



# **“Learning Highly Recursive Input Grammars”: A Replication Study**

Sulav Malla

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

## **Supervisors:**

Rahul Gopinath, The University of Sydney

School of Electrical and Computer Engineering  
The University of Sydney

November 04, 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

# Contents

<b>Abstract</b>	<b>ii</b>
<b>1 INTRODUCTION</b>	<b>2</b>
<b>2 GENERAL REVIEW OF LITERATURE</b>	<b>4</b>
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? . . . . .	9
2.0.5 Why C? . . . . .	10
2.0.6 Why is learning input grammar important? . . . . .	10
2.0.7 Related Work . . . . .	11
2.0.8 TREEVADA . . . . .	11
2.0.9 KEDARVA . . . . .	12
2.0.10 NATGI . . . . .	12
2.0.11 Problem Statement . . . . .	12
<b>3 METHODOLOGIES &amp; IMPLEMENTATION</b>	<b>13</b>
<b>4 EVALUATION</b>	<b>14</b>
<b>5 DISCUSSION</b>	<b>15</b>
<b>6 CONCLUSION</b>	<b>16</b>
<b>References</b>	<b>17</b>



# 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 boolean-value 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 “Synthesizing Input Grammars: A Replication Study” [4] at CISPA, which investigated the accuracy of the results in the original paper. Similarly, as done by the researchers at CISPA, 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. Afterward, 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 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).

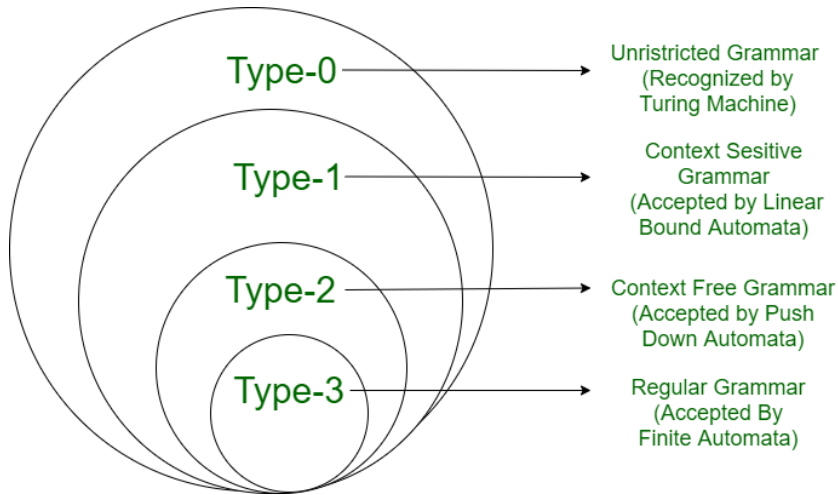
Here, the production rules  $P$  define all symbol substitutions that can be performed recursively to generate different symbol sequences, known as *strings* [6]. These rules are written in the form  $A \rightarrow w$ , where  $A \in N$  and  $w \in (N \cup T)^*$ .

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  $\varepsilon$  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, 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  $\gamma$  only in the specific context where it appears between  $\alpha$  and  $\beta$  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  $\varepsilon$ , 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  $\alpha$  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  $\mathcal{O}$ .

#### 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.

$G_w$	
$start$	$\rightarrow stmt$
$stmt$	$\rightarrow while\_ boolexpr\_ do\_ stmt$ $  if\_ boolexpr\_ then\_ stmt\_ else\_ stmt$ $  L\_ =\_ numexpr$ $  stmt\_ ;\_ stmt$
$boolexpr$	$\rightarrow \sim boolexpr$ $  boolexpr\_ \&\_ boolexpr$ $  numexpr\_ ==\_ numexpr$ $  false$ $  true$
$numexpr$	$\rightarrow (\_ numexpr\_ +\_ numexpr\_ )$ $  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]

For this walkthrough, the simple while grammar  $G_w$  and input string  $S$  provided in figure ?? will be used. Noting 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  $\mathcal{O}$ .

**Algorithm 1** High-level overview of ARVADA [2]**Require:** a set of examples  $S$ , a language oracle  $\mathcal{O}$ .

---

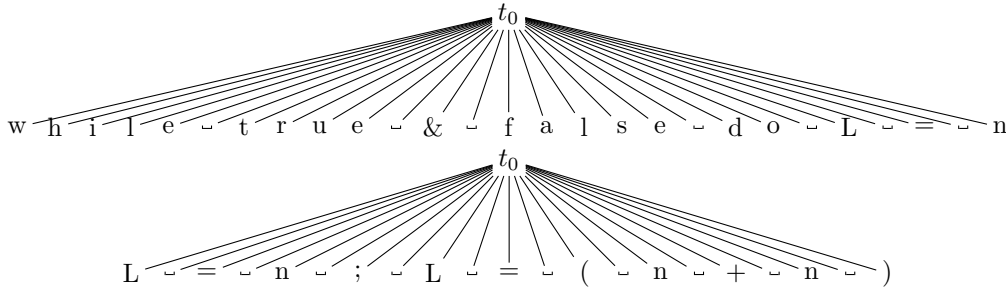
```

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.



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

In this initially naive parse tree, for each string input. There is a non-terminal  $t_0$ , and all characters in the string, including the spaces, are children nodes.

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

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

### 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 reimplementing.

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, there is a possibility 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 prioritized, 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 grammar helps with code comprehension, reverse engineering, detecting and refactoring code smells, transforming source code for optimization or bug fixing, and generating test inputs[12]. However, due to the increasing restrictions in many software systems—caused by global privacy and security regulations—we often do not have access to the source code needed to learn their grammar. Even upon open-source programs and code, many only have closed source prasers, making white-box or grey-box instrumentation difficult.

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

#### GLADE

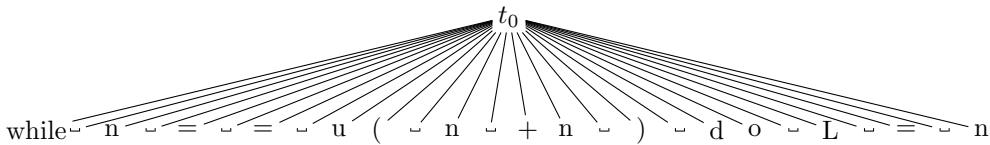
GLADE is an algorithm proposed by Bastani et al., published in Synthesizing Input Grammars at PLDI 2017 [3]. Like ARVADA, GLADE uses a set of valid inputs and black-box access to the program, with the aim of automatically approximating the context-free input grammar of the given program.

Because GLADE and ARVADA share a similar experimental setting, 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.

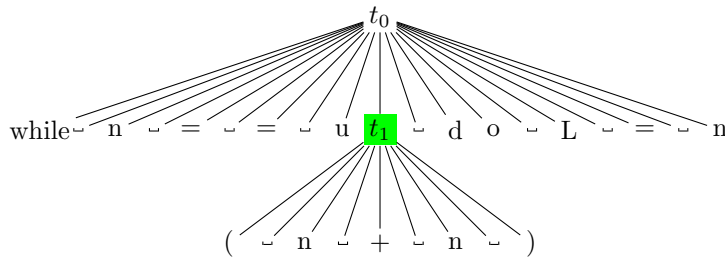
### 2.0.8 TREEVADA

Published in “Fast Deterministic Black-box Context-free Grammar inference” 2023, at the University of Texas[12], TREEVADA 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 reproducible 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 same heuristics used in ARVADA such lowercase, numbers. Essentially adding more structure to the initial parse trees compared to ARVADA, which does not consider brackets. Noting that TREEVADA assumes that program uses ‘ ’ quotes to wrap string and brackets for nesting only.



**Figure 2.4:** Naive Parse trees after pre-tokenisation in ARVADA



**Figure 2.5:** Naive Parse trees after pre-tokenisation in TREEVADA

Secondly, to tackle the non-deterministic aspect of ARVADA, by switching to determin-

istic data-structure for specific operation. Additionally, also adding two new heuristics bubble length and bubble depth when considering bubble ranks during bubbling.

### 2.0.9 KEDARVA

### 2.0.10 NATGI

### 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 explain 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 [15]. Therefore, this study holds significance within the field of computer science.



## Chapter 3

# METHODOLOGIES & IMPLEMENTATION

Methodologies: ARVADA walkthrough

How you did it, and point out any differences?

**Listing 3.1:** Struct used in code to store all *trees*

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

**Listing 3.2:** Struct used in *trees*

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

## Chapter 4

# EVALUATION

## 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] 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.

- [16] M. R. Arefin, S. Rahman, and C. Csallner. “Black-box Context-free Grammar Inference for Readable & Natural Grammars,” pre-published.
- [17] 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.
- [18] A. K. Lampinen. “Can language models handle recursively nested grammatical structures? A case study on comparing models and humans,” pre-published.
- [19] F. Li et al. “Incremental Context-free Grammar Inference in Black Box Settings,” pre-published.
- [20] Open Science Collaboration, “Estimating the reproducibility of psychological science,” *Science*, vol. 349, no. 6251, aac4716, Aug. 28, 2015. DOI: 10.1126/science.aac4716.
- [21] M. Schröder, J. Cito, and T. U. Wien, “Static Inference of Regular Grammars for Ad Hoc Parsers,”