchmarks developed by human experts. BENCHPRESS applies active learning to enhance Grewe’s et al. dataset with benchmarks with missing features and leads to improving the heuristic’s speedup by 50%. We hope this work to demonstrate a sustainable method to direct feature space search of program generation and that BENCHPRESS’s release to researchers will enable research in related domains.

Chapter 6

Deep Directed Language Modeling for Compiler Features

6.1 Introduction

This Chapter presents the design and implementation of a novel, directed language model for compiler benchmarks. The proposed contribution is based on deep, bi-directional, unsupervised learning. In the following Sections, the approach, implementation and evaluation of this technique are discussed.

We develop `BENCHDIRECT` [143], a BERT-based directed language model for OpenCL benchmark generation [46, 139]. Similarly to our previous contribution, namely `BENCHPRESS`, `BENCHDIRECT` targets and synthesizes benchmarks from desired parts of the feature space. The key distinction in `BENCHDIRECT` is a feature encoder integrated into its language model which enables conditioning the token generation on the desired features along with the context of the input source code. `BENCHDIRECT` is able to target areas of the feature space fast, often at a single inference step, contrary to `BENCHPRESS`'s feature agnostic-language model that performs thousands of random inferences.

We compare `BENCHDIRECT` against `BENCHPRESS` in targeting the features of Rodinia benchmarks that are developed by compiler engineering experts. Its directed language model enables it to target features with significantly higher accuracy than `BENCHPRESS` at a fraction of the total inference time needed by the latter. Across three feature spaces, `BENCHDIRECT` matches perfectly the features of Rodinia benchmarks $1.8\times$ more frequently compared to `BENCHPRESS`. It also generates benchmarks that are up to 72% closer to Rodinia target features compared to `BENCHPRESS`'s sam-

ples, while its sampling process is up to 36% faster. Developing a directed program synthesizer with a reduced time complexity paves the way for large scale applications on more difficult compiler optimisation tasks.

In this Chapter, we present the following contributions:

1. We develop the first directed, generative language model for compilers that conditions program generation towards feature space representations.
2. We provide an extensive empirical evaluation on the trade-off between execution time and accuracy in targeting compiler features using `BENCHDIRECT` versus `BENCHPRESS`.

6.2 Approach

`BENCHPRESS`'s synthesizer presented in Chapter 5 is feature agnostic. It infills source code given the input context left and right of the `[HOLE]`. `BENCHPRESS` is only able to steer program generation through a costly beam search on the model's output: we generate a large number of random code candidates and we feed those that are closer to the target features back into the model's input with new holes for further edits. Given `BENCHPRESS`'s language model is undirected, it often needs hundreds of thousands of code candidates to increase the chance of finding a few with the right features. This is inefficient and unsustainable on complex compiler tasks.

Instead of randomly trying to fill the space with new benchmarks to get closer to the target features, a more desirable approach to target them directly during synthesis is needed. Ideally, this would help generate a benchmark with the right features in a single inference. To this end, we develop `BENCHDIRECT`.

`BENCHDIRECT` consists of the first directed language model for compiler benchmark synthesis. Along with the masked source code input, `BENCHDIRECT` also encodes its compiler features before masking. Its classification head selects tokens to fill a `[HOLE]` by jointly observing the code context and the encoded features. This leads to selecting tokens that are likely to generate a kernel that is (a) compiling (similarly to `BENCHPRESS`) but also (b) matching the target features provided in the input.

`BENCHDIRECT` inherits `BENCHPRESS`'s three stage pipeline: (a) We collect data and preprocess them, (b) we train the language model and (c) we sample the model to steer benchmark generation. The same techniques are used for this contribution's

step (a) and (c) with BENCHPRESS, discussed in Chapter 5. BENCHDIRECT is distinguished by its underlying language model that conditions token generation on the targeted features. In the rest of this Section, we discuss BENCHDIRECT’s directed language model component. For further details on the data processing, the masked language model training and the sampling technique for steerable generation, we redirect the reader to Section 5.3.

6.3 Directed Language Modeling

BENCHDIRECT’s feature encoder is based on Transformer [149] and is shown in Figure 6.1. We encode a vector of numerical compiler features using an embedded layer with positional encoding followed by a Transformer-Encoder. We reduce the dimensions of the Transformer’s output using a Fully Connected layer to match BERT language model’s hidden state representation of its input source code. Both hidden states are concatenated and fed to a Fully Connected layer with GELU [68] activation to extract correlated features. Finally, a Decoding Fully Connected layer projects the joint hidden state into the vocabulary space. The feature encoder’s input consists of 134 positions divided into three fixed segments. Each represents one feature space used in our evaluation: (a) 8 positions for Grewe’s et al. features, (b) 56 for Autophase and (c) 70 for InstCount features. BENCHDIRECT can support multiple spaces and it only needs to be trained once to direct benchmark synthesis on any of them. To steer generation in a new feature space, we simply need to extend a new segment in the Transformer-Encoder’s input and apply fine-tuning using the new space’s feature extractor to collect data from our training corpus.

BENCHDIRECT is trained with the same approach described in Subsection 5.3.2. We sample randomly one OpenCL kernel and introduce a [HOLE] to provide it to the language model’s input. The model learns to predict the first token of the hidden sequence using cross categorical entropy loss function. Introducing compiler features in training is the distinction to this process. When one OpenCL kernel is sampled, its compiler features are also collected. The model receives a pair of inputs, (src_i, fv) and one output $token_i$, where i is the index at which the [HOLE] is located.

It is important to note that we do not feed the feature vectors of all three feature spaces to the encoder at the same time. Instead, we uniformly select one, we set its values to the respective segment of the encoder’s input and we [PAD] all other positions such that gradients are not applied. Over training time, the model observes datapoints

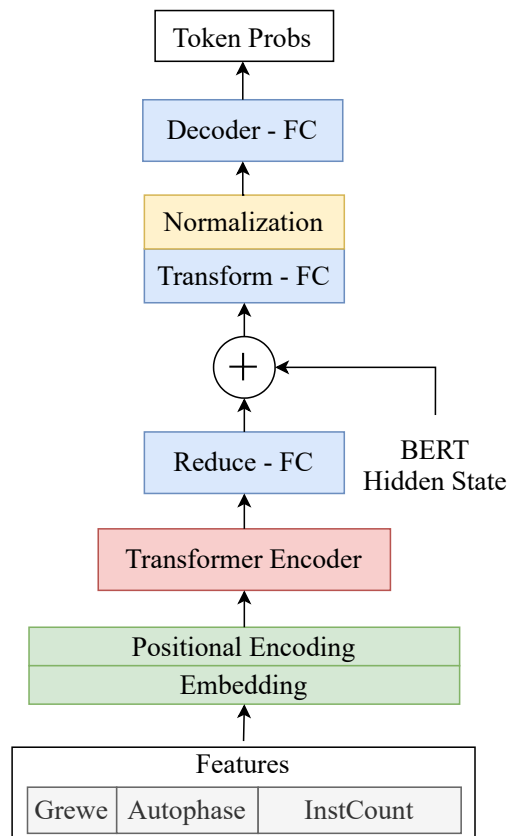


Figure 6.1: BENCHDIRECT's directed language model design.

from all feature spaces for every kernel. Padding all feature spaces but one allows the trained model to learn how to direct synthesis to each one of them independently. Providing vectors from all spaces as one datapoint would possibly allow the model to learn correlations between them but this is not useful to us. What is more, directed synthesis on one of the feature spaces would be impossible. The model would have been trained to observe all three feature vectors for one given source code input. This means we would have to know the mapping function among all feature spaces to translate a target feature vector to all supported ones for the encoder's input. Instead, keeping one feature space per datapoint leads to the encoder's weights to be tuned accordingly to perform accurately on all spaces separately. Parts of the network (e.g. the FC layers) are jointly trained to optimise all feature spaces encoding. Other parts, such as the (Q, K, V) matrices are grouped in vectors, one for each index separately, and are only trained when their respective positions are not padded. An alternative solution would be to use many Transformer-Encoders, one per feature space, and train each separately. During generation, the appropriate Transformer would be manually selected given the desired feature space. Although this is a valid approach, there is no evidence to suggest

it would perform better than one Transformer model large enough to learn all segments separately.

During sampling, `BENCHDIRECT` receives a source code input and the target features as an input. Given the code context and the `[HOLE]` position, the model will attempt to select those tokens that will produce a compiling kernel with features as close as possible to the target in that respective feature space. At its best, we hope `BENCHDIRECT` can receive an empty code input and provide the target benchmark at a single inference step. At the very least, the beam search sampler will go through fewer iterations and fewer inferences per generation compared to `BENCHPRESS`.

6.4 Experimental Setup

We describe the configurations used in training `BENCHDIRECT` and the parameters used in our evaluation, namely (1) Feature Spaces - we use three different representations of program features, (2) Target Benchmarks - We use Rodinia benchmarks [33] and their features as the target for synthesis by `BENCHDIRECT` and (3) Comparison to state of the art - we compare `BENCHDIRECT` against `BENCHPRESS` in targeting the features of Rodinia benchmarks, measuring accuracy and total execution time.

6.4.1 Platforms

We train both models and conduct all our experiments on two 64-bit systems each having one Intel Xeon E5-2620 32-core CPU, 4x Nvidia TITAN X Pascal GPU and 64 Gigabytes of RAM. We use Ubuntu 18.04, PyTorch 1.9.1 [123], CUDA version 11.4 and Nvidia driver version 510.47.03. We use Clang-10 to compile programs and LLVM-10 to compile and execute InstCount and Autophase [65] extracting tools. For compatibility reasons, we are required to use Clang LibTooling from LLVM-6 to execute Grewe’s et al. [61] feature extractor.

6.4.2 Language Modeling for source code

We use the same OpenCL dataset collected from GITHUB that was used to train `BENCHPRESS` in Chapter 5. This dataset consists of 19,637 compiling OpenCL kernels with re-written variables and formatted text. In this evaluation, we train both `BENCHPRESS` and `BENCHDIRECT` on this dataset for 40 epochs of 200,000 steps each using a batch size of 64. For both models’ BERT hyper-parameters, we select 2 hidden layers

and 12 attention heads. We set intermediate size and hidden size to 768 and max position embeddings to 512. We set the maximum length of holes to 100% of a kernel’s token length. We use Adam optimizer to train the network with a learning schedule of 20,000 warmup steps that peaks at 45×10^{-6} . Both models are trained to a final loss value of 0.14. In this evaluation, we train both language models for more epochs compared to the training process described in Section 5.4.2. After experimentation, we observe the quality significantly improves after the loss is reduced below 0.2.

6.4.3 Feature Spaces

Compiler predictive models use static code features to represent programs and learn optimisation heuristics. A vector of independent characteristics represents a single program. Each of them are typically an integer or float value. Features are extracted at the Syntax level by traversing the AST or at the IR level using the compiler’s middle end (e.g. LLVM-IR). A feature space is the collection of all possible program feature vectors.

BENCHDIRECT is a generative model that can be steered to generate samples for a desired part of the feature space. Its language model conditions on compiler features during program generation, enabling the model to steer into the feature space efficiently. We evaluate BENCHDIRECT on three source feature representations we find across the literature, (a) Syntax-level Grewe’s et al. features [61], (b) IR-level LLVM-InstCount [98] and (c) IR-level Autophase [65]. We describe these three feature spaces in detail in our second contribution’s evaluation at Subsection 5.4.3.

6.4.4 Targeted Benchmark Generation

BENCHPRESS, our second contribution, develops strong performance compared to state of the art program synthesizers and its benchmarks outperform even human-written benchmarks from GITHUB in two tasks, (a) targeting the features of Rodinia benchmarks and (b) improving the accuracy of a compiler heuristic model. However, its undirected language model requires up to hundreds of thousands of inferences for its beam search sampler to minimize its samples’ distance from the target features. BENCHDIRECT strives to address this inefficient process, using a directed language model.

We repeat the experiment of Section 5.4.5 to evaluate BENCHDIRECT’s accuracy and execution time in targeting the features of Rodinia benchmarks compared

to BENCHPRESS. We target the features of Rodinia benchmarks in all three feature spaces for a range of different workload sizes: 32, 64, 128, 256, 512, 1024 and 2048. A large workload size leads to a significant time overhead but is required to ensure high accuracy for BENCHPRESS’s undirected language model. This may not be the case for BENCHDIRECT’s directed synthesizer, speeding up directed generation without compensating on its accuracy. In this experiment, we explore how this parameter affects accuracy and total execution time for both models. We set an upper threshold of 6 beam search iterations per target feature vector. We select 6 iterations as they are enough to observe the difference between both models in reducing their distance from the target without introducing a large time overhead to conduct our experiments.

6.5 Results And Analysis

In this Section, we show our experiments’ results in comparing BENCHDIRECT with BENCHPRESS. We present a thorough study of both models’ accuracy and execution time and discuss their trade-offs in targeting the features of human-written, compiler benchmarks.

6.5.1 Targeted Benchmark Generation

We target the features of Rodinia benchmarks using BENCHPRESS and BENCHDIRECT. Both models use beam search over their synthesizer to minimize their samples’ distance from the target features. At the end of each search, we select the generated kernel whose features have the minimum Euclidean distance from the target benchmark. We perform this experiment for multiple beam search candidate sizes: 32, 64, 128, 256, 512, 1024 and 2048. In Figures 6.2a, 6.2b and 6.2c we show the Pareto fronts of the average relative proximity achieved over all Rodinia benchmarks versus the total amount of inferences. Relative proximity is defined in Section 5.5.2 as a percentage of how close a feature vector is to the target features relatively to the axis origins. Inferences are calculated as the number of beam search iterations to target all benchmarks multiplied by the workload size. Each datapoint is annotated with its workload size configuration. In Figures 6.3a, 6.3b and 6.3c, we show BENCHDIRECT’s improvement in accuracy and execution time compared to BENCHPRESS, for each workload size setting.

BENCHDIRECT outperforms BENCHPRESS in average relative proximity and total

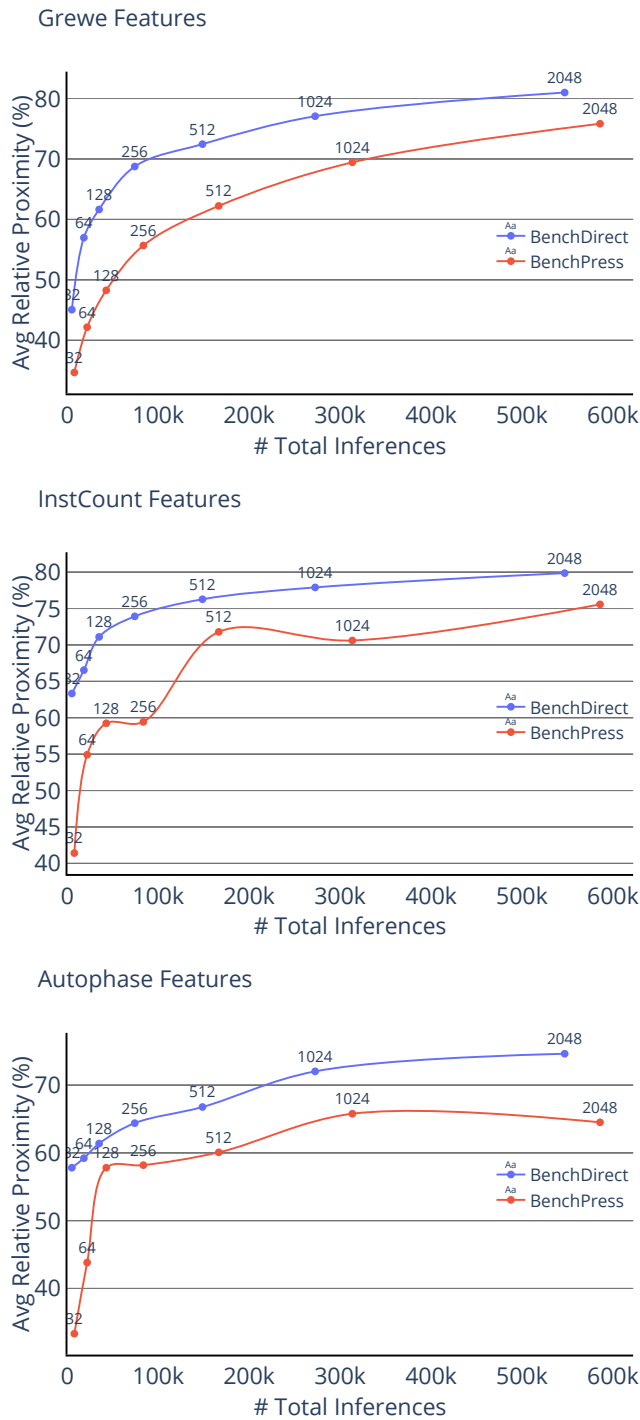


Figure 6.2: Pareto fronts of the average relative proximity versus total inferences for BENCHDIRECT and BENCHPRESS in targeting Rodinia benchmarks over three feature spaces ((a) Grewe’s et al., (b) InstCount and (c) Autophase). Higher relative proximity and fewer inferences are better, therefore optimal points, i.e., Pareto-dominant, are those towards the top left. We annotate the workload size configuration per Pareto point.

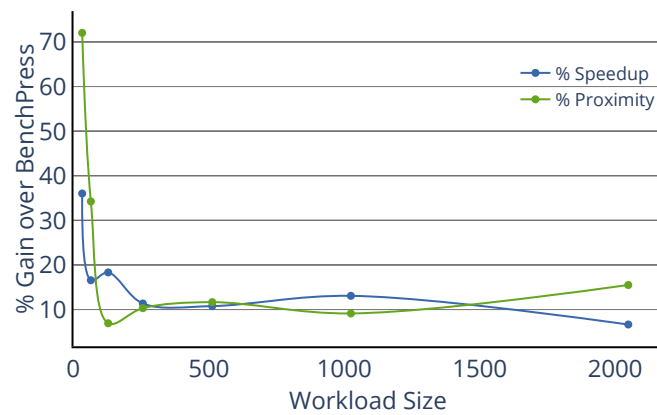
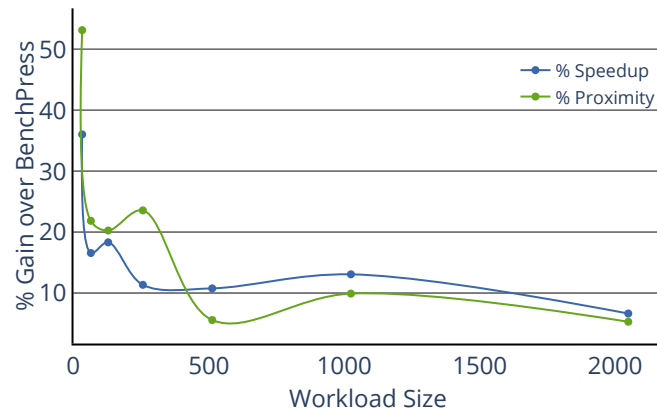
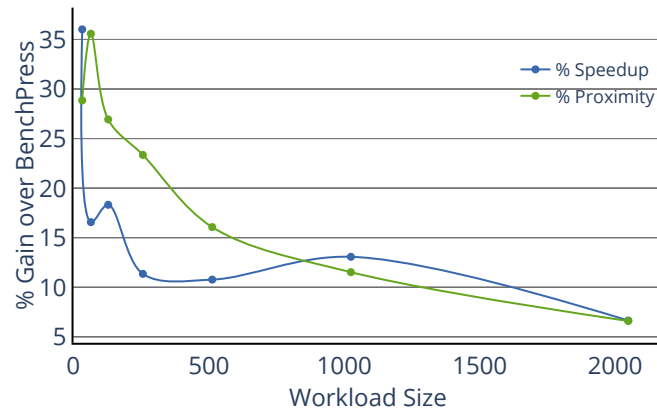


Figure 6.3: BENCHDIRECT's acquired execution time speedup and relative proximity improvement over BENCHPRESS per workload size configuration for (a) Grewe's et al., (b) InstCount and (c) Autophase feature spaces.

inferences for all workload size configurations, across all three feature spaces. Taking the average proximity and the execution time as a design space, the datapoints that are optimal with respect to these two metrics belong exclusively to BENCHDIRECT, while there are no configurations for BENCHPRESS that optimise either metric compared to the former. The effect BENCHDIRECT’s directed language model has in targeting features is especially highlighted when the workload size is small. BENCHDIRECT’s synthesizer conditions directly on the target features and provides, in very few attempts, candidates that match or are very close to them. This means a dramatic reduction in the amount of benchmarks per beam search does not drastically hamper the model’s accuracy. The same is not true for BENCHPRESS. While BENCHDIRECT offers an average speedup of 10.2% and an improvement in average relative proximity of 10.1% for workloads greater or equal to 512, the speedup reaches up to 36% in all three feature spaces and the accuracy gain up to 72% on InstCount features for smaller workloads. This indicates BENCHDIRECT remains consistent in the amount of iterations needed to achieve high accuracy, while BENCHPRESS suffers in both areas.

Both models achieve a peak accuracy when they use a workload size of 2048. This is expected as generating more candidates increases the probability of getting closer to the target features. Using this configuration on both models, we show in Figures 6.4, 6.5 and 6.6 the best relative proximity achieved for each target benchmark in all three feature spaces. Similarly to Figures 5.7, 5.8 and 5.9, candidates whose euclidean distance from the target is 0 (i.e., perfect match feature-wise) are marked with a white asterisk (*). For a selection of Rodinia target benchmarks, we show how the minimum distance from the target is reduced over the course of 5 beam search iterations for both models in Figures 6.7, 6.8 and 6.9.

BENCHDIRECT generates $1.8\times$ more candidates that match exactly the target features compared to BENCHPRESS. Specifically, it matches 21 targets on Grewe’s et al. features, 14 on InstCount and 10 targets on Autophase, compared to BENCHPRESS’s 17, 3 and 5 exact matches respectively. Overall, BENCHDIRECT gets closer to the target compared to BENCHPRESS. Its samples are closer, or as close, for 45 out of 58 Rodinia targets on Grewe’s et al. features, 47 out of 52 on InstCount and 49 out of 52 on Autophase. BENCHPRESS provides better candidates for 13, 5 and 3 targets on Grewe’s et al., InstCount and Autophase features respectively. Even though it is expected for BENCHDIRECT to miss some target features due to the experiment’s randomness, we pick out a few such examples to discuss why this happens.

The largest performance gap in favour of BENCHPRESS is observed on ellipse

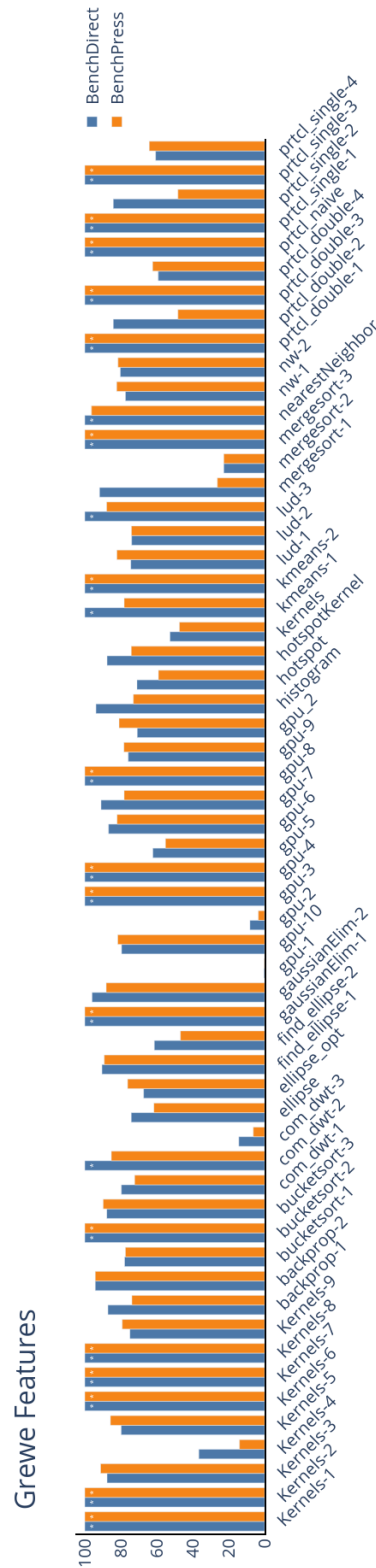


Figure 6.4: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on Grewe's et al. feature space. We show the best match for BENCHDIRECT and BENCHPRESS.

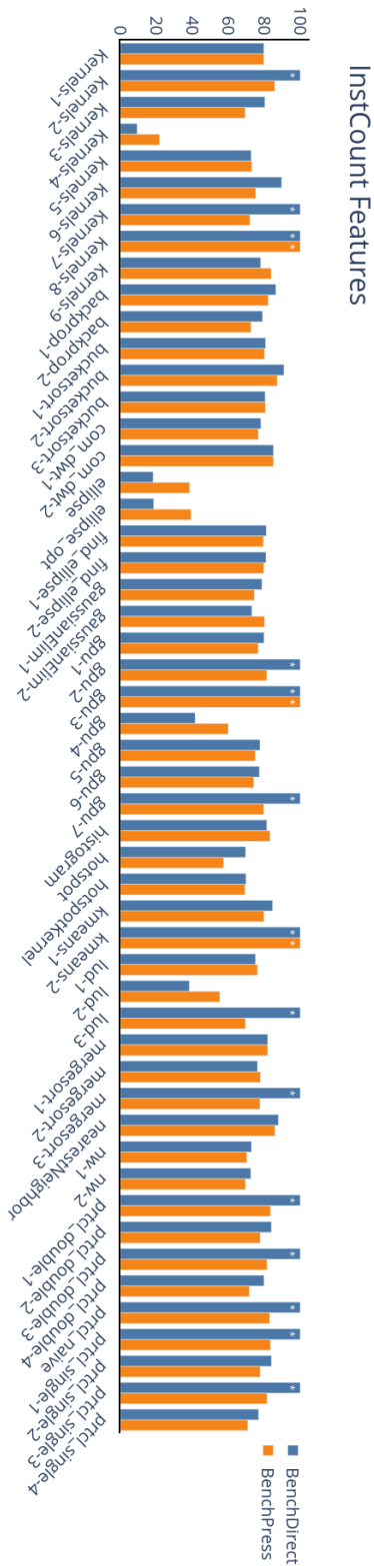


Figure 6.5: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on InstCount feature space. We show the best match for BENCHDIRECT and BENCHPRESS.

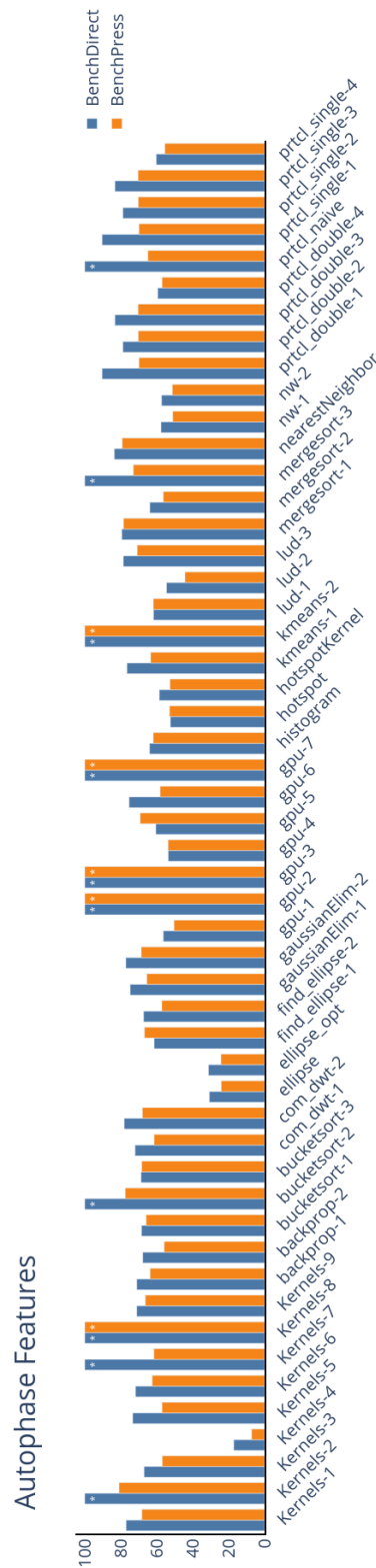


Figure 6.6: Relative proximity to each Rodinia benchmark of the candidate kernel with the closest features on Autophase feature space. We show the best match for BENCHDIRECT and BENCHPRESS.

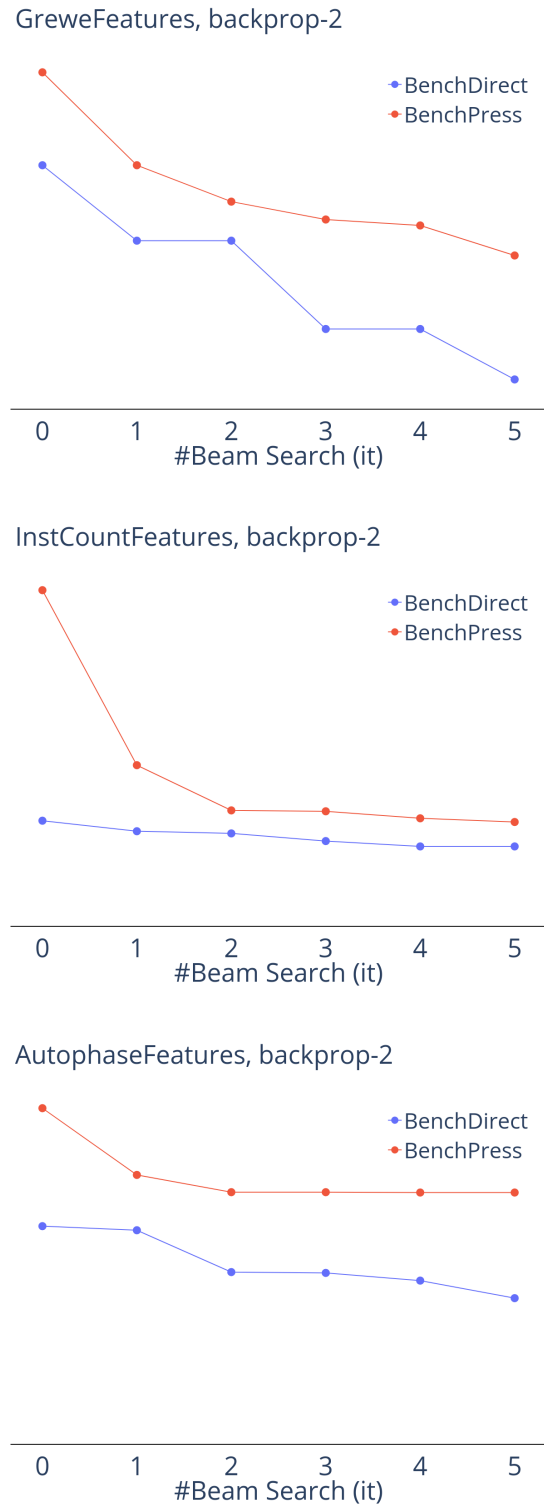


Figure 6.7: A comparative visualization of the minimum distance achieved (y-axis) from backprop-2 target benchmark over the course of six beam search iterations (x-axis) for all three feature spaces.

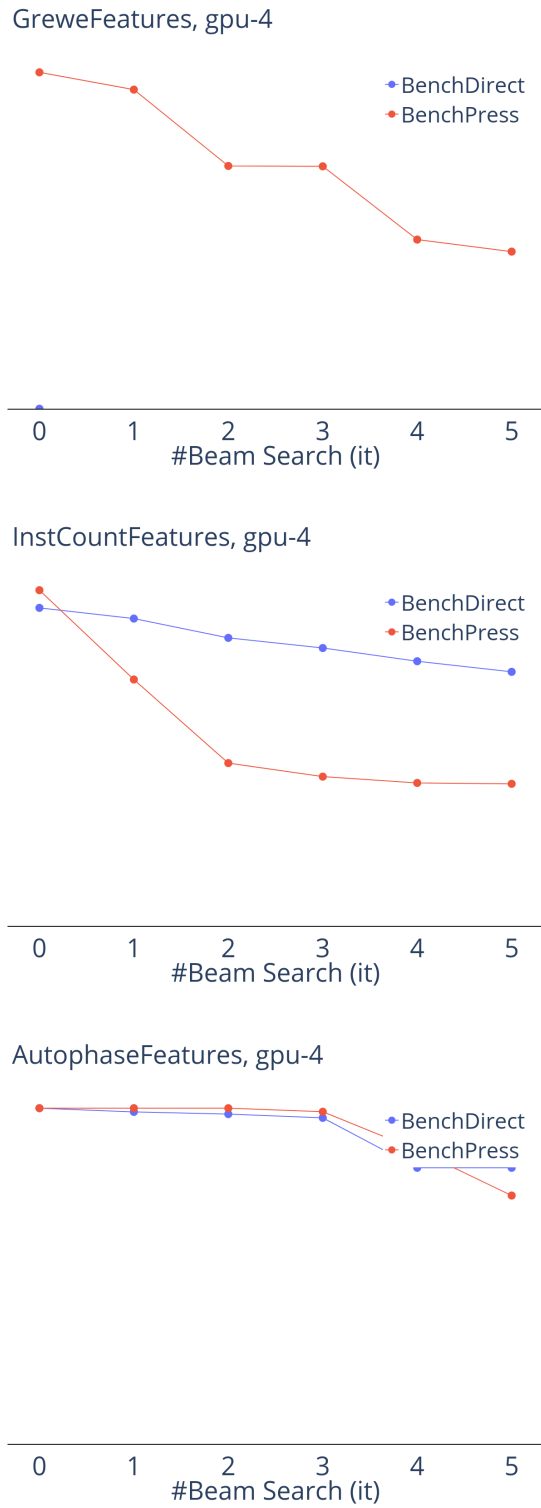


Figure 6.8: A comparative visualization of the minimum distance achieved (y-axis) from `gpu-4` target benchmark over the course of six beam search iterations (x-axis) for all three feature spaces.

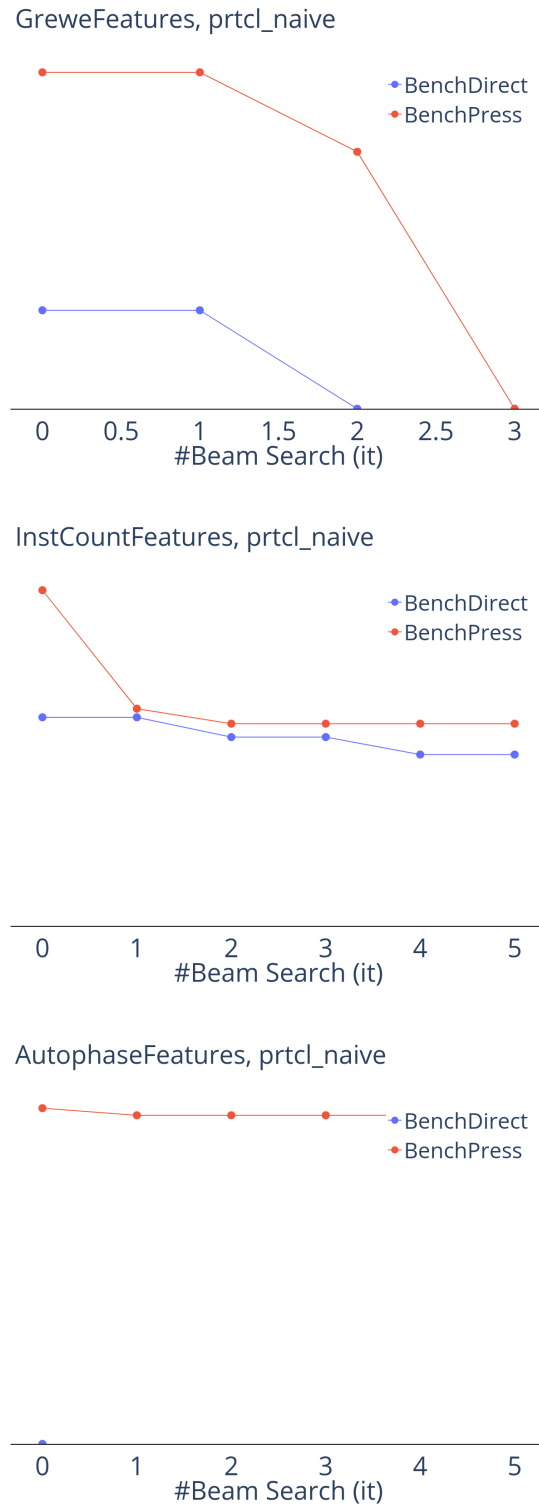


Figure 6.9: A comparative visualization of the minimum distance achieved (y-axis) from `particle_naive` target benchmark over the course of six beam search iterations (x-axis) for all three feature spaces.

and `ellipse_opt` on `InstCount` features. These two benchmarks are very large, containing multiple thousands of instructions, therefore they are difficult kernels to target. We examine both models' generated samples over all 6 beam search iterations. In both cases, we find `BENCHDIRECT`'s closest candidate on the first iteration to be 8% closer to the target compared to `BENCHPRESS`'s. After measuring the distance distribution from the target for both models' samples, we find `BENCHDIRECT` is 93% more likely to generate a sample whose distance is lower compared to `BENCHPRESS` on the first beam search iteration. `BENCHDIRECT` seems to succeed in these two target benchmarks indeed. However, at every inference step `BENCHDIRECT` tries to match the target features in a single [HOLE] infill. As these two kernels are very large, this is a challenging task leading to most of its produced candidates to have syntactic errors, leaving it with only a few benchmarks that compile. Even though its first iteration's samples are closer compared to `BENCHPRESS`, all successive iterations are becoming increasingly difficult for `BENCHDIRECT` to produce a compiling kernel which also reduces the minimum distance. For that reason, `BENCHPRESS`'s random and cautious steps lead to benchmarks that are eventually closer. We notice this pattern to happen in all targets where `BENCHPRESS` produced a better candidate. For these targets, it is likely that if we break down the difficulty into smaller steps by using intermediate feature vectors, this would have helped `BENCHDIRECT` to get to the target features gradually but more accurately. A potential downside of this approach is that identifying the granularity of such steps requires intuition and knowledge of the model. Breaking down a difficult target into too many, easy steps introduces an extra time overhead. On the other hand, if the intermediate vectors are too few and there is a large distance between them, `BENCHDIRECT` may not be able to travel from one step to the next.

6.6 Summary

Predictive models for compilers are crucial in constructing optimal heuristics for program optimisation efficiently. However, they are restricted by the quality of training data they are trained on. `BENCHPRESS`, our second contribution, proposes a steerable program generator that improves the Grewe's et al. CPU vs GPU heuristic model by 50% when it extends its training dataset with synthetic benchmarks. However, its undirected language model requires thousands of random inferences to reach close to the target, which is inefficient.

In this work, we present `BENCHDIRECT` which addresses this inefficiency.

BENCHDIRECT consists of a generative model based on BERT which has been extended to condition code synthesis on both input source code context and the features that are targeted. BENCHDIRECT's synthesizer is trained to perform directed synthesis on three different feature spaces that are commonly used in the domain of compiler optimisation: (a) Grewe's et al. features, (b) Autophase features and (c) InstCount features. BENCHDIRECT matches exactly the features of Rodinia benchmarks $1.8\times$ more frequently compared to BENCHPRESS, it is up to 72% more accurate and up to 36% faster in targeting their features. BENCHDIRECT improves significantly the process of generating compiler benchmarks that are directed into compiler feature spaces which is crucial to improve predictive modeling for compiler construction. Overall, it extends significantly BENCHPRESS's performance and we encourage developers to use this model for targeted benchmark generation. As BENCHDIRECT's synthesis depends on targeting features, BENCHPRESS is still the relevant model to use for general, undirected code generation. We hope this work to stimulate further research in the domain of directed program synthesis.

Chapter 7

Conclusion

This thesis proposes three novel methods to enable automation in software testing and compiler optimisation, two fields that still remain largely manual. Our first approach addresses the test oracle problem with deep neural networks used to summarize and classify program runtime correctness. Our second contribution tackles the compiler benchmarks shortage using deep unsupervised language models for program synthesis. Our approach improves predictive models for compiler optimisations whose performance was previously hindered by the training data shortage. Our third contribution proposes the first generative model for source code that is directed within compiler feature spaces. It improves significantly our second contribution’s accuracy and inference time in targeting the feature space. All three proposed techniques attempt to address challenges that had remained unsolved by existing research.

7.1 Contributions

This thesis makes three main contributions:

7.1.1 Automate the Test Oracle

The test oracle problem is a long-standing challenge in the field of software testing with experts spending time and effort to compute the expected behaviour of programs for thousands of randomly generated test cases. This thesis proposes a novel technique that uses deep learning to encode execution traces and classify them as “pass” or “fail”. The empirical evaluation on 15 open-source, industry-scale codebases illustrates our approach achieves maximum accuracy in classifying test executions of a codebase by only training on a fraction of them, specifically 14%.

7.1.2 Steerable Program Generation of Compiler Benchmarks

Predictive models for compilers have shown to outperform human experts in finding the right optimisation heuristics that enable programs to make use of all available resources of a target architecture. However, they are restricted by the limited amount and feature diversity of compiler benchmarks available to train on. This thesis proposes BENCHPRESS, a steerable program generator using active learning and deep unsupervised language modeling to tackle this compiler benchmarks shortage. BENCHPRESS finds important areas of compiler feature spaces that are likely to improve the performance of predictive models. Provided a set of desired features, BENCHPRESS synthesizes compiler benchmarks in OpenCL that target these features. Its benchmarks compile at a rate of 87% compared to the state of the art's 2.3%, while they can also be longer and more diverse. BENCHPRESS outperforms humans in targeting the features of high-quality benchmarks across three feature spaces, while also it improves the accuracy of a device mapping predictive model by 50%, by improving its training data.

7.1.3 Directed Language Modeling for Compiler Benchmark Generation

BENCHPRESS is the first steerable program generator for compiler benchmarks, but its inference time overhead can be quite high, as it needs to infer thousands of random benchmarks to get incrementally close to the target features. This can be prohibitive for large-scale experiments. In our third contribution, we propose BENCHDIRECT, a directed, bi-directional language model that infills programs by judging the left and right context of the code and also attending on the targeted compiler features. BENCHDIRECT increases BENCHPRESS's accuracy in targeting programs by up to 36%, while it reduces the total inference time by up to 72%. BENCHDIRECT matches exactly the features of Rodinia benchmarks in three feature spaces $1.8\times$ more frequently and it is 82% more likely to produce a benchmark closer to the target features compared to BENCHPRESS.

7.2 Critical Analysis

This Section contains a critical analysis of the techniques presented in this work.

7.2.1 Using Machine Learning as a Test Oracle

In Chapter 4, this thesis discusses the use of supervised learning as a test oracle for runtime behaviour classification. Although this contribution reduces the costly computation of expected outputs for test inputs by 86%, it has its limitations.

First, the machine learning model in this technique needs to be re-trained for each subject program. Although the training process does not require effort from the developer, it may be an extra overhead depending on the size of the model used. Also, the training data collection through instrumentation can pose several challenges. Our technique's interface to collect data is based on CMAKE and LLVM-IR. This means codebases that do not support these tools cannot be used through our model and execution traces will have to be collected manually. Maintainability is another issue. Such technologies evolve fast, therefore our instrumentation tools may need to be regularly maintained to keep up with modern compiler versions, otherwise it may be deprecated for new subject programs.

Second, we use recurrent neural networks, the LSTM specifically, to embed and classify program execution traces. Even though the LSTM is an effective architecture in encoding arbitrary sequences into fixed vector representations, there have been numerous new language models that significantly outperform it since publishing this contribution. Given the increasing complexity of software, it is unclear whether the LSTM can keep achieving maximum classification accuracy with execution traces becoming larger and more complex. Maintaining a high accuracy is imperative for our approach to be used by developers effectively.

7.2.2 Generative Modeling for Compiler Benchmarks

Chapter 5 presents BENCHPRESS, a technique for guided/steerable synthesis of compiler benchmarks. Using this technique improves the quality of ML-based heuristics for compilers by providing a fine-grained exploration of the space of representative programs. One limitation of this technique is that the synthesizer is restricted into generating functions. Although we keep a large collection of custom structs and data types, the model is restricted in invoking only them but cannot generate new ones,

limiting the space of possible program functionalities that can be produced.

Although BENCHPRESS is steerable thanks to beam search sampling applied over a workload of produced programs per generation, the underlying language model itself is unguided. This can lead to inefficiency as many thousands of inferred programs may be needed over tens of generations to target a single feature vector. For difficult parts of the feature space, this can amount to several minutes or even hours of GPU time to collect a single benchmark. However, there is no other existing work directing program synthesis in compiler feature spaces. What is more, we strive to mitigate BENCHPRESS's time overhead in our third contribution, BENCHDIRECT.

Another concern for BENCHPRESS's samples is how much they resemble code produced by humans. BENCHPRESS performs iterative edits on its own produced samples to guide their features towards the desired part of the feature space. Often, this leads to repetitive statements that will be eliminated by the compiler or misplaced blocks of code that are not executed on runtime e.g., computation inside an `if` statement that is always evaluated to `false`. Synthetic benchmarks that are not human-likely are difficult for developers to interpret and debug and may lead to execution patterns that deviate significantly from what is expected.

7.2.3 Directed Program Synthesis

Chapter 6 presents BENCHDIRECT, a directed language model for compiler benchmarks synthesis. BENCHDIRECT is an extension of BENCHPRESS and improves its accuracy of targeting benchmarks by up to 36%, while reducing its inference time overhead by up to 72%.

BENCHDIRECT uses a Transformer-Encoder to influence the language model's token predictions with respect to the desired features. The Transformer receives a concatenation of all supported feature spaces as an input. We use three feature spaces that are commonly used in compiler optimisation to train and evaluate BENCHDIRECT. One limitation is how BENCHDIRECT's sampling process scales with respect to the input features' number of dimensions. While usually features for compilers are low-dimensional (i.e. under 100 features), this architecture is likely to have a considerable training and sampling overhead if many feature spaces are integrated, increasing the input sequence's length. As such, it is unlikely a many-dimensional space can be supported, as it would increase quadratically the training complexity of the Transformer. A different technique is needed to represent effectively a large number of different features without making the training and inference process prohibitively large.

7.3 Future Work

This Section outlines 2 promising directions for future research enabled by this thesis.

7.3.1 Efficient Directed Program Synthesis

As discussed in Chapters 5 and 6, directed synthesizers strongly outperform undirected approaches in delivering high-quality source code. However, the existing technique for guided synthesis using language models is inefficient, requiring thousands of benchmarks to slowly converge to the target features. `BENCHDIRECT` offers a more efficient approach to steerable generation compared to `BENCHPRESS` but the improvement is not exponential.

In most generative language models, there is a trend of steadily increasing the sizes of the architectures and the amount of training data. While larger models e.g., Large Language Models (LLMs), do perform significantly better, we suggest generative models have to be re-designed to emulate the efficiency of humans. This is currently allowed with Reinforcement Learning (RL), a method that emulates the way humans learn, by taking actions and observing their consequences through a reward and punishment system. This helps a RL agent understand the series of actions that will maximize its rewards given the environment.

For the task of directed program synthesis, a RL agent can start from an empty sequence (i.e. the starting state of the environment) and iteratively select actions to maximize its reward. In this case, the action space can be the type of edit performed e.g., ‘ADD’, ‘REMOVE’, ‘REPLACE’ or ‘COMPILE’ along with the index of the sequence at which the action takes place. This assimilates the actions taken by human developers when writing software. The reward function could take into account if the current program compiles and how far it is from the target features to guide the agent into the correct solution.

We feel such an approach is a promising way to make directed program synthesis more accurate and at the same time reduce its inference overhead to a tiny fraction of what it currently is. We are very eager to work in this direction in the future.

7.3.2 Autonomous Predictive Models

To maximize their performance, existing predictive models require human expertise for the task of (a) feature selection, i.e. design their input feature space and (b) curation

of their training data, i.e. ensure they are trained on high-quality datasets. Feature selection is based on human intuition and expertise which often leads to inaccurate predictive models because important features have been removed from the original training data. Finding lots of training data is also a great challenge. Even popular machine learning applications, such as semantic segmentation, require several hours of human labour for labelling.

Both feature selection and dataset curation require expertise and manual effort but they can be automated through autonomously-trained predictive models. In Chapter 5, we show how the programs generated by BENCHPRESS improve a predictive model's accuracy in predicting the fastest execution device for OpenCL workloads. As such, a generative model can be used as a teacher for the predictive model's learning process. To explore the feature space and find these areas that are likely to improve the predictive model, active learning can be used in conjunction with a generative model that will label the active learner's queries.

Neural networks have shown to be successful in representing efficiently raw information into fixed vectors on the latent space. Instead of performing manual feature selection for a predictive model, the feature extraction can be left as a machine learning task, where the model itself decides those features that are sensitive to the classification decision. Taking the device mapping predictive model as an example, instead of hand-crafting the input feature space of the model, it can receive BENCHPRESS's hidden state representation of a benchmark. An active learning algorithm such as EER (Expected Error Reduction) can sample the predictive model with unlabelled data and compute a query. The query can be satisfied by BENCHPRESS, which will synthesize a compiler benchmark with the desired features. Such an approach would require minimal supervision from a human expert and no starting training data. A pre-trained generative model could teach autonomously a predictive model for any downstream task applied to source code.

We believe this is an exciting direction for machine learning and we aim to explore its applications to compiler optimisation.

7.4 Concluding Remarks

The increasing scale of software's complexity today leads to testing and optimising programs being two challenging and time-consuming processes that rely on manual effort and expertise. Addressing these challenges requires developing new automated

tools that will make software faster and safer and will boost developers' productivity.

This thesis leverages the growth of deep language models to propose novel approaches for program correctness classification, program generation and compiler optimisation that are significantly simpler and more effective than established approaches. The methodologies are applied across two domains: software testing and compiler optimisation. In both cases, the proposed techniques outperform state-of-the-art methods that have been developed through years of research.

The demonstrated outcomes open new lines of research into programming language modeling, generation and testing of programs through deep learning. Although our results look promising, there is still much work to be done. Promising future research directions include adapting these techniques to problems that have a wider impact further enabling the automation of domains that still rely on manual effort. Such examples are generating general purpose programming languages that are widely used to create software, or enabling our test oracle approach to work across different languages and technologies. I hope the work of this thesis will be useful to researchers in the future in discovering new domains in the field of artificial intelligence for source code modeling. To this direction, I publicly release all the code and data associated with the contributions of this thesis ¹ ².

¹<https://github.com/fivosts/Learning-over-test-executions>

²<https://github.com/fivosts/BenchPress>

Bibliography

- [1] ABDELBAKY, M., PARASHAR, M., KIM, H., JORDAN, K., SACHDEVA, V., SEXTON, J., JAMJOOM, H., SHAE, Z.-Y., PENCHEVA, G., TAVAKOLI, R., AND WHEELER, M. Enabling high-performance computing as a service. *Computer* 45 (10 2012), 72–80.
- [2] AGARAP, A. F. Deep learning using rectified linear units (relu), 2018.
- [3] AGGARWAL, ET AL. A neural net based approach to test oracle. *ACM SIGSOFT Software Engineering Notes* 29, 3 (2004), 1–6.
- [4] AHAMED, S. S. R. Studying the feasibility and importance of software testing: An analysis.
- [5] ALETI, A., BUHNOVA, B., GRUNSKÉ, L., KOZIOLEK, A., AND MEEDENIYA, I. Software architecture optimization methods: A systematic literature review. *IEEE Transactions on Software Engineering* 39, 5 (2013), 658–683.
- [6] ALLAMANIS, M., BARR, E. T., BIRD, C., AND SUTTON, C. Learning Natural Coding Conventions . In *FSE* (2014), pp. 281–293.
- [7] ALLAMANIS, M., BARR, E. T., DEVANBU, P., AND SUTTON, C. A survey of machine learning for big code and naturalness. *ACM Comput. Surv.* 51, 4 (jul 2018).
- [8] ALLAMANIS, M., BROCKSCHMIDT, M., AND KHADEMI, M. Learning to represent programs with graphs. In *ICLR* (2018).
- [9] ALLAMANIS, M., PENG, H., AND SUTTON, C. A Convolutional Attention Network for Extreme Summarization of Source Code . *arXiv:1602.03001* (2016).

- [10] ALLAMANIS, M., PENG, H., AND SUTTON, C. A Convolutional Attention Network for Extreme Summarization of Source Code. In *ICML (2016)*.
- [11] ALMAGHAIRBE, R., AND ROPER, M. Separating passing and failing test executions by clustering anomalies. *Software Quality Journal* 25, 3 (2017), 803–840.
- [12] ALON, U., ET AL. code2vec: Learning distributed representations of code. *arXiv preprint arXiv:1803.09473* (2018).
- [13] ALON, U., LEVY, O., AND YAHAV, E. code2seq: Generating sequences from structured representations of code. *arXiv preprint arXiv:1808.01400* (2018).
- [14] ALON, U., ZILBERSTEIN, M., LEVY, O., AND YAHAV, E. Code2vec: Learning distributed representations of code. *Proc. ACM Program. Lang.* 3, POPL (jan 2019).
- [15] AMMANN, P., AND OFFUTT, J. *Introduction to software testing*. Cambridge Univ. Press, 2016.
- [16] ANDREWS, J. H., BRIAND, L. C., LABICHE, Y., AND NAMIN, A. S. Using mutation analysis for assessing and comparing testing coverage criteria. *IEEE Transactions on Software Engineering* 32, 8 (2006), 608–624.
- [17] Commons lang. <https://commons.apache.org/proper/commons-lang/>, 2020.
- [18] BAGRODIA, R., MEYER, R., TAKAI, M., CHEN, Y.-A., ZENG, X., MARTIN, J., AND SONG, H. Y. Parsec: a parallel simulation environment for complex systems. *Computer* 31, 10 (1998), 77–85.
- [19] BALOG, M., GAUNT, A. L., BROCKSCHMIDT, M., NOWOZIN, S., AND TARLOW, D. Deepcoder: Learning to write programs. In *International Conference on Learning Representations (2017)*.
- [20] BARR, E., ET AL. The oracle problem in software testing: A survey. *IEEE TSE* 41, 5 (2015), 507–525.
- [21] BEBIS, G., AND GEORGIPOULOS, M. Feed-forward neural networks. *IEEE Potentials* 13, 4 (1994), 27–31.

- [22] BENCHMARKS, R. <http://lava.cs.virginia.edu/Rodinia/download.htm>. [Online; accessed 25-Apr-2022].
- [23] BENGIO, Y., DUCHARME, R., VINCENT, P., AND JANVIN, C. A neural probabilistic language model. *J. Mach. Learn. Res.* 3, null (mar 2003), 1137–1155.
- [24] BERTOLINO, A. Software testing research: Achievements, challenges, dreams. In *Future of Software Engineering* (2007), IEEE Computer Society, pp. 85–103.
- [25] BOWRING, J., ET AL. Active learning for automatic classification of software behavior. In *ACM SIGSOFT Software Engineering Notes* (2004), pp. 195–205.
- [26] BRANTS, T., POPAT, A. C., XU, P., OCH, F. J., AND DEAN, J. Large language models in machine translation. In *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)* (Prague, Czech Republic, June 2007), Association for Computational Linguistics, pp. 858–867.
- [27] BRIAND, L. C. Novel applications of machine learning in software testing. In *QSIC'08* (2008), IEEE, pp. 3–10.
- [28] BRUN, Y., AND ERNST, M. D. Finding latent code errors via machine learning over program executions. In *Proceedings of the 26th ICSE* (2004), pp. 480–490.
- [29] BUDUMA, N., AND LOCASCIO, N. *Fundamentals of Deep Learning: Designing Next-Generation Machine Intelligence Algorithms*, 1st ed. O'Reilly Media, Inc., 2017.
- [30] BUTERIN, V. *Ethereum Project (release 3.5)*, 2019. <https://github.com/ethereum/aleth>.
- [31] CHABOT, M., MAZET, K., AND PIERRE, L. Automatic and configurable instrumentation of c programs with temporal assertion checkers. In *2015 ACM/IEEE International Conference on Formal Methods and Models for Code-sign (MEMOCODE)* (2015), pp. 208–217.
- [32] CHAUDHARI, S., MITHAL, V., POLATKAN, G., AND RAMANATH, R. An attentive survey of attention models. *ACM Trans. Intell. Syst. Technol.* 12, 5 (oct 2021).

- [33] CHE, S., BOYER, M., MENG, J., TARJAN, D., SHEAFFER, J. W., LEE, S.-H., AND SKADRON, K. Rodinia: A benchmark suite for heterogeneous computing. In *2009 IEEE International Symposium on Workload Characterization (IISWC)* (2009), pp. 44–54.
- [34] CHEN, T., ET AL. An orchestrated survey on automated software test case generation. *Journal of Systems and Software* (2013).
- [35] L7-filter, application layer packet classifier for linux. <http://l7-filter.clearos.com/>, 2013.
- [36] COLLIE, B., GINSBACH, P., WOODRUFF, J., RAJAN, A., AND O’BOYLE, M. F. M3: Semantic api migrations. In *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)* (2020), IEEE, pp. 90–102.
- [37] COLLINS, A., FENSCH, C., AND LEATHER, H. Auto-Tuning Parallel Skeletons . *Parallel Processing Letters* 22, 02 (6 2012), 1240005.
- [38] CONG, J., HUANG, M., WU, D., AND YU, C. H. Invited - heterogeneous datacenters: Options and opportunities. In *Proceedings of the 53rd Annual Design Automation Conference* (New York, NY, USA, 2016), DAC ’16, Association for Computing Machinery.
- [39] CUMMINS, C., PETOUMENOS, P., STEUWER, M., AND LEATHER, H. Towards Collaborative Performance Tuning of Algorithmic Skeletons . In *HLPGPU* (2016).
- [40] CUMMINS, C., PETOUMENOS, P., WANG, Z., AND LEATHER, H. Synthesizing benchmarks for predictive modeling. In *2017 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)* (2017), pp. 86–99.
- [41] CUMMINS, C., WASTI, B., GUO, J., CUI, B., ANSEL, J., GOMEZ, S., JAIN, S., LIU, J., TEYTAUD, O., STEINER, B., TIAN, Y., AND LEATHER, H. Compilergym: Robust, performant compiler optimization environments for ai research, 2021.
- [42] DA SILVA, A. F., KIND, B. C., DE SOUZA MAGALHÃES, J. W., ROCHA, J. N., FERREIRA GUIMARÃES, B. C., AND QUINÃO PEREIRA, F. M. Ang-

- habenbench: A suite with one million compilable c benchmarks for code-size reduction. In *2021 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)* (2021), pp. 378–390.
- [43] DAVID, Y., ALON, U., AND YAHAV, E. Neural Reverse Engineering of Stripped Binaries. *arXiv:1902.09122* (2019).
- [44] DE BRUIN, S., LIVENTSEV, V., AND PETKOVIĆ, M. Autoencoders as tools for program synthesis, 2021.
- [45] DENG, J., DONG, W., SOCHER, R., LI, L.-J., LI, K., AND FEI-FEI, L. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition* (2009), pp. 248–255.
- [46] DEVLIN, J., CHANG, M.-W., LEE, K., AND TOUTANOVA, K. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL* (2019).
- [47] DEVLIN, J., UESATO, J., BHUPATIRAJU, S., SINGH, R., MOHAMED, A.-R., AND KOHLI, P. Robustfill: Neural program learning under noisy i/o. In *Proceedings of the 34th International Conference on Machine Learning - Volume 70* (2017), ICML'17, JMLR.org, p. 990–998.
- [48] DINELLA, E., RYAN, G., MYTKOWICZ, T., AND LAHIRI, S. K. Toga: A neural method for test oracle generation. In *Proceedings of the 44th International Conference on Software Engineering* (New York, NY, USA, 2022), ICSE '22, Association for Computing Machinery, p. 2130–2141.
- [49] DO, H., ELBAUM, S., AND ROTHERMEL, G. Supporting controlled experimentation with testing techniques: An infrastructure and its potential impact. *Empirical Software Engineering* 10, 4 (2005), 405–435.
- [50] ELMAN, J. L. Finding structure in time. *Cognitive Science* 14, 2 (1990), 179–211.
- [51] FALCH, T. L., AND ELSTER, A. C. Machine Learning Based Auto-tuning for Enhanced OpenCL Performance Portability . In *IPDPSW* (2015), IEEE.
- [52] FENG, Z., GUO, D., TANG, D., DUAN, N., FENG, X., GONG, M., SHOU, L., QIN, B., LIU, T., JIANG, D., AND ZHOU, M. Codebert: A pre-trained model for programming and natural languages, 2020.

- [53] FIVOSTS. <https://github.com/fivosts/BenchPress.git>. [Online; accessed 1-Sept-2022].
- [54] FRIED, D., AGHAJANYAN, A., LIN, J., WANG, S., WALLACE, E., SHI, F., ZHONG, R., YIH, W.-T., ZETTLEMOYER, L., AND LEWIS, M. Incoder: A generative model for code infilling and synthesis, 2022.
- [55] FURSIN, G., AND TEMAM, O. Collective optimization: A practical collaborative approach. *ACM Trans. Archit. Code Optim.* 7, 4 (dec 2011).
- [56] GALASSI, A., LIPPI, M., AND TORRONI, P. Attention in natural language processing. *IEEE Transactions on Neural Networks and Learning Systems* 32, 10 (2021), 4291–4308.
- [57] GAO, J., TSAO, J., WU, Y., AND JACOB, T. H.-S. Testing and quality assurance for component-based software.
- [58] GITHUB. <https://docs.github.com/en/rest>. [Online; accessed 25-Apr-2022].
- [59] GOENS, A., BRAUCKMANN, A., ERTEL, S., CUMMINS, C., LEATHER, H., AND CASTRILLON, J. A case study on machine learning for synthesizing benchmarks. In *Proceedings of the 3rd ACM SIGPLAN International Workshop on Machine Learning and Programming Languages* (New York, NY, USA, 2019), MAPL 2019, Association for Computing Machinery, p. 38–46.
- [60] GOOGLE. <https://cloud.google.com/bigquery>. [Online; accessed 25-Apr-2022].
- [61] GREWE, D., WANG, Z., AND O’BOYLE, M. F. P. Portable mapping of data parallel programs to opencl for heterogeneous systems. In *Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO)* (2013), pp. 1–10.
- [62] GUO, D., SVYATKOVSKIY, A., YIN, J., DUAN, N., BROCKSCHMIDT, M., AND ALLAMANIS, M. Learning to complete code with sketches. In *International Conference on Learning Representations* (2021).
- [63] GUO, Y., LI, P., LUO, Y., WANG, X., AND WANG, Z. Exploring gnn based program embedding technologies for binary related tasks. In *Proceedings*

of the 30th IEEE/ACM International Conference on Program Comprehension (New York, NY, USA, 2022), ICPC '22, Association for Computing Machinery, p. 366–377.

- [64] GUPTA, K., CHRISTENSEN, P. E., CHEN, X., AND SONG, D. Synthesize, execute and debug: Learning to repair for neural program synthesis. In *Proceedings of the 34th International Conference on Neural Information Processing Systems* (Red Hook, NY, USA, 2020), NIPS'20, Curran Associates Inc.
- [65] HAJ-ALI, A., HUANG, Q. J., XIANG, J., MOSES, W., ASANOVIC, K., WAWRZYNEK, J., AND STOICA, I. Autophase: Juggling hls phase orderings in random forests with deep reinforcement learning. In *Proceedings of Machine Learning and Systems* (2020), I. Dhillon, D. Papailiopoulos, and V. Sze, Eds., vol. 2, pp. 70–81.
- [66] HAN, J., AND MORAGA, C. The influence of the sigmoid function parameters on the speed of backpropagation learning. In *From Natural to Artificial Neural Computation* (Berlin, Heidelberg, 1995), J. Mira and F. Sandoval, Eds., Springer Berlin Heidelberg, pp. 195–201.
- [67] HARIRI, F., AND SHI, A. Srciror: A toolset for mutation testing of c source code and llvm intermediate representation. In *2018 33rd IEEE/ACM International Conference on Automated Software Engineering (ASE)* (2018), pp. 860–863.
- [68] HENDRYCKS, D., AND GIMPEL, K. Gaussian error linear units (gelus), 2016.
- [69] HENNING, J. L. Spec cpu2006 benchmark descriptions. *SIGARCH Comput. Archit. News* 34, 4 (sep 2006), 1–17.
- [70] HIERONS, R. M. Verdict functions in testing with a fault domain or test hypotheses. *ACM TOSEM* 18, 4 (2009), 14.
- [71] HIERONS, R. M. Oracles for distributed testing. *IEEE TSE* 38, 3 (2012), 629–641.
- [72] HINDLE, A., BARR, E. T., GABEL, M., SU, Z., AND DEVANBU, P. On the naturalness of software. *Commun. ACM* 59, 5 (apr 2016), 122–131.

- [73] HINDLE, A., BARR, E. T., SU, Z., GABEL, M., AND DEVANBU, P. On the naturalness of software. In *Proceedings of the 34th International Conference on Software Engineering* (2012), ICSE '12, IEEE Press, p. 837–847.
- [74] HINTON, G., ET AL. A fast learning algorithm for deep belief nets. *Neural computation* 18, 7 (2006), 1527–1554.
- [75] HINTON, G. E., MCCLELLAND, J. L., AND RUMELHART, D. E. *Distributed Representations*. MIT Press, Cambridge, MA, USA, 1986, p. 77–109.
- [76] HOCHREITER, S. The vanishing gradient problem during learning recurrent neural nets and problem solutions. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems* 6 (04 1998), 107–116.
- [77] HOCHREITER, S., AND SCHMIDHUBER, J. Long short-term memory. *Neural Comput.* 9, 8 (nov 1997), 1735–1780.
- [78] HOCHREITER, S., AND SCHMIDHUBER, J. Long short-term memory. *Neural Comput.* 9, 8 (nov 1997), 1735–1780.
- [79] HOCHREITER, S., AND SCHMIDHUBER, J. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [80] HOLLINGSWORTH, J., MILLER, B., AND CARGILLE, J. Dynamic program instrumentation for scalable performance tools. In *Proceedings of IEEE Scalable High Performance Computing Conference* (1994), pp. 841–850.
- [81] HORVÁTH, F., GERGELY, T., ÁRPÁD BESZÉDES, TENGERI, D., BALOGH, G., AND GYIMÓTHY, T. Code coverage differences of java bytecode and source code instrumentation tools. *Software Quality Journal* 27, 1 (2019), 79–123.
- [82] HOWDEN, W. Theoretical and empirical studies of program testing. *IEEE Transactions on Software Engineering SE-4*, 4 (1978), 293–298.
- [83] HYUNSOOK DO, AND ROTHERMEL, G. On the use of mutation faults in empirical assessments of test case prioritization techniques. *IEEE Transactions on Software Engineering* 32, 9 (2006), 733–752.
- [84] JALOTE, P. *An Integrated Approach to Software Engineering*, 3rd ed. Springer Publishing Company, Incorporated, 2010.

- [85] JIA, Y., AND HARMAN, M. An analysis and survey of the development of mutation testing. *IEEE Transactions on Software Engineering* 37, 5 (2011), 649–678.
- [86] JIA, Y., AND HARMAN, M. An analysis and survey of the development of mutation testing. *IEEE transactions on software engineering* 37, 5 (2011), 649–678.
- [87] JIN, H., ET AL. Artificial neural network for automatic test oracles generation. In *Proceedings of CSSE* (2008), vol. 2, IEEE, pp. 727–730.
- [88] JING, K., AND XU, J. A survey on neural network language models, 2019.
- [89] JOUPPI, N. P., YOUNG, C., PATIL, N., PATTERSON, D., AGRAWAL, G., BAJWA, R., BATES, S., BHATIA, S., BODEN, N., BORCHERS, A., BOYLE, R., CANTIN, P.-L., CHAO, C., CLARK, C., CORIELL, J., DALEY, M., DAU, M., DEAN, J., GELB, B., GHAEMMAGHAMI, T. V., GOTTIPATI, R., GULLAND, W., HAGMANN, R., HO, C. R., HOGBERG, D., HU, J., HUNDT, R., HURT, D., IBARZ, J., JAFFEY, A., JAWORSKI, A., KAPLAN, A., KHAITAN, H., KILLEBREW, D., KOCH, A., KUMAR, N., LACY, S., LAUDON, J., LAW, J., LE, D., LEARY, C., LIU, Z., LUCKE, K., LUNDIN, A., MACKEAN, G., MAGGIORE, A., MAHONY, M., MILLER, K., NAGARAJAN, R., NARAYANASWAMI, R., NI, R., NIX, K., NORRIE, T., OMERNICK, M., PENUKONDA, N., PHELPS, A., ROSS, J., ROSS, M., SALEK, A., SAMADIANI, E., SEVERN, C., SIZIKOV, G., SNELHAM, M., SOUTER, J., STEINBERG, D., SWING, A., TAN, M., THORSON, G., TIAN, B., TOMA, H., TUTTLE, E., VASUDEVAN, V., WALTER, R., WANG, W., WILCOX, E., AND YOON, D. H. In-datacenter performance analysis of a tensor processing unit. *SIGARCH Comput. Archit. News* 45, 2 (jun 2017), 1–12.
- [90] JUST, R., JALALI, D., AND ERNST, M. D. Defects4J: A database of existing faults to enable controlled testing studies for java programs. In *ISSTA 2014, Proceedings of the 2014 International Symposium on Software Testing and Analysis* (2014), pp. 437–440.
- [91] KANADE, A., MANIATIS, P., BALAKRISHNAN, G., AND SHI, K. Learning and evaluating contextual embedding of source code, 2020.

- [92] KANER, C., FALK, J. L., AND NGUYEN, H. Q. *Testing Computer Software, Second Edition*, 2nd ed. John Wiley; Sons, Inc., USA, 1999.
- [93] KELLERER, H., PFERSCHY, U., AND PISINGER, D. *Introduction*. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004, pp. 1–14.
- [94] KINGMA, D. P., AND BA, J. Adam: A method for stochastic optimization. In *3rd International Conference for Learning Representations (2015)*.
- [95] KOROTEEV, M. V. Bert: A review of applications in natural language processing and understanding. *ArXiv abs/2103.11943 (2021)*.
- [96] LANGDON, W., ET AL. Inferring automatic test oracles. In *Proceedings of the 10th SBST (2017)*, pp. 5–6.
- [97] LATTNER, C. LLVM: An Infrastructure for Multi-Stage Optimization. Master’s thesis, Computer Science Dept., University of Illinois at Urbana-Champaign, Urbana, IL, Dec 2002. See <http://llvm.cs.uiuc.edu>.
- [98] LATTNER, C., AND ADVE, V. LLVM: A compilation framework for lifelong program analysis and transformation. In *CGO (San Jose, CA, USA, Mar 2004)*, pp. 75–88.
- [99] LEATHER, H., BONILLA, E., AND O’BOYLE, M. Automatic Feature Generation for Machine Learning Based Optimizing Compilation . *TACO 11 (2014)*.
- [100] LECUN, Y., BENGIO, Y., AND HINTON, G. Deep learning. *Nature 521 (05 2015)*, 436–44.
- [101] LI, L. H., YATSKAR, M., YIN, D., HSIEH, C.-J., AND CHANG, K.-W. Visualbert: A simple and performant baseline for vision and language, 2019.
- [102] LI, Y., CHOI, D., CHUNG, J., KUSHMAN, N., SCHRITTWIESER, J., LEBLOND, R., ECCLES, T., KEELING, J., GIMENO, F., LAGO, A. D., HUBERT, T., CHOY, P., DE MASSON D’AUTUME, C., BABUSCHKIN, I., CHEN, X., HUANG, P.-S., WELBL, J., GOWAL, S., CHEREPANOV, A., MOLLOY, J., MANKOWITZ, D. J., ROBSON, E. S., KOHLI, P., DE FREITAS, N., KAVUKCUOGLU, K., AND VINYALS, O. Competition-level code generation with alphacode. *Science 378, 6624 (2022)*, 1092–1097.

- [103] LIDBURY, C. <https://github.com/ChrisLidbury/CLSmith>. [Online; accessed 30-Jan-2023].
- [104] Sed, linux stream editor. <https://linux.die.net/man/1/sed>, 2009.
- [105] LIU, H., KUO, F., TOWEY, D., AND CHEN, T. Y. How effectively does metamorphic testing alleviate the oracle problem? *IEEE Transactions on Software Engineering* 40, 1 (2014), 4–22.
- [106] <https://clang.llvm.org/docs/LibTooling.html>. [Online; accessed 30-Jan-2023].
- [107] MA, W., ZHAO, M., SOREMEKUN, E., HU, Q., ZHANG, J. M., PAPADAKIS, M., CORDY, M., XIE, X., AND TRAON, Y. L. Graphcode2vec: Generic code embedding via lexical and program dependence analyses. In *Proceedings of the 19th International Conference on Mining Software Repositories* (New York, NY, USA, 2022), MSR '22, Association for Computing Machinery, p. 524–536.
- [108] MAGNI, A., DUBACH, C., AND O'BOYLE, M. Automatic Optimization of Thread-Coarsening for Graphics Processors. In *PACT (2014)*, ACM, pp. 455–466.
- [109] MAHRENHOLZ, D., SPINCZYK, O., AND SCHRODER-PREIKSCHAT, W. Program instrumentation for debugging and monitoring with aspectc++. In *Proceedings Fifth IEEE International Symposium on Object-Oriented Real-Time Distributed Computing. ISIRC 2002* (2002), pp. 249–256.
- [110] MELO, L. T. C., RIBEIRO, R. G., GUIMARÃES, B. C. F., AND PEREIRA, F. M. Q. A. Type inference for c: Applications to the static analysis of incomplete programs. *ACM Trans. Program. Lang. Syst.* 42, 3 (nov 2020).
- [111] MICOLET, P., SMITH, A., AND DUBACH, C. A Machine Learning Approach to Mapping Streaming Workloads to Dynamic Multicore Processors. In *LCTES* (2016).
- [112] MIKOLOV, T., CHEN, K., CORRADO, G., AND DEAN, J. Efficient estimation of word representations in vector space, 2013.
- [113] MIKOLOV, T., KOPECKY, J., BURGET, L., GLEMBEK, O., AND ?CERNOCKY, J. Neural network based language models for highly inflective languages. In

- 2009 *IEEE International Conference on Acoustics, Speech and Signal Processing* (2009), pp. 4725–4728.
- [114] MISRA, J., AND SAHA, I. Artificial neural networks in hardware: A survey of two decades of progress. *Neurocomputing* 74, 1 (2010), 239–255. Artificial Brains.
- [115] MOHRI, M., ROSTAMIZADEH, A., AND TALWALKAR, A. *Foundations of Machine Learning*. The MIT Press, 2012.
- [116] MURPHY, K. P. *Machine Learning: A Probabilistic Perspective*. The MIT Press, 2012.
- [117] NARAYANAN, A., CHANDRAMOHAN, M., VENKATESAN, R., CHEN, L., LIU, Y., AND JAISWAL, S. graph2vec: Learning distributed representations of graphs. *ArXiv abs/1707.05005* (2017).
- [118] NARDI, P. A., AND DAMASCENO, E. A survey on test oracles. *Advances in Theoretical and Applied Informatics* 1, 2 (2015), 50–59.
- [119] NYE, M., HEWITT, L. B., TENENBAUM, J. B., AND SOLAR-LEZAMA, A. Learning to infer program sketches. In *ICML* (2019).
- [120] OGILVIE, W. F., PETOUMENOS, P., WANG, Z., AND LEATHER, H. Fast Automatic Heuristic Construction Using Active Learning . In *LCPC* (2014).
- [121] OPENCL-SPECIFICATION. https://www.khronos.org/registry/OpenCL/specs/3.0-unified/html/OpenCL_C.html. [Online; accessed 25-Apr-2022].
- [122] PAPINENI, K., ROUKOS, S., WARD, T., AND ZHU, W.-J. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for Computational Linguistics (USA, 2002)*, ACL '02, Association for Computational Linguistics, p. 311–318.
- [123] PASZKE, A., GROSS, S., MASSA, F., LERER, A., BRADBURY, J., CHANAN, G., KILLEEN, T., LIN, Z., GIMELSHEIN, N., ANTIGA, L., DESMAISON, A., KOPE, A., YANG, E., DEVITO, Z., RAISON, M., TEJANI, A., CHILAMKURTHY, S., STEINER, B., FANG, L., BAI, J., AND CHINTALA, S. Pytorch: An imperative style, high-performance deep learning library. In *Advances*

- in Neural Information Processing Systems 32*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, Eds. Curran Associates, Inc., 2019, pp. 8024–8035.
- [124] PENNINGTON, J., SOCHER, R., AND MANNING, C. GloVe: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)* (Doha, Qatar, Oct. 2014), Association for Computational Linguistics, pp. 1532–1543.
- [125] PETERS, M., NEUMANN, M., IYYER, M., GARDNER, M., CLARK, C., LEE, K., AND ZETTLEMOYER, L. Deep contextualized word representations.
- [126] PODGURSKI, A., ET AL. Automated support for classifying software failure reports. In *Proceedings of 25th ICSE 2003*. (2003), IEEE, pp. 465–475.
- [127] PRADEL, M., AND SEN, K. Deepbugs: a learning approach to name-based bug detection. *Proceedings of the ACM on Programming Languages 2*, OOPSLA (2018), 147.
- [128] RADFORD, A., AND NARASIMHAN, K. Improving language understanding by generative pre-training.
- [129] RAYCHEV, V., VECHEV, M., AND KRAUSE, A. Predicting Program Properties from "Big Code". In *POPL* (2015).
- [130] RIGGER, M., AND SU, Z. Intramorphic testing: A new approach to the test oracle problem. In *Proceedings of the 2022 ACM SIGPLAN International Symposium on New Ideas, New Paradigms, and Reflections on Programming and Software* (New York, NY, USA, 2022), Onward! 2022, Association for Computing Machinery, p. 128–136.
- [131] ROY, N., AND MCCALLUM, A. Toward optimal active learning through monte carlo estimation of error reduction.
- [132] RUMELHART, D. E., HINTON, G. E., AND WILLIAMS, R. J. Learning representations by back-propagating errors. In *Nature* 323, 533–536 (1986). <https://doi.org/10.1038/323533a0> (1986).
- [133] RUMELHART, D. E., AND MCCLELLAND, J. L. *Learning Internal Representations by Error Propagation*. 1987, pp. 318–362.

- [134] Microsoft SEAL (release 3.2). <https://github.com/Microsoft/SEAL>, 2019. Microsoft Research, Redmond, WA.
- [135] SEUNG, H. S., OPPER, M., AND SOMPOLINSKY, H. Query by committee. In *Proceedings of the Fifth Annual Workshop on Computational Learning Theory* (New York, NY, USA, 1992), COLT '92, Association for Computing Machinery, p. 287–294.
- [136] SOYDANER, D. Attention mechanism in neural networks: where it comes and where it goes. *Neural Computing and Applications* 34, 16 (may 2022), 13371–13385.
- [137] SRIVASTAVA, N., HINTON, G., KRIZHEVSKY, A., SUTSKEVER, I., AND SALAKHUTDINOV, R. Dropout: A simple way to prevent neural networks from overfitting. *J. Mach. Learn. Res.* 15, 1 (jan 2014), 1929–1958.
- [138] STEPHENSON, M., SASTRY HARI, S. K., LEE, Y., EBRAHIMI, E., JOHNSON, D. R., NELLANS, D., O'CONNOR, M., AND KECKLER, S. W. Flexible software profiling of gpu architectures. *SIGARCH Comput. Archit. News* 43, 3S (jun 2015), 185–197.
- [139] STONE, J. E., GOHARA, D., AND SHI, G. Opencl: A parallel programming standard for heterogeneous computing systems. *Computing in Science Engineering* 12, 3 (2010), 66–73.
- [140] TAI, K. S., SOCHER, R., AND MANNING, C. D. Improved semantic representations from tree-structured long short-term memory networks. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)* (Beijing, China, July 2015), Association for Computational Linguistics, pp. 1556–1566.
- [141] TEMPLER, K., AND JEFFERY, C. A configurable automatic instrumentation tool for ansi c. In *Proceedings 13th IEEE International Conference on Automated Software Engineering (Cat. No.98EX239)* (1998), pp. 249–258.
- [142] TSIMPOURLAS, F., PAPADOPOULOS, L., BARTSOKAS, A., AND SOUDRIS, D. A design space exploration framework for convolutional neural networks

- implemented on edge devices. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* 37, 11 (2018), 2212–2221.
- [143] TSIMPOURLAS, F., PETOUMENOS, P., XU, M., CUMMINS, C., HAZELWOOD, K., RAJAN, A., AND LEATHER, H. Benchdirect: A directed language model for compiler benchmarks, 2023.
- [144] TSIMPOURLAS, F., PETOUMENOS, P., XU, M., CUMMINS, C., HAZELWOOD, K., RAJAN, A., AND LEATHER, H. Benchpress: A deep active benchmark generator. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques* (New York, NY, USA, 2023), PACT '22, Association for Computing Machinery, p. 505–516.
- [145] TSIMPOURLAS, F., RAJAN, A., AND ALLAMANIS, M. Supervised learning over test executions as a test oracle. In *Proceedings of the 36th Annual ACM Symposium on Applied Computing* (New York, NY, USA, 2021), SAC '21, Association for Computing Machinery, p. 1521–1531.
- [146] TSIMPOURLAS, F., ROOIJACKERS, G., RAJAN, A., AND ALLAMANIS, M. Embedding and classifying test execution traces using neural networks. *IET Software* 16, 3 (2022), 301–316.
- [147] VALLÉE-RAI, R., CO, P., GAGNON, E., HENDREN, L., LAM, P., AND SUNDARESAN, V. Soot - a java bytecode optimization framework. In *Proceedings of the 1999 Conference of the Centre for Advanced Studies on Collaborative Research* (1999), CASCON '99, IBM Press, p. 13.
- [148] VANMALI, M., ET AL. Using a neural network in the software testing process. *International Journal of Intelligent Systems* 17, 1 (2002), 45–62.
- [149] VASWANI, A., SHAZEER, N., PARMAR, N., USZKOREIT, J., JONES, L., GOMEZ, A. N., KAISER, L., AND POLOSUKHIN, I. Attention is all you need. In *Proceedings of the 31st International Conference on Neural Information Processing Systems* (Red Hook, NY, USA, 2017), NIPS'17, Curran Associates Inc., p. 6000–6010.
- [150] WANG, K., AND CHRISTODORESCU, M. Coset: A benchmark for evaluating neural program embeddings, 2019.

- [151] WANG, K., SINGH, R., AND SU, Z. Dynamic neural program embedding for program repair. *ICLR* (2018).
- [152] WANG, K., AND SU, Z. Blended, precise semantic program embeddings. In *Proceedings of the 41st ACM SIGPLAN Conference on Programming Language Design and Implementation* (New York, NY, USA, 2020), PLDI 2020, Association for Computing Machinery, p. 121–134.
- [153] WANG, Z., AND O’BOYLE, M. Mapping Parallelism to Multi-cores: A Machine Learning Based Approach . In *PPoPP* (2009), no. 15, ACM, pp. 75–84.
- [154] WANG, Z., AND O’BOYLE, M. Partitioning Streaming Parallelism for Multi-cores: A Machine Learning Based Approach . In *PACT* (2010), ACM, pp. 307–318.
- [155] WANG, Z., AND O’BOYLE, M. Machine learning in compiler optimization. *Proceedings of the IEEE* 106, 11 (2018), 1879–1901.
- [156] WANG, Z., SANCHEZ, A., AND HERKERSDORF, A. Scisim: A software performance estimation framework using source code instrumentation. In *Proceedings of the 7th International Workshop on Software and Performance* (New York, NY, USA, 2008), WOSP ’08, Association for Computing Machinery, p. 33–42.
- [157] WANG, Z., TOURNAVITIS, G., FRANKE, B., AND O’BOYLE, M. Integrating Profile-driven Parallelism Detection and Machine-learning-based Mapping . *TACO* (2014).
- [158] WEN, Y., WANG, Z., AND O’BOYLE, M. Smart Multi-Task Scheduling for OpenCL Programs on CPU/GPU Heterogeneous Platforms . In *HiPC* (2014), IEEE.
- [159] WOLVERTON, R. The cost of developing large-scale software. *IEEE Transactions on Computers* C-23, 6 (1974), 615–636.
- [160] WONG, E., YANG, J., AND TAN, L. AutoComment: Mining Question and Answer Sites for Automatic Comment Generation . In *ASE* (2013), IEEE, pp. 562–567.

- [161] XU, X., WANG, X., AND XUE, J. M3v: Multi-modal multi-view context embedding for repair operator prediction. In *Proceedings of the 20th IEEE/ACM International Symposium on Code Generation and Optimization (2022)*, CGO '22, IEEE Press, p. 266–277.
- [162] YANEVA, V., KAPOOR, A., RAJAN, A., AND DUBACH, C. Accelerated finite state machine test execution using gpus. In *APSEC (2018)*.
- [163] YANEVA, V., RAJAN, A., AND DUBACH, C. Compiler-assisted test acceleration on gpus for embedded software. In *Proceedings of the 26th ACM SIGSOFT International Symposium on Software Testing and Analysis (2017)*, pp. 35–45.
- [164] YANG, X., CHEN, Y., EIDE, E., AND REGEHR, J. Finding and understanding bugs in c compilers. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language Design and Implementation (New York, NY, USA, 2011)*, PLDI '11, Association for Computing Machinery, p. 283–294.
- [165] YU, S., AND ZHOU, S. A survey on metric of software complexity. In *2010 2nd IEEE International Conference on Information Management and Engineering (2010)*, pp. 352–356.