



“Learning Highly Recursive Input Grammars”: A Replication Study

Sulav Malla
SID:500495980

A thesis presented in partial fulfilment of the requirements for the degree of
Bachelor of Engineering Honours (Software)

Supervisors:

Dr. Rahul Gopinath, The University of Sydney
Dr. Yash Shirvastava, The University of Sydney

School of Electrical and Computer Engineering
The University of Sydney

November 07, 2025

No part of this work may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the prior permission of the author or The University of Sydney.

Abstract

Knowing the source code of a program greatly aids program comprehension, testing techniques such as fuzzing, optimisation, and debugging. However, due to various restrictions and other external factors, accessing source code is often not possible. To address this limitation, recent research has focused on inferring input grammars and execution behaviour directly from valid program inputs. The ARVADA algorithm, developed by Lemieux C. and Kulkarni N. at UCB in 2021, is a black-box approach designed to infer the grammar of a program using only valid inputs and a black-box oracle.

Motivated by the lack of formal guarantees and concerns about the selectiveness of the original study, this thesis attempts to re-implement the ARVADA algorithm in a clean-room environment using the C programming language. Although a complete implementation was not achieved, the work highlights the challenges faced and insights gained during the re-implementation. These include the inherent complexity of ARVADA, the difficulty of implementing such an algorithm in C, and the limited clarity in the explanation provided in the original paper.

Statement of Achievement

Acknowledgement

I would like to acknowledge my supervisors for this project. Firstly, I would like to thank **Dr. Rahul Gopinath** for providing me with the opportunity to undertake this project and for his continued support when I needed it. I would also like to thank **Dr. Yash Shrivastava** for agreeing to be my ECE supervisor.

Additionally, I wish to express my gratitude to my friends and family for their constant support, patience, and understanding throughout this. There were many high moments, low moments, and times when I felt completely stuck, but at the end am very glad to have taken this project.

Thank you, Dr. Rahul Gopinath and Dr. Yash Shrivastava.

Contents

Abstract	ii
Statement of Achievement	iii
Acknowledgement	iv
Glossary	vii
1 INTRODUCTION	2
2 GENERAL BACKGROUND & REVIEW OF LITERATURE	4
2.0.1 What are Grammars?	4
2.0.2 What are Context-Free Grammars?	4
2.0.3 What is ARVADA?	7
2.0.4 Why replicate ARVADA?	14
2.0.5 Why C?	15
2.0.6 Why is learning input grammar important?	15
2.0.7 Related Work	15
2.0.8 TREEVADA	15
2.0.9 KEDAVRA	16
2.0.10 NATGI	18
2.0.11 Problem Statement	19
3 METHODOLGIES & IMPLEMENTATION	20
3.0.1 Data structures, and Initial Parse Trees	20
3.0.2 Pre-tokenisation	24
3.0.3 MergeAllValid	24
4 EVALUATION	25
5 DISCUSSION	26
6 CONCLUSION	27
References	28
Appendix	30

List of Figures

2.1	Chomsky hierarchy [10]	5
2.2	Definition a simple while grammar G_w , sample input strings S , and oracle \mathcal{O} [2]	9
2.3	Initial set of flat naive parse trees ARVADA builds given inputs S , where each t_i is a non-terminal	10
2.4	ARVADA pre-tokenisations example	10
2.5	Examples of a MERGEALLVALID run, tokenised and non tokenisation	11
2.6	Possible bubble example	12
2.7	Example full run of ARVADA bubbling using non-tokenisation nor optimised tree, and the grammar induced (part 1)[2]	13
2.8	Example full run of ARVADA bubbling using non-tokenisation nor optimised tree, and the grammar induced (part 2)[2]	14
2.9	Naive Parse trees after pre-tokenisation in ARVADA	16
2.10	Naive Parse trees after pre-tokenisation in TREEVADA [12]	16
2.11	Result during KEDAVRA pre-tokenisation given input <code>if(a); else b = 5;</code> . .	17
2.12	Result after KEDAVRA pre-tokenisation given input <code>if(a); else b = 5;</code> , and example generalisation of t_1	17
2.13	Decomposed Sequences of tokenised example from 2.12	18
3.1	Data Structure to store and track trees	20
3.2	Data Structure to build each tree	21
3.3	building of Nodes Structure	22
3.4	Building each initial parse tree	23
3.5	Pre-tokenisation of each root node	24
6.1	Code for building a basic node	30
6.2	Code for changing list capacities appropriately	31
6.3	Code for freeing tree given a root node	32
6.4	Pre-tokeniser code	33
6.5	Pre-tokenisation helper code	34

Glossary

Chapter 1

INTRODUCTION

In the field of software testing, generating test inputs for a program (fuzzing) is a well-known and popular technique. To improve the effectiveness of fuzzers, recent research has focused on recovering input grammars from existing programs, as incorporating knowledge about the input language and grammar of the program under test drastically improves the effectiveness of fuzzers [1].

Learning Highly Recursive Input Grammars, authored by Lemieux C., Sen K., and Kulkarni N., and published in 2021 at UCB [2], presents an algorithm called ARVADA, developed for this purpose. ARVADA attempts to learn context-free grammars (CFGs) of a specified program given a set of valid inputs and a black-box oracle \mathcal{O} . Along with presenting the algorithm, the paper evaluates it by comparing ARVADA to GLADE [3], a previously developed state-of-the-art algorithm for the same task in a similar setting. During the evaluation, it was observed that the F1 scores of GLADE were much lower than those reported in the official GLADE paper [3]. This led to a replication study of the original GLADE paper, conducted by Gopinath R., Bachir B., and Zeller A., and published as “Synthesising Input Grammars: A Replication Study” [4] at CISP, which investigated the accuracy of the results in the original paper. Similarly, as done by the researchers at CISP, this thesis aims to replicate ARVADA to reproduce and investigate the results presented in the original paper [2]. Additionally, also attempting to make improvements in the paper, algorithm, and its evaluation.

This thesis is an attempt to reproduce ARVADA in a clean-room environment, meaning with no reference to or knowledge of the original implementation, but only the abstraction and explanations provided in the paper [2]. The implementation language of choice is C.

First, this paper will introduce general concepts of grammar and parsing, the importance of learning input grammars, a deeper explanation of ARVADA, and the motivation for conducting a replication study. This is followed by a brief explanation of ARVADA and how it works according to the original paper [2], then a more in-depth explanation of the algorithm and how it was reproduced in C with reference to the code. Afterwards, an evaluation compares the results of the reproduced implementation with those from the original paper [2]. Finally, a discussion highlights and comments on the original algorithm, the reproduced results, and the

overall process of conducting the replication study, followed by a brief conclusion.

- All code and implementation are open source and can be found here: <https://github.com/Stainima/ARVADA>

Chapter 2

GENERAL BACKGROUND & REVIEW OF LITERATURE

2.0.1 What are Grammars?

Grammars in computer science can be defined as a set of rules by which valid sentences in a language are constructed [5], serving as a blueprint for the language. Beginning with a start symbol, which is a single non-terminal, production rules are applied sequentially, adding symbols from the alphabet according to the grammar's production rules, to derive a string that is valid in the language [6].

A grammar can be represented by the tuple $\langle N, T, P, S \rangle$, where:

- N is a finite, non-empty set of non-terminal symbols/alphabet.
- T is a finite set of terminal symbols/alphabet.
- P is a finite set of production rules.
- S is the start symbol (a non-terminal symbol/alphabet).

Here, the production rules P define all symbol substitutions that can be performed recursively to generate different symbol sequences, known as *strings* [6]. These rules are written in the form $A \rightarrow w$, where $A \in N$ and $w \in (N \cup T)^*$.

These sets of rules are used by parsers to parse a string of tokens and analyse its syntax against this set of rules (Grammars) [7], [8].

2.0.2 What are Context-Free Grammars?

Depending on the set of production rules, each grammar can be classified according to the *Chomsky hierarchy* [9].

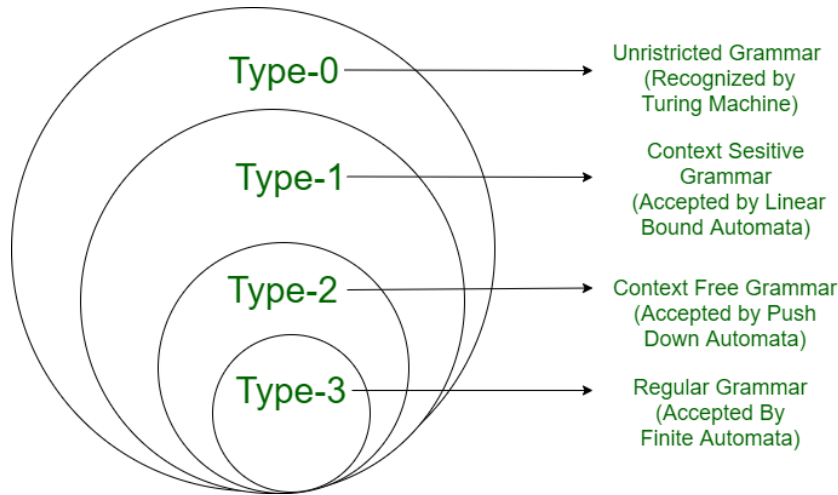


Figure 2.1: Chomsky hierarchy [10]

The hierarchy outlines four distinct types of grammars, ranging from **Type-3** (the most restrictive) to **Type-0** (the most general and unstructured). **Type-2** grammars in this hierarchy are known as *context-free grammars*.

Type-3

Type-3 grammars, known as regular grammars, are used to generate regular languages. A grammar is classified as a regular grammar if its production rules follow from $X \rightarrow a$ or $X \rightarrow aY$. The left-hand side must consist of a single T and the right-hand side must consist of a single T or a single T and a single N .

Regular grammar can take 2 forms, right linear and left linear.

Right-linear takes the following form, where the N on the right-hand side of the arrow is on the far right.

$$X \rightarrow a$$

$$X \rightarrow aB$$

$$X \rightarrow \varepsilon$$

Left-linear take the following form, where the N on the right-hand side of the arrow is on the far left.

$$X \rightarrow a$$

$$X \rightarrow Ba$$

$$X \rightarrow \varepsilon$$

Note that right-linear and left-linear forms cannot be mixed, Doing so may generate languages that are not regular. Additionally, production of ε is allowed only if the corresponding nonterminal does not appear on the right-hand side of any production rule [8], [11]

Type-2

Type-2 grammars, commonly known as context-free grammars (CFGs), have production rules of the form:

$$A \rightarrow \alpha$$

where $A \in N$ and $\alpha \in (T \cup N)^*$, meaning any string composed of terminal and nonterminal symbols [8], [11].

Due to the nature of the right-hand side, nonterminals are allowed to recursively expand, which can lead to repetition. For example:

$$A \rightarrow sAb$$

$$A \rightarrow \varepsilon$$

can produce derivations of the form

$$A \Rightarrow sAb \Rightarrow ssAbb \Rightarrow \dots \Rightarrow s^n b^n,$$

until eventually $A \Rightarrow \varepsilon$.

This property is particularly useful because it allows the grammar to enforce well-formed parentheses expressions and other recursive structures purely through production rules. As a result, context-free grammars are of great importance in the design and specification of programming languages[8], [11].

Type-1

Type-1 grammars, also known as context-sensitive grammars, have production rules of the form:

$$\alpha A \beta \rightarrow \alpha \gamma \beta$$

where $A \in N$ and $\alpha, \beta, \gamma \in (N \cup T)^*$. This means that A can be expanded to γ only in the specific context where it appears between α and β in the given order [8], [11].

A key property of context-sensitive grammars is that production rules must be *non-contracting*: the length of the right-hand side must be greater than or equal to the length of the left-hand side. This implies that no nonterminal on the right-hand side can derive ε , meaning nullable productions are not allowed.

An example derivation using a context-sensitive production is:

$$A \Rightarrow aA\beta \Rightarrow aaA\beta\beta \Rightarrow \dots$$

Type-0

Type-0 grammars, also known as unrestricted grammars, sit at the top of Chomsky's hierarchy. Their production has a complete lack of restrictions and takes the general form:

$$\alpha \rightarrow \beta$$

where $\alpha, \beta \in (N \cup T)^+$, meaning they are strings consisting of terminal and nonterminal symbols. The only requirement is that α must contain at least one nonterminal symbol to allow for further derivations [8], [11].

Type-0 grammars are the most expressive class in the Chomsky hierarchy and can generate all recursively enumerable languages.

2.0.3 What is ARVADA?

Introduction

ARVADA is an algorithm published in “Learning Highly Recursive Input Grammars” [2] at the University of California, Berkeley in 2021. It is designed to learn context-free grammars from a set of positive examples and a Boolean-valued oracle \mathcal{O} . Starting from initially flat parse trees, ARVADA repeatedly applies two specialised operations, **bubbling** and **merging**, to incrementally add structure to these trees. From this structured representation, it extracts the smallest possible set of context-free grammar rules that accommodate all the given examples. The algorithm aims to generalise the language as much as possible without overgeneralising beyond what is accepted by \mathcal{O} .

Like GLADE [3], ARVADA operates under the assumption of a black-box oracle \mathcal{O} . This means that ARVADA has no access to or knowledge of the internal workings of the oracle and can only observe the Boolean values returned by \mathcal{O} .

GLADE

GLADE is an algorithm proposed by Bastani et al., published in Synthesising Input Grammars at PLDI 2017 [3]. Like ARVADA, GLADE uses a set of valid inputs and black-box access to the program, with the aim of automatically approximating the context-free input grammar of the given program.

Because GLADE shares a similar experimental setting as ARVADA and is a predecessor, GLADE was used as a benchmark for ARVADA during its evaluation. It was shown that GLADE, on average, had a faster runtime compared to ARVADA. However, in terms of generalisation, following the original grammar of the oracle more closely, ARVADA outperformed GLADE across the 11 benchmarks, achieving a higher F1 score on 9 of the 11 benchmarks.

Explanation

ARVADA takes as input the oracle \mathcal{O} and a set of positive, valid oracle inputs S . For each string $s \in S$, querying $\mathcal{O}(s)$ returns **True**. The algorithm begins by constructing a flat parse tree for each string in S . Each tree has a single root node t_0 whose children correspond to the individual characters of the input string s .

Next, ARVADA performs the **bubbling** operation. In this step, a sequence of sibling nodes in the tree is selected and replaced with a new non-terminal node. This new node takes the

selected sibling nodes as its children, thereby introducing an additional level of structure. Essentially, ARVADA transforms sequences of terminal nodes in the flat parse tree into subtrees by introducing new non-terminal nodes and progressively adding structure to the tree.

ARVADA then decides whether to accept or reject each bubble by checking whether the newly bubbled structure enables a sound generalisation of the learned grammar. Each non-leaf node in the tree can be viewed as a non-terminal in the emerging grammar. To determine whether a bubble should be accepted, ARVADA checks whether replacing any node in the tree with the new bubbled subtree results in the generation of valid input strings according to \mathcal{O} . If the replacement produces valid strings, the bubble is accepted, and the tree is restructured so that both the bubbled subtree and the replaced node share the same non-terminal label.

The addition of new non-terminal nodes expands the language defined by the learned grammar, since any string derivable from the same label can now be substituted interchangeably. This relabeling of the bubbled subtree and the replaced node is called a **merge**, as it merges the labels of two previously distinct nodes in the tree. If a bubble is not accepted, it is discarded, and none of the trees are affected or structurally modified.

Walkthrough

This walkthrough will follow very closely to the examples provided in the original paper [2], and use a concrete example to provide an in-depth understanding of ARVADA.

For this walkthrough, the simple while grammar G_w and input strings S provided in figure 2.2 will be used. Clarifying again that ARVADA treats \mathcal{O} as a black box, meaning it has no structural knowledge of G_w , and G_w is only shown to clarify the behaviour of \mathcal{O} for understanding and comprehension.

G_w

```

start  → stmt
stmt   → while_ boolexp_ do_ stmt
        | if_ boolexp_ then_ stmt_ else_ stmt
        | L_ =_ numexpr
        | stmt_ ;_ stmt
boolexp → ~boolexp
        | boolexp_ &_ boolexp
        | numexpr_ ==_ numexpr
        | false
        | true
numexpr → ( _ numexpr_ +_ numexpr_ )
        | L
        | n

```

$$S = \{\text{while true \& false do L = n, L = n ; L = (n+n)}\}$$

$$O(i) = \begin{cases} \text{True} & \text{if } i \in \mathcal{L}(G_w) \\ \text{False} & \text{otherwise} \end{cases}$$

Figure 2.2: Definition a simple while grammar G_w , sample input strings S , and oracle \mathcal{O} [2]

Algorithm 1 High-level overview of ARVADA [2]

Require: a set of examples S , a language oracle \mathcal{O} .

```

1:  $bestTrees \leftarrow \text{NAIVEPARSETREES}(S)$ 
2:  $bestTrees \leftarrow \text{MERGEALLVALID}(bestTrees, \mathcal{O})$ 
3:  $updated \leftarrow \text{True}$ 
4: while  $updated$  do
5:    $updated \leftarrow \text{False}$ 
6:    $allBubbles \leftarrow \text{GETBUBBLES}(bestTrees)$ 
7:   for  $bubble$  in  $allBubbles$  do
8:      $bldTrees \leftarrow \text{APPLY}(bestTrees, bubble)$ 
9:      $accepted, mergedTs \leftarrow \text{CHECKBUBBLE}(bldTrees, \mathcal{O})$ 
10:    if  $accepted$  then
11:       $bestTrees \leftarrow mergedTs$ 
12:       $updated \leftarrow \text{True}$ 
13:    break
14:  end if
15: end for
16: end while
17:  $G \leftarrow \text{INDUCEDGRAMMAR}(bestTrees)$ 
18: Return  $G$ 

```

ARVADA follows the high-level overview provided in algorithm 1, and begins by taking all the input strings in S and building a naive flat parse tree for each input string.

In this initially naive parse tree, for each string input. There is a non-terminal t_0 , and all

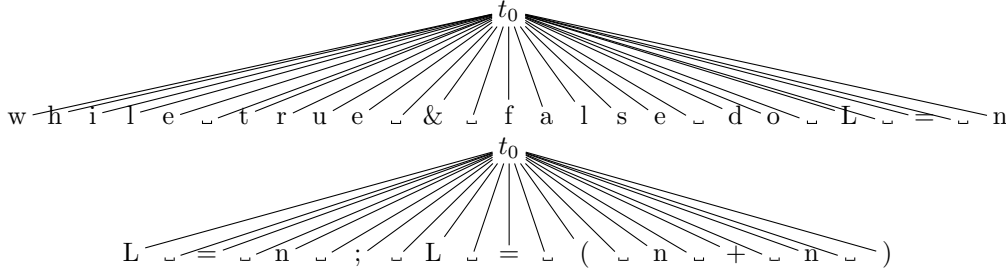


Figure 2.3: Initial set of flat naive parse trees ARVADA builds given inputs S , where each t_i is a non-terminal

characters in the string, including the spaces, are child nodes. From this alone, the following grammar can be induced.

$$t_0 \rightarrow \text{while } _ \text{true } _ \& _ \text{false } _ \text{do } _ L _ = _ n$$

$$t_0 \rightarrow L _ = _ n _ ; _ L _ = _ (_ n _ + _ n _)$$

The string derivable from a node N , in this case t_0 , is the concatenation of all its leaf nodes or itself, if N happens to be a leaf node.

Next, although not present in pseudocode 1, ARVADA has a functionality called pre-tokenisation, which can be toggled on or off. Pre-tokenisation is a step after all the naive trees are built, where the algorithm groups together sequences of contiguous characters of the same class (lowercase, uppercase, whitespaces, digits) into leaf tokens and all punctuations are kept separate. Essentially, grouping at the leaf node level.

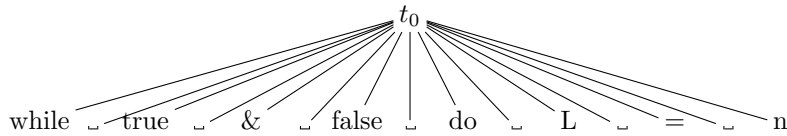


Figure 2.4: ARVADA pre-tokenisations example

Next, ARVADA attempts its first optimisation, performing MERGEALLVALID, which attempts to generalise and add structure to the tree by seeing if any of the leaf nodes/tokens can already be merged and put under new non-terminal labels. MERGEALLVALID looks at 2 leaf nodes if it is non-tokenised and/or 2 non-terminals if tokenised at a time, let these be t_a and t_b . Then it goes through all the parse trees, creating a replica of each tree with all instances of t_b replaced with t_a . Let the original parse trees be T and the replica with replacements be T' . From the replica, new candidate strings can be derived; these candidate strings are fed into \mathcal{O} . The same process is done, where all instances of t_a are replaced with t_b . Now having 2 sets of candidate strings, where in 1 set, t_b replaces t_a , and other set, t_a replaces t_b . If all the strings in both sets are accepted by \mathcal{O} , then t_a and t_b are put under a new non-terminal labelled t_c , and a successful merge has occurred. Note roots t_0 are also non-terminal and can be t_a or t_b .

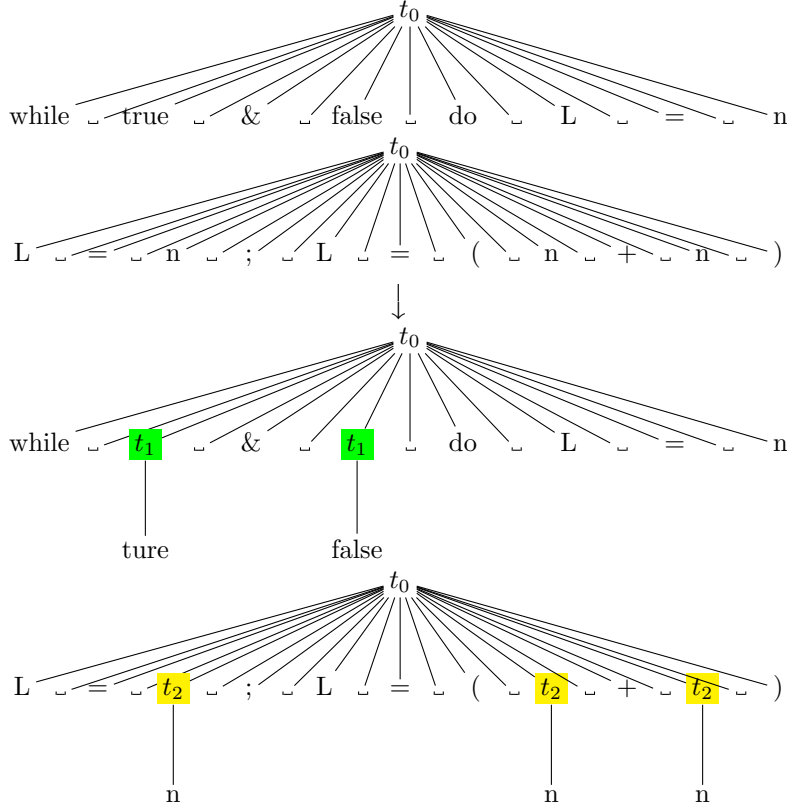


Figure 2.5: Examples of a MERGEALLVALID run, tokenised and non tokenisation

As MERGEALLVALID only looked at 1 leaf-node or non-terminal on its own (no sequence of leaf-node and/or non-terminals), ARVADA now moves onto bubbling. Where bubbling now looks at a sequence of leaf-nodes and/or non-terminal nodes.

The fundamental next step ARVADA perform is to bubble up a sequence of sibling nodes in the current trees T into a new non-terminal. To bubble up a sequence s_1 in trees T a new non-terminal is created t_{s_1} with children s_1 and all occurrences of s_1 in each tree t in T with t_{s_1} . Example shown in figure 2.6, where t_3 is a bubbled up. After bubbling up a sequence s_1 , ARVADA either accepts or rejects the bubble. A bubble is only accepted if it enables a valid generalisation of the example string inputs. Meaning if the relabeling of the bubbled non-terminal, merging its label with the label of another existing node, expands the language accepted by the induced grammar while still remaining valid to \mathcal{O} . In practice, merging here is done in a similar fashion to the merging in MERGEALLVALID. Given 2 labels t_a and t_b , ARVADA mutates trees from the current trees in T such that subtrees rooted at t_a are replaced by subtrees rooted at t_b and vice versa, which, when concatenated, gives a candidate string which can be tested against the \mathcal{O} .

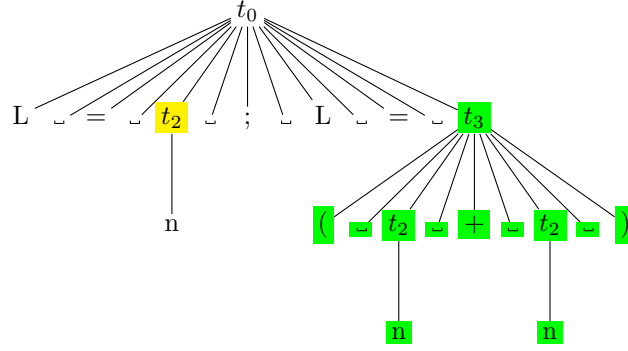


Figure 2.6: Possible bubble example

Referring to the pseudocode in 1, from the current trees T after `MERGEALLVALID`, ARVADA performs `GETBUBBLES`. For each tree $t \in T$, `GETBUBBLES` collects all proper contiguous subsequences of children in t . That is, if a tree contains a node t_i with children $C = c_1, c_2, \dots, c_n$, the potential bubble for this tree includes all sequences of length greater than 1 and less than n . Note that if the children are a subtree, the structure is maintained, like in figure 2.6. `GETBUBBLES` returns all these as 1-bubble, and all non-conflicting pairs of these as 2-bubbles. Two subsequences are non-conflicting if they do not strictly overlap: they can be disjoint or a proper subsequence of the other. So $((c_1, c_2, c_3), (c_4, c_5, c_6))$ and $((c_1, c_2, c_3), (c_2, c_3))$ are fine, however $((c_1, c_2, c_3), (c_2, c_3, c_4))$ is not. The purpose of 2 bubbles is to account for the limitations of 1-bubble. At some point, ARVADA will reach a stage where no 1-bubble can be merged with an existing label of a node in the tree, but more generalisations can still be made. Consider figure 2.8 (3), at this point in the run, ARVADA finds the single bubble $true \rightarrow t_5$ cannot be validly merged with any other existing node in the trees T . To cope with this, ARVADA bubbles up 2 distinct sequences (distinct or sub/super sets), only accepting the 2-bubble if they can merge, as seen in the example. $t_6 \rightarrow false$.

Now, ARVADA goes through 1-bubbles, taking 1 at a time and attempts to see if it can validly merge the current bubble with another non-terminal in the original trees. This is done in `CHECKBUBBLE`. If the merge is accepted, the trees are updated to reflect the successful merge; no further bubbles are checked, all current bubbles are discarded, and the process is repeated with the new updated trees. Meaning, new bubbles are formed and checked.

ARVADA to optimise checking bubbles that are more likely to get accepted first uses certain heuristics to order the exploration of the bubbles. Primarily sorting by similarity in context first, and frequency second and taking the top 100. Technical details are explained in the implementation.

After all the bubbles, 1-bubbles and 2-bubbles are exhausted, it looks at the final stage of the parse trees and returns the induced grammar. More technical explanations are provided in the implementation section of this paper.

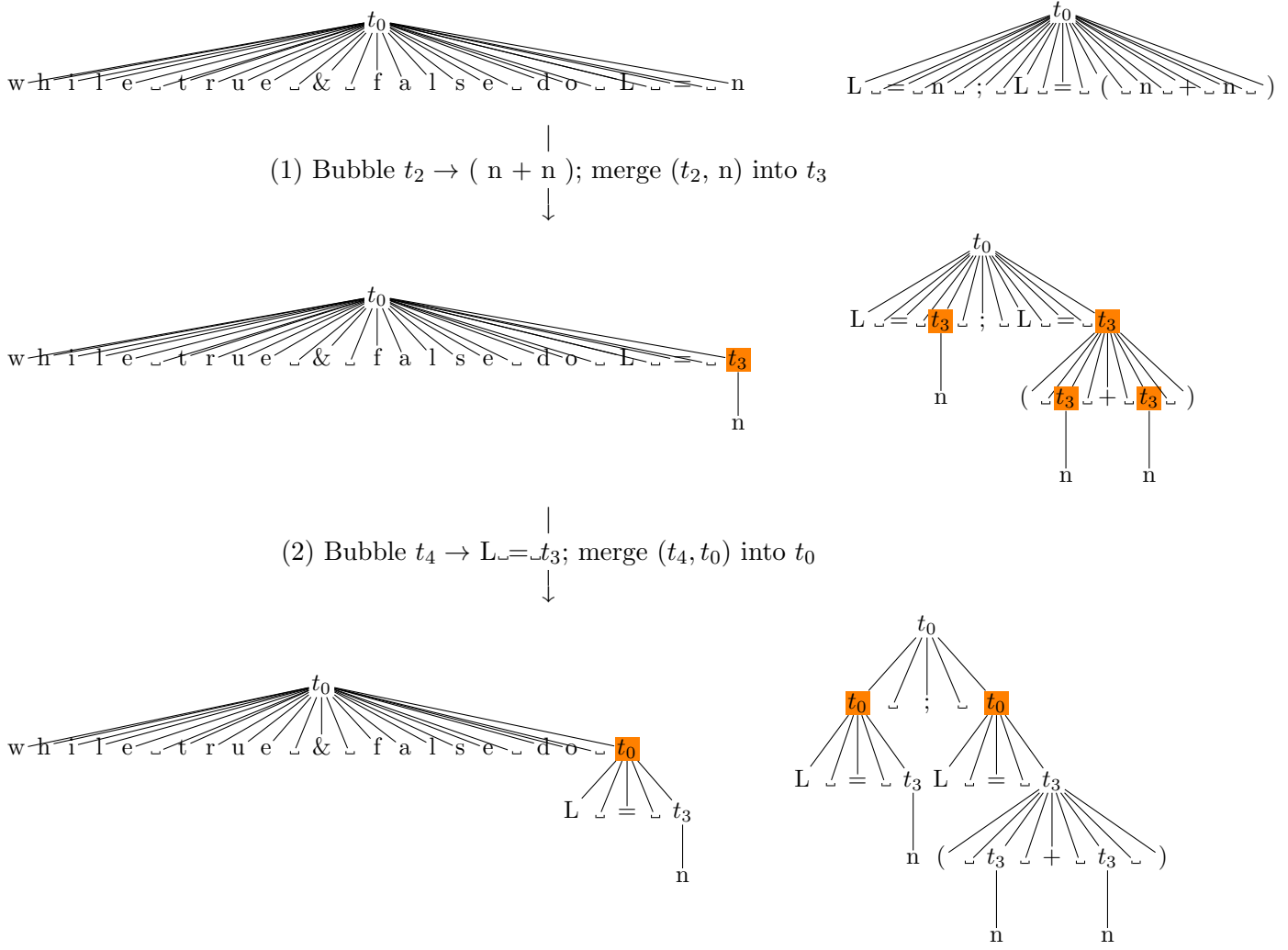


Figure 2.7: Example full run of ARVADA bubbling using non-tokenisation nor optimised tree, and the grammar induced (part 1)[2]

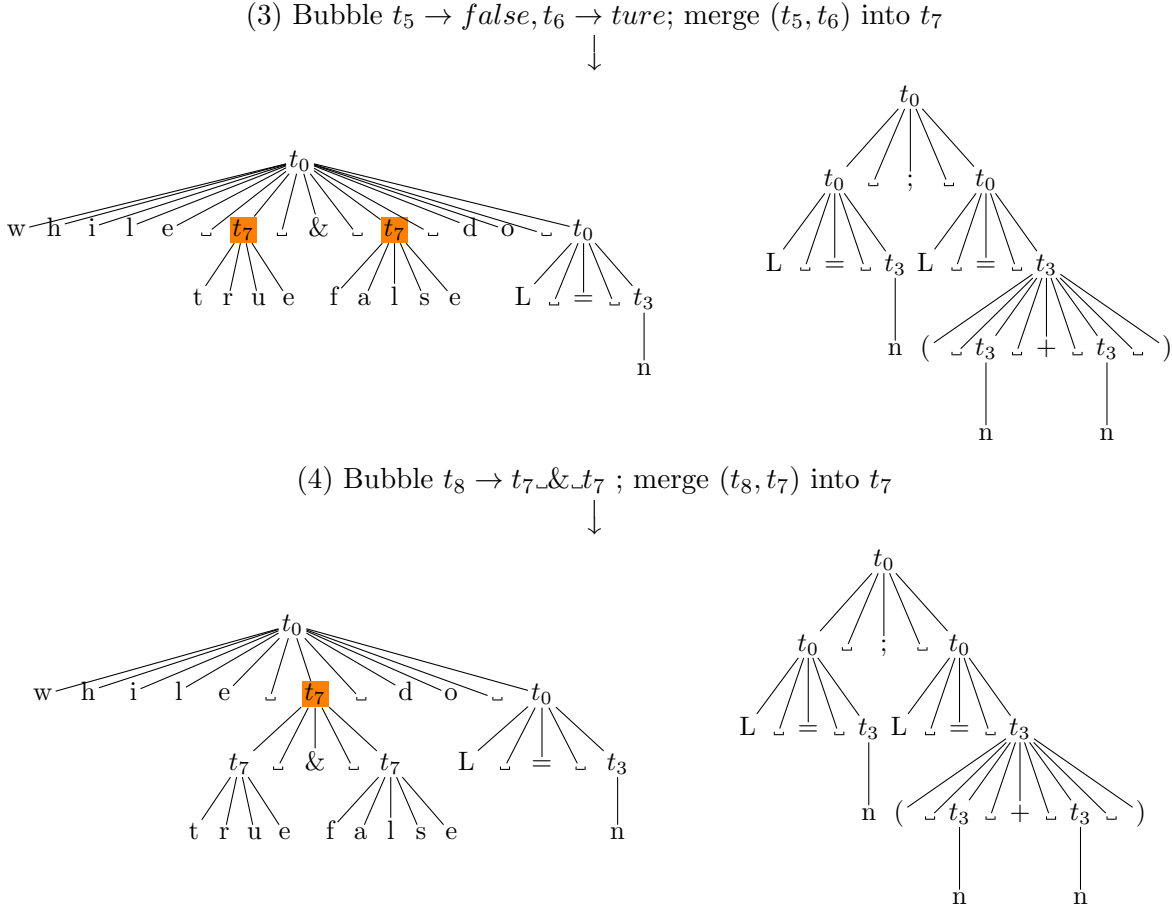


Figure 2.8: Example full run of ARVADA bubbling using non-tokenisation nor optimised tree, and the grammar induced (part 2)[2]

2.0.4 Why replicate ARVADA?

The replication study of GLADE [3] conducted by researchers at CISPA [4] in 2022 reported results that were inconsistent with those presented in the original paper [3]. Since ARVADA is closely related to GLADE, sharing a similar purpose and experimental setting, there is reason to suspect that ARVADA may also yield inconsistent results upon reimplementations.

Another motivation for replication arises from the incomplete nature of the ARVADA study and its evaluation. Although the paper presents a reasonable amount of testing and statistical data, compares the algorithm with GLADE, and provides a detailed explanation of ARVADA, it does not offer formal guarantees. It remains unclear whether ARVADA is applicable in all scenarios or what subset of scenarios it can reliably handle. Furthermore, when comparing ARVADA with GLADE, the possibility exists that the grammars used for testing were selected based on performance considerations. In later research, it has been stated that ARVADA only has high F1 scores when input strings are relatively small, and can vary widely on each run [12], raising suspicions that the study may have been selective.

2.0.5 Why C?

In the original study [2], the ARVADA algorithm was implemented in Python. When compared to GLADE [3], which was implemented in Java, ARVADA exhibited a slower average runtime across all benchmarks. In the study, this was attributed to the natural runtime disadvantage of Python compared to Java, which is valid; however, the possibility that ARVADA itself might be inherently slow was not acknowledged.

In a comparative study, A Pragmatic Comparison of Four Different Programming Languages [13], it was found that when speed and efficiency are prioritised, C is a better choice than Python. As a mid-level, statically typed, and structured language that runs under a compiler, C consistently outperforms dynamically typed, interpreted languages such as Python [14]. Moreover, C’s provision of only essential features contributes to its efficiency but also increases programming complexity compared to Python [13], [14].

Therefore, to investigate and potentially address runtime bottlenecks, C was selected as the implementation language for this replication study. The added complexity inherent to programming in C, relative to Python, was also acknowledged.

2.0.6 Why is learning input grammar important?

Grammar inference is important for many software engineering tasks, as knowledge about a program’s grammar helps with code comprehension, reverse engineering, detecting and refactoring code smells, transforming source code for optimisation or bug fixing, and generating test inputs [cite arefinFastDeterministicBlackbox2024](#). However, due to the increasing restrictions in many software systems—caused by global privacy and security regulations—we often lack access to the source code needed to learn their grammar. Even with open-source programs and code, many only have closed-source parsers, making white-box or grey-box instrumentation difficult [12], [15].

This limitation highlights the importance of learning input grammars, which are formal grammars defining valid program inputs, and approaches that treat a program as a black box and recursively apply input grammars to infer or reverse-engineer the underlying grammar of the program.

2.0.7 Related Work

2.0.8 TREEVADA

Published in “Fast Deterministic Black-box Context-free Grammar inference” 2023, at the University of Texas[12], TREEAVADA is an algorithm that is based on ARVADA[2], aiming to solve the non-deterministic and speed-related limitations of ARVADA. The paper claims, TREEVADA yields better quality grammar in a single run, with faster run time.

To achieve such results, TREEVADA uses a few different techniques. First, during pre-

tokenisation, TREEVADA pre-structures its parse according to nesting rules induced by balanced brackets common in many grammars, uses ‘ ’ quotes grouping and identifying string literals, and the same heuristics used in ARVADA, such as lowercase and numbers. Essentially, adding more structure to the initial parse trees compared to ARVADA, which does not consider brackets. Noting that TREEVADA assumes that the program uses ‘ ’ quotes to wrap a string and brackets for nesting only.

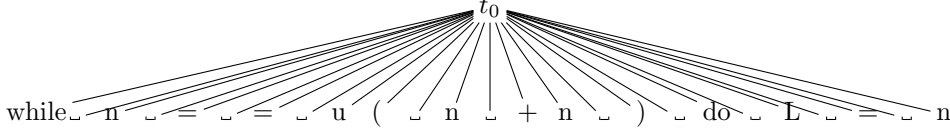


Figure 2.9: Naive Parse trees after pre-tokenisation in ARVADA

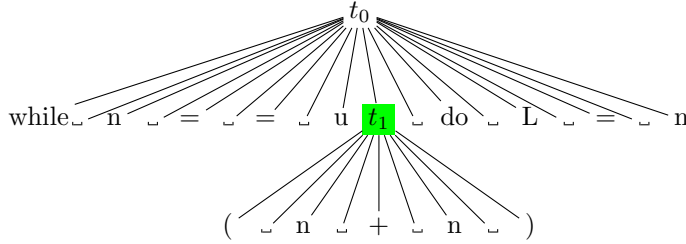


Figure 2.10: Naive Parse trees after pre-tokenisation in TREEVADA [12]

Secondly, to address the non-deterministic aspect of ARVADA, TREEAVADA switches to a deterministic data structure for specific operations and discards all bubbles with unmatched parentheses when bubbling. Additionally, two new heuristics, bubble length and bubble depth, are considered when considering bubble ranks during bubbling.

2.0.9 KEDAVRA

Proposed in “Incremental Context-free Grammar Inference in Black Box Settings” [15] in 2024, KEDAVRA is an algorithm designed to improve upon TREEVADA, in the same setting. Although TREEAVADA did improve upon ARVADA, it still had limitations of low accuracy, slow processing speeds, and limited readability due to complex grammar structures, which were inherited from ARVADA [2], [15]. The paper highlights that KEDAVRA outperforms ARVADA and TREEVADA in terms of grammar precision, recall, runtime/computational efficiency, and readability while maintaining similar memory usage.

The approach taken by KEDAVRA consists of 3 main parts:

- Tokenisation
- Data Decomposition
- Incremental grammar inference

Tokenisation

Similar to ARVADA and TREEVADA, KEDAVRA performs a tokenisation step. However, unlike ARVADA, which is tokenised based on class, and TREEAVADA, which extends it to consider brackets. KEDAVRA tokenises based on common lexical rules such as identifiers, strings, and numbers, and only considers whitespace if the given \mathcal{O} is sensitive to it.

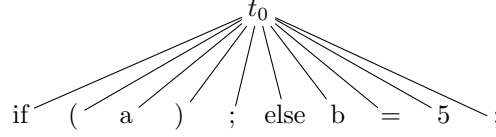
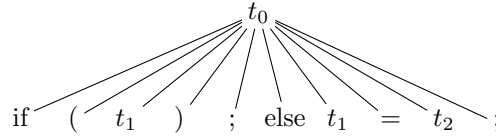


Figure 2.11: Result during KEDAVRA pre-tokenisation given input `if(a); else b = 5;`

Furthermore, KEDAVRA during grammar inference performs **character-level generalisation** for token values. Given a token character type, it computes all possible characters corresponding to that character type. For instance, if a token is a lowercase letter, all lowercase letters are included in the set.



$$t_1 \rightarrow \text{"a"} | \text{"b"} | \dots | \text{"z"}$$

Figure 2.12: Result after KEDAVRA pre-tokenisation given input `if(a); else b = 5;`, and example generalisation of t_1

Data Decompostion

KEDAVRA then breaks down all complex example strings into simple components, while collectively preserving all grammatical structure of the original sequence and essentially, breaking the input string into smaller valid input strings. This is done to overcome the slowdown and lack of readability of generated grammars due to the complexity inherent in the entire example strings.

Incremental Grammar inference

After Decomposition, KEDAVRA orders all the decomposed sequences by length and begins to infer grammar from the simplest token sequence, iterating over all token sequences derived from the input strings. Done similarly to ARVADA and TREEAVADA, KEDAVRA uses bubbling and merging. Also, like TREEVADA, only consider bubbles that have matched parentheses if parentheses are included.

```

;
t2 ;
t1 ;
t1 = t2 ;
if ( t1 ) ;
while ( t2 ) ;
if ( t1 ) ; else ;

```

Figure 2.13: Decomposed Sequences of tokenised example from 2.12

2.0.10 NATGI

Published in “Black-Box Context-free Grammar Inference for Readable & Natural Grammars” at the University of Texas 2025[16], NATGI is a novel LLM-guided grammar inference framework that extends TREEAVADA’s parse tree recovery. Aiming to improve upon the difficulties of human readability and the lack of guarantees on larger and realistic languages of previous work.

NATGI, as mentioned above, builds upon TREEAVADA’s idea of bracket-induced tree structuring. It follows the same pre-tokenising step, adding two extra modifications. First, there is no separation of class between lowercase and uppercase letters, treating them as a single class. Second, it attempts to remove redundant whitespaces by iteratively removing whitespaces and checking if the resulting program is still valid. During this process, and each time a non-terminal is created via a merge, NATGI collects the subtrees that are about to be relabelled and prompts LLM models, given the sub-trees about to be merged for a more descriptive label compared to $t_1, t_2 \dots$.

Then, before using any bubble heuristics process like ARVADA, TREEAVADA or KE-DAVRA, NATGI further exploits the brackets. It automatically treats any sequence within enclosed brackets as a bubble and attempts to further structuralise the parse trees. This partially structured tree is then given to an LLM model with a specialised prompt to finish. If the LLM terminates due to LLM’s limitations, such as hallucinations [16], [17], it then uses TREEAVADA heuristics and bubbling to compensate.

To allow a broader range of rules compared to the ones induced solely from the parse tree, NATGI decomposes each parse tree into smaller fragments. Allowing the resulting grammar to capture a broader set of language rules. This is done via effective pruning through the use of Hierarchical Delta Debugging (HDD). Pruning is the selective deletion of branches of a tree. HDD is an algorithm that iteratively prunes the trees as long as it fails a certain criterion [18].

Finally, NATGI performs lexical expansion. Since the leaf nodes are limited by the token in the input string, NATGI takes these leaf nodes and tries to expand the lexicon accepted. For

example, if leaf token is a single character, NATGI attempts to see if multiple characters are accepted.

2.0.11 Problem Statement

Due to the unknowns and gaps in the original paper [2]—including the lack of guarantees regarding the kinds of grammars ARVADA can handle and the complexity of the grammars it can learn—along with the suspicion that the original study may have been selective in nature, this thesis aims to achieve the following:

A replication of the study “Learning Highly Recursive Input Grammars” [2], through a clean-room re-implementation of ARVADA, with the intention of expanding its evaluation to include a wider variety of grammars and identifying its key characteristics.

The clean-room re-implementation will mean development is entirely from scratch in a new environment, based solely on the explanations provided in the original paper. This will also enable an assessment of the adequacy of those explanations—specifically, whether a complete re-implementation is possible based on the information given or if the paper lacks essential details for reproducibility.

Although the primary goal of this study is to re-implement, assess, and expand upon the original work, I acknowledge that this study may be incomplete or no new knowledge may be discovered. While this may seem trivial—merely confirming the findings of the previous study without adding novel insights [8], replication studies play a vital role in reinforcing the trustworthiness and confidence of empirical results, which is a central tenet of the scientific method. Replication can increase certainty when findings are reproduced and promote innovation when they are not [19]. Therefore, this study holds significance within the field of computer science.

Chapter 3

METHODOLOGIES & IMPLEMENTATION

This section of the thesis aims to provide a clear understanding of how the re-implementation was conducted in C. Providing all the steps taken in this attempt, with as much technical detail as possible and referencing appropriate sections of the original paper.

3.0.1 Data structures, and Initial Parse Trees

Data Structures

During this process of re-implementation, 2 structures were initialised. Structures are a way to group several related variables into one place in C [20]. One structure to represent a node in the parse trees, and another structure to hold and reference each parse tree.

Nodes Structure for storing trees

```
1      typedef struct nodes{
2          int capacity;
3          int count;
4          struct node **rootNodes;
5      } Nodes;
```

Figure 3.1: Data Structure to store and track trees

Figure 3.1 shows the structure *Nodes*, used to hold all the parse trees. It has variables called *rootNodes*, which is a list containing pointers to the roots of all the parse trees. Now, as this implementation is in C and memory management is a dynamic and manual process, this structure holds 2 more variables *capacity* and *count*. Both of these variables are integers that are used to help manage the size of *rootNodes*. Variable *capacity* tracks the number of root node pointers *rootNodes* can currently hold, and variable *count* tracks the number of actual root node pointers *rootNode* currently has.

Node Structure to build each tree

```

1      typedef struct node{
2          int capacity;
3          char character;
4          struct node *parent;
5          int t_label;
6          int num_child;
7          int pos;
8          struct node **children;
9      } Node;

```

Figure 3.2: Data Structure to build each tree

Figure 3.2 shows the structure *Node*, used for each individual node in the trees. It has 7 variables inside it.

- **character:** stores a C character if the node is a leaf/terminal node, else a null character if the node is a non-terminal.
- **t_label:** stores an integer. The integer is a positive value if the node is a non-terminal, else -1 if it is a terminal.
- **parent:** hold a pointer to the current node's parent, which is a *Node* as well. If root, parent is a null pointer.
- **pos:** is an integer which represents the index of the current node when the parse trees are initially constructed for all terminal nodes. For non-terminal nodes, it will be a null value.
- **children:** hold a list of pointer, which point to *Node*.
- **capacity:** current size of the list *children*.
- **num_child:** current number of *Node*(s) pointer in list *children*.

Building Initial parse trees

This part of the thesis is derived from section III-A of the original paper [2]. Using the 2 structures described, next, the initial parse trees were built. First, we build the *Nodes* structure, which houses all the parse tree pointers. This is done by using C's built-in dynamic memory function *malloc*, where *capacity* was set randomly to 4, *count* to 0, as it does not house any parse tree pointers yet, and list *rootNodes* was given enough memory to house 4 *Node* pointers, where 4 is the current *capacity*.

```

1    // Keeping track of all root trees (each string in example
    file)
2    Nodes *root_trees= malloc(sizeof(Nodes));
3    root_trees->capacity = 4; // randomly assigned
4    root_trees->count = 0;
5    root_trees->rootNodes= malloc(root_trees->capacity * sizeof
    (Node*));

```

Figure 3.3: building of Nodes Structure

Next, the file containing all the valid example strings is read as standard. Refer to figure 3.4, then from the top of the file, each valid string is read and each naive flat parse tree is built one by one. After each tree is built, the pointer to the tree's root is given and housed in `root_trees` from figure 3.3.

Memory management

As the number of valid example strings and the length of each example string are not predetermined, the number of `Nodes` (each parse tree), and the number of leaf nodes each non-terminal root node is going to is also non-determined. This means memory has to be allocated dynamically. To do this, two functions are used: one function to check and increase memory allocated to `root_nodes` in `Nodes`, and another to check and increase memory allocated to `children` in `Node`. Both these functions check the `capacity` of their respective structures, and if `count` for `Nodes` and `num_child` for `Node` equals `capacity`, they increase the memory allocated to their respective lists by 1.5 times the `capacity`, updating other variables appropriately. Function code in the appendix 6.2

Since memory is being assigned dynamically, it also has to be freed dynamically. To help with this, a help function is used. This helper function is given a node that does a depth-first search, and frees nodes as it searches, working its way bottom up. This is used to free memory as needed throughout, and all memory at the end. Refer to appendix 6.3

```

1  // Buffer to store the current example read.
2  char *current_line = NULL; //
3  size_t line_buffer_len = 0;
4  ssize_t read_line_len = 0;
5
6  //Read line by line until the end of the file has been
   reached.
7  while((read_line_len = getline(&current_line, &
   line_buffer_len, file_ptr)) != -1){
8
9      current_line[read_line_len - 1] = '\\0';
10
11     //building the naive parse tree for each string
12     // So the root node for each parse tree
13     Node *current_tree = build_basic_node(); // appendix
        6.1
14     current_tree->capacity = 10; // randomly assigned
15     current_tree->t_label = 0;
16     current_tree->children = calloc(current_tree->capacity,
        sizeof(Node*));
17
18     // Going through all the characters in the string 1 by
        1.
19     // Building the terminal nodes, assign a character
        value to the node.
20     for(int i = 0; i < (read_line_len - 1); i ++ ){
21         Node * node = build_basic_node();
22         node->parent = current_tree;
23         node->character = current_line[i];
24         node->pos = i;
25         current_tree->children[current_tree->num_child] =
            node;
26         current_tree->num_child ++;
27         // checking current trees' capacity
28         // increase appropriately if needed
29         check_node_capacity(current_tree);
30
31     }
32
33     // Check current capacity of the root node.
34     // Increase appropriately
35     root_trees->rootNodes[root_trees->count] = current_tree
        ;
36     root_trees->count = root_trees->count + 1;
37     check_nodes_capacity(root_trees);
38
39 }
40
41 // Free buffer
42 free(current_line);

```

Figure 3.4: Building each initial parse tree

3.0.2 Pre-tokenisation

```

1      // Step 2: Toggeling pre-tokenisation
2      // Section III-E of Original paper
3      int tokenise = 0;
4
5      if(tokenise){
6          for (int i = 0; i < root_trees->count; i ++){
7              pre_tokenise(root_trees->rootNodes[i]);
8          }
9      }

```

Figure 3.5: Pre-tokenisaiton of each root node

After the initial parse trees , pre-tokenisation, from section III-A of the original paper [2] was implemented. This involved looping through each parse tree in the *Nodes* structure, and for each prase tree, looping through all the leaf nodes, grouping all consective sequence of leaf nodes of the same class and putting them under a intermediate node sitting between the root node and a sequence of nodes, also known as a label node. Pre-tokenisation of each parse tree was encapsulated inside a function, which performed the pre-tokenisation given the root t_0 of a prase tree.

Given a root t_0 , the function first detaches all the children (c_0, c_1, \dots, c_n) form the root t_0 . This is just done by assigning an empy list to the root's children list pointer, and storing the list of children c_0, c_1, \dots, c_n into a new variable. It then builds an empty new node, the intermediate non-terminal label node t_i , and starts iterating over the children c_0, c_1, \dots, c_n . It takes the frist child node c_0 , assign it as a child of t_i keeping track of the class and the number of nodes assigned to t_i (lenght of sequence of same class). Then it will go through $c_1, c_2 \dots c_n$, and as long as the current child c_i is of the same class it will continue assign it as a child of t_i . When a punction or a c_i of a different class is found, if t_i houses a sequence of children greater than 1, it will then assign t_i as a child of root t_0 and a new non-terminal node t_j is built. After, same process is repeated. If t_i ever houses a sequence of lenght 1 when class changes, and the sequence of different class begin, t_i is dissolved, and child c_p , is directly linked with the root t_0 , and it is freshly restarted with a new non-terminal label node t_k , same is done for punction. Essentially, only having an intermediate non-terminal label nodes for leaf nodes in the same class and that form a consective sequence of length 2 or greater. Refer to appendix 6.4 and 6.5.

3.0.3 MergeAllValid

The first optimisation.

Chapter 4

EVALUATION

Due to the incompleteness of a full re-implementation in the given time frame, no tests have been conducted and no evaluations have been made. However, reasons for incompleteness and challenges faced during the process are discussed in the discussion section of this paper.

Chapter 5

DISCUSSION

Chapter 6

CONCLUSION

References

- [1] R. Gopinath and A. Zeller. “Building Fast Fuzzers,” pre-published.
- [2] N. Kulkarni, C. Lemieux, and K. Sen, “Learning Highly Recursive Input Grammars,” in *2021 36th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, Melbourne, Australia: IEEE, Nov. 2021, pp. 456–467, ISBN: 978-1-6654-0337-5. DOI: 10.1109/ASE51524.2021.9678879.
- [3] O. Bastani, R. Sharma, A. Aiken, and P. Liang, “Synthesizing Program Input Grammars,”
- [4] B. Bendrissou, R. Gopinath, and A. Zeller, ““Synthesizing input grammars”: A replication study,” in *Proceedings of the 43rd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, San Diego CA USA: ACM, Jun. 9, 2022, pp. 260–268, ISBN: 978-1-4503-9265-5. DOI: 10.1145/3519939.3523716.
- [5] T. Jiang, M. Li, B. Ravikumar, and K. W. Regan, “Formal Grammars and Languages,”
- [6] “Grammar in Theory of Computation,” Grammar in Theory of Computation.
- [7] S. Mulik, S. Shinde, and S. Kapase, “Comparison of Parsing Techniques For Formal Languages,”
- [8] M. Hendriks and V. Zaytsev, “Consider it Parsed!” Thesis, University of Twente, Enschede, Netherlands, 39 pp.
- [9] N. Chomsky, “THREE MODELS FOR THE DESCRIPTION OF LANGUAGE,” Paper, MIT, Massachusetts, 1956.
- [10] “Chomsky Hierarchy in Theory of Computation.”
- [11] Z. Shi, *Intelligence Science*. 2021, 215-266.
- [12] M. R. Arefin, S. Shetiya, Z. Wang, and C. Csallner. “Fast Deterministic Black-box Context-free Grammar Inference,” pre-published.
- [13] S. Ali and S. Qayyum. “A Pragmatic Comparison of Four Different Programming Languages,” pre-published.
- [14] R. Kumar, S. Chander, and M. Chahal, “Python versus C Language: A Comparison,” vol. 9, 2022.
- [15] F. Li et al. “Incremental Context-free Grammar Inference in Black Box Settings,” pre-published.
- [16] M. R. Arefin, S. Rahman, and C. Csallner. “Black-box Context-free Grammar Inference for Readable & Natural Grammars,” pre-published.

- [17] P. Orvalho and M. Kwiatkowska. “Are Large Language Models Robust in Understanding Code Against Semantics-Preserving Mutations?” Pre-published.
- [18] G. Mishherghi and Z. Su, “HDD: Hierarchical delta debugging,” in *Proceedings of the 28th International Conference on Software Engineering*, Shanghai China: ACM, May 28, 2006, pp. 142–151, ISBN: 978-1-59593-375-1. DOI: 10.1145/1134285.1134307.
- [19] M. Shepperd, “Replication studies considered harmful,” in *Proceedings of the 40th International Conference on Software Engineering: New Ideas and Emerging Results*, May 27, 2018, pp. 73–76. DOI: 10.1145/3183399.3183423.
- [20] W3School. “C Structures (structs),” C Structures (structs).
- [21] R. Gopinath, B. Mathis, and A. Zeller, “Mining input grammars from dynamic control flow,” in *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, Virtual Event USA: ACM, Nov. 8, 2020, pp. 172–183, ISBN: 978-1-4503-7043-1. DOI: 10.1145/3368089.3409679.
- [22] A. K. Lampinen. “Can language models handle recursively nested grammatical structures? A case study on comparing models and humans,” pre-published.
- [23] Open Science Collaboration, “Estimating the reproducibility of psychological science,” *Science*, vol. 349, no. 6251, aac4716, Aug. 28, 2015. DOI: 10.1126/science.aac4716.
- [24] M. Schröder, J. Cito, and T. U. Wien, “Static Inference of Regular Grammars for Ad Hoc Parsers,”

Appendix

```
1  Node *build_basic_node(){
2
3      Node * node = malloc(sizeof(Node));
4      node->parent = NULL;
5      node->capacity = 0;
6      node->character = '\\0';
7      node->t = -1;
8      node->num_child = 0;
9      node->pos = -1;
10     node->children = NULL;
11     return node;
12 }
```

Figure 6.1: Code for building a basic node

```

1
2 // Function that checks Node capacity and count.
3 void check_nodes_capacity(Nodes *nodes){
4
5     // If number for Node's is less than current capacity,
6     return.
7     if(nodes->capacity > nodes->count){
8         return;
9     }
10
11 // If number of children has reached capacity, update
12 cap.
13 int cur_cap = nodes->capacity;
14 nodes->capacity = (cur_cap + (cur_cap/2) + 1);
15 void *new_space = realloc(nodes->rootNodes, nodes->
16     capacity * sizeof(Node*));
17 if(errno == ENOMEM || new_space == NULL){
18     fprintf(stderr, "Error increasing the capacity of
19     nodes.");
20     return;
21 }
22 nodes->rootNodes = new_space;
23 }
24
25 // Function that checks Node capacity and num_child.
26 void check_node_capacity(Node *node){
27
28     if(node->capacity > node->num_child){
29         return;
30     }
31
32 // If number for Node's is less than current capacity,
33 return.
34 int cur_cap = node->capacity;
35 node->capacity = (cur_cap + (cur_cap/2) + 1);
36 void *new_space = realloc(node->children, node->capacity
37     * sizeof(Node*));
38 if(errno == ENOMEM || new_space == NULL){
39     fprintf(stderr, "Error increasing the capacity of
40     nodes.");
41     return;
42 }
43 node->children = new_space;
44 }

```

Figure 6.2: Code for changing list capacities appropriately

```
1      // Function to free all the nodes given a root Node
2      void free_tree(Node *root){
3
4          // dfs, freeing children nodes first
5          if (root->num_child > 0){
6              for( int i = 0; i < root->num_child; i++){
7                  free_tree(root->children[i]);
8              }
9              free(root->children);
10         }
11         free(root);
12     }
```

Figure 6.3: Code for freeing tree given a root node

```

1 void pre_tokenise(Node* root){
2
3     // Flags to count contiguous sequences.
4     int sequence_begun = 0;
5     int sequence_length = 0;
6     // class 1 = Whitespaces
7     // class 2 = Uppercase
8     // class 3 = Lowercase
9     // class 4 = Digits
10
11     // Set basic node with a list up.
12     // parent and children of the node assigned in progress.
13     Node *tmp = build_basic_node_with_list();
14
15     // num of child and parent save in a local var,
16     // helps later for memory management
17     int cur_children_num = root->num_child;
18     Node ** cur_children = root->children;
19     root->capacity = 1;
20     root->num_child = 0;
21     root->children = calloc(root->capacity, sizeof(Node*));
22
23
24     for( int i = 0; i < cur_children_num; i++){
25
26         Node * cur_node = cur_children[i];
27         // Current character is a whiteSpace
28         if(cur_node->character == ' '){
29             check_and_tokenise(1, &sequence_length, &sequence_begun, &tmp, cur_node);
30
31             // Current character is an uppercase character
32         } else if(isupper(cur_node->character)){
33             check_and_tokenise(2, &sequence_length, &sequence_begun, &tmp, cur_node);
34
35             // Current character is a lower character.
36         } else if(islower(cur_node->character)){
37             check_and_tokenise(3, &sequence_length, &sequence_begun, &tmp, cur_node);
38
39             // Current character is a digit.
40         } else if(isdigit(cur_node->character)){
41             check_and_tokenise(4, &sequence_length, &sequence_begun, &tmp, cur_node);
42
43             // Current character is punctuation
44         } else {
45
46             // sequence of size 1 in progress, remove warpper
47             if(sequence_length == 1){
48                 root->children[root->num_child] = tmp->children[0];
49                 root->num_child++;
50                 check_node_capacity(root);
51                 tmp->num_child = 0;
52                 tmp->children[0] = NULL;
53                 free((tmp)->children);
54                 free_tree(tmp);
55
56             } else if(sequence_length > 1){
57                 // multi character sequence, close and attach
58                 tmp->parent = tmp->children[0]->parent;
59                 (*tid)++;
60                 tmp->t_label = *tid;
61
62                 Node *p = cur_node->parent;
63                 if (p) {
64                     p->children[p->num_child] = tmp;
65                     (p->num_child)++;
66                     check_node_capacity(p);
67                 }
68
69             }
70
71             tmp = build_basic_node_with_list();
72             root->children[root->num_child] = cur_node;
73             root->num_child++;
74             check_node_capacity(root);
75             sequence_begun = 0;
76             sequence_length = 0;
77         }
78     }
79
80     // End of loop is reached, deal with the left over unclosed sequence
81     if(sequence_length == 1){
82         root->children[root->num_child] = tmp->children[0];
83         root->num_child++;
84         check_node_capacity(root);
85         tmp->num_child = 0;
86         tmp->children[0] = NULL;
87         free((tmp)->children);
88         free_tree(tmp);
89     } else if(sequence_length == 0){
90         free_tree(tmp);
91     } else {
92         root->children[root->num_child] = tmp;
93         root->num_child++;
94         tmp->parent = root;
95         check_node_capacity(root);
96     }
97     free(cur_children);
98 }

```

Figure 6.4: Pre-tokeniser code

```

1
2 void check_and_tokenise(int cur_class, int *sequence_length, int *sequence_begun, Node **tmp,
3   Node *cur_node){
4 // A sequence of the same class as currnet node is in progress.
5 if (*sequence_begun == cur_class && *sequence_length > 0){
6
7     // Because the parent of the first child is never assigned,
8     // this is done to deal with cases where a sequence of just
9     // length 1 is found.
10    if (*sequence_length == 1 && (*tmp)->num_child > 0) {
11        (*tmp)->children[0]->parent = *tmp;
12    }
13    // Assigning the current node to tmp,
14    // increase sequence number.
15    cur_node->parent = *tmp;
16    (*tmp)->children[(*)tmp->num_child] = cur_node;
17    (*tmp)->num_child ++;
18    check_node_capacity(*tmp);
19    (*sequence_length)++;
20
21 // no sequence or a different sequence has begun
22 }else{
23
24     // different sequence of length 1 is in progress.
25     if(*sequence_length == 1){
26         // previous sequence was just one node, discard wrapper
27         Node *p = cur_node->parent;
28         if (p) {
29             p->children[p->num_child] = (*tmp)->children[0];
30             (p->num_child)++;
31             check_node_capacity(p);
32         }
33
34         (*tmp)->num_child = 0;
35         (*tmp)->children[0] = NULL;
36         free((*tmp)->children);
37         free_tree(*tmp);
38         *tmp = build_basic_node_with_list();
39         (*tmp)->children[(*)tmp->num_child] = cur_node;
40         (*tmp)->num_child ++;
41         check_node_capacity(*tmp);
42
43     // different sequence of length greater than 1 has begun.
44     }else if(*sequence_length > 1){
45         // finalise previous multi-node sequence
46         (*tmp)->parent = (*tmp)->children[0]->parent;
47         *tid = *tid + 1;
48         (*tmp)->t_label = *tid;
49
50         Node *p = cur_node->parent;
51         if (p) {
52             p->children[p->num_child] = *tmp;
53             (p->num_child)++;
54             check_node_capacity(p);
55         }
56
57         *tmp = build_basic_node_with_list();
58         (*tmp)->children[(*)tmp->num_child] = cur_node;
59         (*tmp)->num_child ++;
60         check_node_capacity(*tmp);
61         *sequence_length = 1;
62
63     // no sequence in progress, assign this node as first node,
64     // and begin the sequence.
65     }else if(*sequence_begun == 0) {
66         (*tmp)->children[(*)tmp->num_child] = cur_node;
67         (*tmp)->num_child ++;
68         check_node_capacity(*tmp);
69         *sequence_length = 1;
70     }
71 }
72 }
73 *sequence_begun = cur_class;
74 }

```

Figure 6.5: Pre-tokenisation helper code