

UNIVERSIDADE FEDERAL DE MINAS GERAIS
Instituto de Ciências Exatas
Programa de Pós-Graduação em Ciência da Computação

Saffran de Rezende, João

Reactive Methodologies to Infinite Text Processing

Belo Horizonte
2021

Saffran de Rezende, João

Reactive Methodologies to Infinite Text Processing

Final Version

Thesis presented to the Graduate Program in Computer Science of the Federal University of Minas Gerais in partial fulfillment of the requirements for the degree of Master in Computer Science.

Advisor: Haniel Barbosa

Co-Advisor: Fernando Magno Quintão Pereira

Belo Horizonte
2021

Rezende, João Saffran de.

R467m Metodologias reativas para processamento de textos
infinitos [manuscrito] / João Saffran de Rezende – 2021.
xxvi, 76 f. il.

Orientador: Haniel Moreira Barbosa.

Coorientador: Fernando Magno Quintão Pereira.

Dissertação (mestrado) - Universidade Federal de Minas
Gerais, Instituto de Ciências Exatas, Departamento de Ciência
da Computação.

Referências: f.67-72

1. Computação – Teses. 2. Linguagem de programação
(Computadores) – Teses. 3. Programação reativa – Teses. 4.
Análise (Gramática de computador) – Teses. I. Barbosa, Haniel
Moreira. II. Pereira, Fernando Magno Quintão. III. Universidade
Federal de Minas Gerais, Instituto de Ciências Exatas,
Departamento de Ciência da Computação. IV. Título.

CDU 519.6*33(043)



UNIVERSIDADE FEDERAL DE MINAS GERAIS
INSTITUTO DE CIÊNCIAS EXATAS
PROGRAMA DE PÓS-GRADUAÇÃO EM CIÊNCIA DA COMPUTAÇÃO

FOLHA DE APROVAÇÃO

METODOLOGIAS REATIVAS PARA PROCESSAMENTO DE TEXTOS
INFINITOS

JOÃO SAFFRAN DE REZENDE

Dissertação defendida e aprovada pela banca examinadora constituída pelos Senhores:

Prof. HANIEL MOREIRA BARBOSA - Orientador
Departamento de Ciência da Computação - UFMG

Prof. FERNANDO MAGNO QUINTÃO PEREIRA - Coorientador
Departamento de Ciência da Computação - UFMG

Prof. MÁRIO SÉRGIO FERREIRA ALVIM JÚNIOR
Departamento de Ciência da Computação - UFMG

Prof. RODRIGO GERALDO RIBEIRO
Departamento de Computação - UFOP

Belo Horizonte, 23 de Março de 2021.

“If I have seen further it is by standing on the shoulders of Giants.”
(Isaac Newton (1675))

Resumo

Um *evento de string* é a ocorrência de um padrão específico na saída textual de um programa. A captura e tratamento de eventos de string tem várias aplicações, como anonimização de logs, tratamento de erros e notificação de usuário, implementação de web crawlers e refatoração de código. No entanto, não há hoje uma abordagem sistemática para identificar e tratar eventos de string. Este trabalho define formalmente eventos de string e apresenta a teoria e prática de um framework para tratá-los. Demonstramos a eficácia deste framework propondo duas implementações. Primeiro, apresentamos ZHEFUSCATOR, um sistema que edita ocorrências de informações confidenciais em logs de banco de dados. ZHEFUSCATOR é implementado como uma extensão da Java Virtual Machine (JVM). Ele intercepta padrões de interesse em tempo real e não requer intervenções no código-fonte do programa a ser protegido. Demonstramos que o ZHEFUSCATOR é até 14x mais rápido do que uma abordagem força bruta, convergindo para uma gramática que descreve o formato do log de um banco de dados mysql depois de observar menos de 10 exemplos deste logs. Demonstramos também que este processo de inferir formatos de log e capturar eventos de string pode ser implementado com mínimo *overhead*. Em segundo lugar, apresentamos uma notação geral para o tratamento de texto infinito. Essa notação destaca semelhanças em tarefas que, embora em princípio diferentes, codificam os mesmos desafios essenciais. Nós combinamos essa notação propondo ZHELANG, uma linguagem reativa que permite os usuários combinarem operações básicas para identificar e tratar eventos de string. Como prova de conceito, demonstramos como os operadores de ZHELANG podem ser combinados para implementar aplicativos como: ofuscadores de log e máquinas de busca.

Palavras-chave: Linguagens de Programação, Programação Reativa, Parsing, Síntese de Gramáticas

Abstract

A *string event* is the occurrence of a specific pattern in the textual output of a program. The capture and treatment of string events has several applications, such as log anonymization, error handling and user notification, implementing web crawler and performing code refactoring. However, there is no systematic approach to identify and treat string events today. This work formally defines string events and brings forward the theory and practice of a general framework to handle them. We demonstrate the effectiveness of this framework by presenting two implementations that use it.

First we introduce ZHEFUSCATOR, a system that redacts occurrences of sensitive information in database logs. ZHEFUSCATOR is implemented as an extension to the Java Virtual Machine (JVM). It intercepts patterns of interest on-the-fly and does not require interventions in the source code of the protected program. It can infer log formats and capture string events with minimal performance overhead. As an illustration, it is up to 14x faster than an equivalent brute-force approach, converging to a definitive grammar after observing less than 10 examples from typical logs.

Second we introduce a general notation to the handling of infinite text processing. This notation highlights commonalities in tasks that, although in principle different, encode the same essential challenges. We have concretized this notation into ZHELANG, a reactive language that lets users combine basic operations to identify and treat string events. As a proof of concept, we demonstrate how ZHELANG operators can be combined to implement applications as disparate as log obfuscators and search engines.

Keywords: Programming Languages, Reactive Programming, Parsing, Grammar Synthesis

Contents

1	Introduction	9
1.1	String Events in the Context of Data Protection	9
1.2	ZHEFUSCATOR Contributions	10
1.3	ZHELANG Contributions	11
2	Background	13
2.1	Event Recognition: Challenges	13
3	ZHEFUSCATOR	16
3.1	Grammar Synthesis	16
3.1.1	Synthesizing the Grammar for the Host Language	17
3.1.2	On-Line Grammar Synthesis	17
3.1.3	Grammar Synthesis from Examples	19
3.1.3.1	Merging grammars	23
3.1.3.2	Limitations: False Positives	25
3.2	Case Study: the ZHEFUSCATOR	26
3.2.1	User Interface	27
3.2.2	Engineering	28
3.2.3	Discussion	29
3.2.3.1	Lack of Negative Examples	30
3.2.3.2	Expressiveness	30
3.3	Evaluation	30
3.3.1	RQ1—Convergence	31
3.3.1.1	Logs from Database	31
3.3.1.2	Logs from the Operating System	32
3.3.2	RQ2—Effectiveness	33
3.3.2.1	Parsing Effectiveness	34
3.3.2.2	Effectiveness on an Actual Log File	35
3.3.2.3	Increased Effectiveness via Amortized Cost	36
3.3.2.4	Impact of the Tokenizer on Runtime	37
3.3.3	RQ3—Practicality	39
3.3.3.1	Overhead of Treating one String Event	40
3.3.3.2	Deploying on Java Dacapo	41

4	ZHELANG	43
4.1	Language Specification	43
4.1.1	Core Operators	43
4.1.2	Syntax	45
4.1.2.1	Visual Representation	45
4.1.3	Semantics	46
4.1.4	Asymptotic Complexity	49
4.2	Implementation	50
4.2.1	Overview of the Implementation	51
4.2.1.1	The Implementation of Observers	51
4.2.1.2	The ZHSTREAM	53
4.2.1.3	The ZHEPANEL	54
4.2.2	Core Property: Isolation	54
4.3	Use Cases	55
4.3.1	Log Obfuscation	56
4.3.2	Web Crawling	59
4.3.3	Program Refactoring	61
4.3.4	Analysis of the Runtime Behavior of ZHELANG	64
5	Literature Review	66
5.1	Parsing	66
5.2	Inductive Grammar Synthesis	67
5.3	Program Fuzzing	68
5.4	Interactive Grammar Inference	69
5.5	String Events	69
6	Conclusion	70
	Appendix A Proofs of Lemmas and Theorems	72

Chapter 1

Introduction

We define a *string event* as the occurrence of some pattern of interest in the output of a program. Events can be produced automatically, for instance, as part of a log, or due to interactions between programs and users, as in a chat system. Examples of events of interest include the output of sensitive information that must be redacted or occurrences of notifications requiring immediate attention. Since there is no unified framework for capturing and treating string events, each software application handles them in specific ways. Nevertheless, the building blocks to construct such infrastructure are already in place: grammar synthesis [?, ?, ?], function interception [?, ?], parsing [?, ?] and reactive programming [?]. This work uses this body of knowledge to propose two approaches that handle string events. First we define ZHEFUSCATOR, a framework that handles string events for applications running in the Java Virtual Machine, as a way to anonymize sensitive data in logs. Second, we propose a domain specific language, called ZHELANG, to handle string events on infinite text streams.

1.1 String Events in the Context of Data Protection

In this work, we call a *generator* a computer program that produces a string t_i at each time slot i . Software that produces logs, like database servers and operating systems, or content providers, such as e-mail and news services, can be understood as generators. Usually, when part of the output of a generator is analyzed, this analysis is performed *off-line*, i.e., after such text has been produced and stored. However, there are situations in which such analysis must be carried out *on-line*, i.e., while it is being produced.

Data protection laws are one of the forces driving the need for on-line analyses. As an example, the *General Data Protection Regulation* (GDPR)¹, valid in the European Economic Area since 2016, requires companies to anonymize personal data, whenever this data is amenable to be used in ways not foreseen by the company's terms of use [?].

¹<https://eugdpr.org/>

Discussions involving the European GDPR have inspired similar laws in other regions, such as the *California Consumer Privacy Act*², effective since January of 2020 in the American state of California, and the *General Law of Personal Data Protection* [?], taking effect since August of 2020 in Brazil.

Data protection laws bear an impact on log generation, since logs should not leak personal data. However, many software systems have been designed and implemented before the advent of these laws. Adapting these systems to accommodate privacy is an expensive endeavor inasmuch as such adaptation entails modifications in legacy code. However, in this work, we demonstrate that it is possible to filter logs while they are produced, by projecting this problem onto the general framework of string events. The appearance of sensitive information in a log is a string event. Given the right framework, this event can be detected and treated on-the-fly. Nevertheless, the creation and deployment of this framework involves theoretical and practical challenges, which we discuss in the next sections.

1.2 ZHEFUSCATOR Contributions

We describe in Chapter 3 an on-line grammar synthesis algorithm that incrementally over-approximates a grammar for any language (Section 3.1). Our grammars fit into a format henceforth called *Heap-Chomsky Normal Form* (Section 3.1.1), a restriction of Chomsky Normal Form. Heap-CNF grammars recognize a regular language; hence, they can be represented as regular automata. Therefore, these grammars are never ambiguous and admit LL(1) parsers. LL(1) parsers can run in linear time on the input size [?], and admit formal proofs of correctness, as recently shown by Edelmann *et al* [?]. We have implemented a system that uses our theory to anonymize sensitive information in logs while treating the log generator as a black-box (Section 3.2).

Summary of Results. We implemented the above techniques in a tool, the ZHEFUSCATOR, that redacts sensitive data in SQL queries found in logs created by Java-based systems. ZHEFUSCATOR implements a form of reactive programming, which, in the words of Ramson and Hirschfeld [?, p12-2], “*consist of two parts: detection of change and reaction to change.*” Detection is the topic of Section 3.1, whereas reaction is discussed in Section 3.2. In Section 3.3 we evaluate properties of this tool. We summarize the results of this evaluation as follows:

- Section 3.3.1 shows that we can construct a grammar for typical database logs

²AB-375 Privacy: personal information: businesses.(2017-2018)

(MySQL and PostgreSQL) after observing less than 10 examples of outputs. Exercising ZHEFUSCATOR on more complex logs, e.g., files in the `/var/log` directory of MacOS, then convergence requires more examples, but still a small proportion compared to the size of the log. Our worst case performance required 170 examples in a log containing 6,579 entries.

- Section 3.3.2 demonstrates that our on-line approach can be up to 14x faster than a brute-force event detection system that does not synthesize grammars. Performance is important because our techniques are meant to be used in tandem with a running application. If its overhead is prohibitive, then chances are that users would not employ it. Furthermore, the more complex is the language that generates the logs, the larger is the improvement of ZHEFUSCATOR over its trivial counterpart.
- Section 3.3.3 shows that our event handler does not add statistically significant overhead onto 11 out of 15 benchmarks from DaCapo [?] when building a grammar for the entire output of each benchmark. Furthermore, in the four benchmarks where overhead is noticeable, in only one case (`luindex`) it reaches 50%.

Software ZHEFUSCATOR is open software, distributed through the GPLv3 license, and publicly available at <https://github.com/lac-dcc/Zhe>. As of today, it is embedded in products of at least one data-protection company: Cyral Inc. (<https://www.cyral.com/>). And has been published at the Journal of Computer Languages [?]

1.3 ZHELANG Contributions

We propose a new methodology to process infinite text streams, which will be further discussed in Chapter 4. This methodology is centered around a minimal set of operators. These operators—six in total—define a reactive domain specific language, henceforth called ZHELANG. In our vision, the reactive paradigm is an essential metaphor to process infinite text. In the words of Ramson and Hirschfeld [?, p12-2], Reactive Programming “*consists of two parts: detection of change and reaction to change.*” In our context, events are triggered by the occurrence of particular patterns in the text stream. Detection is, thus, the act of recognizing such patterns. Reaction, in turn, is the invocation of predetermined actions over said patterns.

ZHELANG is a scripting language. Programs written as combinations of its six operators must be linked against an *action language*. This language will let programmers specify the right reaction to the occurrence of each event. In this thesis, we evaluate an

implementation of ZHELANG that runs on the Java Virtual Machine and hosts actions written in the Kotlin programming language (<https://kotlinlang.org/>). Additionally, we have an implementation of ZHELANG in Haskell, which we use to specify its semantics. The theory behind ZHELANG, together with its different implementations, yield a number of contributions, which we summarize as follows:

- **Language:** Section 4.1 describes the syntax and the semantics of ZHELANG. ZHELANG is a reactive language, formed by an ensemble of operators. Each operator specifies a particular text pattern, and the action triggered, once such pattern is recognized in the text. Operators can be chained in a way that permits the emergence of complex text processing applications out of simple rules.
- **Implementation:** Section 4.2 describes the implementation of ZHELANG. Currently, we have two prototypes of it: one in Haskell, the other in Kotlin. The former is used for reference purposes: it encodes literally the semantics introduced in Section 4.1.3. The latter is a more mature implementation that supports the craft of high-performance applications.
- **Applications:** Section 4.3 describes three use cases of ZHELANG. In Section 4.3.1, we show how this language implements a system to obfuscate sensitive information in logs; in Section 4.3.2 we use it to implement a web crawler; and in Section 4.3.3 we show how ZHELANG supports the creation of static program analyses. In each application we compare ZHELANG's implementation with a state-of-the-art tool for the respective application.

Chapter 2

Background

In Section 2.1 it is discussed the challenges related to the automatic treatment of string events.

2.1 Event Recognition: Challenges

Handling string events while treating the event generator as a black box is challenging for three reasons, which we discuss in this section. To make this presentation more concrete, we relate the challenges to the following real-world problem:

Example 1 (Concrete Problem) *Consider a log-producing database server running on the Java Virtual Machine. The grammar that describes the log syntax is unknown. Logs might contain SQL queries. Some queries contain sensitive information. Design a system that intercepts strings in the log, before they are printed, and anonymizes particular literals embedded in the SQL queries. A literal is any constant in the SQL query, e.g., integer values, quoted strings, and so on. The users specifying which data must be elided are not necessarily programmers.*

Challenge 1 (Grammar Synthesis) *How to efficiently identify SQL queries within the log, when the log grammar is not known?*

Each generator has its own log format. Part of this log uses the SQL syntax. If we call L the language of log strings, then each string $t \in L$ might contain SQL and noSQL substrings, as Example 2 shows. In this combination of two languages, we call L the *host language* and SQL the *event language*.

Example 2 *Figure 2.1 shows part of a log taken from an actual application (literals have been replaced with fake surrogates). Strings in the target language, SQL, are shown in red. This log contains five examples, one per line. Each example is produced by the generator in*

```

82 Query SELECT * FROM Clts WHERE SSN='078-05-1120' 0
83 Init DB grossi
11 8:02 84 Query SELECT * FROM Byrs WHERE name='J.Generics' 1
85 Connect mysqldumpuser@localhost on
12 8:11 86 Query DELETE * FROM Clts WHERE name='J.Generics'

```

Figure 2.1: Snippet of log with five examples.

successive moments in time. A solution to Challenge 1 amounts to synthesizing a parser for this log.

Requiring a parser for the host language L would complicate the deployment of the obfuscator, as this requirement forces users to be aware of L 's format. It is possible to separate host and event languages via a brute-force approach considering every token of the host language as the potential starting point of a sentence in the event language. However, as we show in Section 3.3.2.1, this approach does not scale well with the number of tokens in the string $t \in L$. The generator produces an infinite stream of strings; hence, Challenge 1 involves inferring a grammar *in the limit*, that is, from an infinite number of examples. Even though this problem is undecidable even for regular or superfinite languages, as shown by Edward Gold [?], ZHEFUSCATOR can efficiently build unambiguous grammars that recognize, in a scalable manner, the subset of the host language defined by all examples seen up to a point. We detail this process in Section 3.1.

Challenge 2 (Interface) *Which interface should users who are not programmers use to specify sensitive patterns?*

Obfuscating the log in Figure 2.1 requires knowing which SQL literals must be redacted. It is up to users of the obfuscator to specify such literals. However, information can be sensitive when used in some types of queries, and innocuous when used in others, as Example 3 illustrates.

Example 3 *Consider an instance of the concrete problem (Ex. 1) that requires redacting occurrences of SSN in the pattern: `SELECT * FROM Clts WHERE SSN='??'`. Occurrences of SSN in other patterns, such as `DELETE FROM Clts WHERE SSN='000-00-0000'`, must be preserved.*

One problem of defining a domain specific language (DSL) to let users specify patterns to obfuscate, such as ZHELANG, is that it prevents users of the log-producing system—usually non-programmers—from using our tool. In Section 3.2.1 we describe a programming-by-examples approach, inspired by the Parsimony IDE [?], which provides users a simple but effective interface for specifying sensitive data. From this interface, we derive an *event grammar* that specifies which queries should have their literals redacted. This grammar feeds ZHEFUSCATOR with knowledge to distinguish sensitive and innocuous

queries. It will redact every literal within the former group, while preserving occurrences of the same literal in the latter. What distinguishes one type of query from the other? Syntax! And this syntax is specified by the user, when building the event grammar (following steps yet to be introduced in Section 3.2.1). Notice that the user will never have to deal with the format of the tokens, e.g., the SSN format in Example 3. All that she must do is to highlight examples of sensitive queries.

Challenge 3 (Engineering) *How to intercept the generator's output without changing its implementation?*

Challenge 3 is an engineering problem specific to the log-generation application. In Section 3.2.2 we describe ZHEFUSCATOR solution for systems running on the Java Virtual Machine. In contrast to our solutions to the other challenges, the approach adopted in Section 3.2.2 is not general—a natural consequence of the fact that Challenge 3 is technology specific. In Section 4.2.1 we will describe ZHELING approach to handle user defined actions, this makes Challenge 3 simpler to solve.

Chapter 3

ZHEFUSCATOR

3.1 Grammar Synthesis

Context-free grammars. Let $G = \langle S, N, T, P \rangle$ be a *context-free grammar*, with *non-terminals* N , *terminals* T , a *start symbol* $S \in N$ and *production rules* $P \subseteq N \times (N \cup T)^*$. The set $V = N \cup T$ is G 's *vocabulary*. A *sentence* is a string of terminals. A sentence t is *generated* from a grammar G if there is a sequence of applications of production rules that transforms S in t . This sequence of applications is called a *derivation*. In a *leftmost derivation* the leftmost non-terminal is always reduced first. The concatenation of strings p and q is $p \bullet q$. If t and t' are strings, and t is a substring of t' , we write $t \in \text{subs}(t')$. A context-free grammar G is in *Chomsky normal form* if all of its production rules are of the form $A ::= BC$, $A ::= \mathbf{a}$, or $S ::= \epsilon$, where ϵ is the empty string, in which A , B and C are non-terminals, \mathbf{a} is a terminal and S is the start symbol. The language that G recognizes, denoted $\text{lang}(G)$, is the set of all strings generated from G . Given a string t , it can be *generated ambiguously* by a grammar G if G allows two different derivations that generate t . If G generates any string ambiguously, then G is *ambiguous*.

String events. Let L be a language. A *text* over L is a sequence of strings t_0, t_1, \dots , such that $t_i \in L$. A *generator* for L is a Turing Machine that generates this text. We say that t_i is the *text generated at time* i . We allow $t_i = t_j, i \neq j$. No function from time to strings is assumed; however, we assume that on the limit the text covers L . Notice that the existence of a generator, coupled with this last assumption, implies that L is recursively enumerable. From these notions, we define string events as follows:

Definition 1 (String Event) A *string event* $\langle s, G_e, t_i, L \rangle$, parameterized by a context-free grammar G_e , which we call the *event grammar*, occurs at time $i, i > 0$, on the text t_i produced by a language L , which we call the *host language*, if there exists $s \in \text{lang}(G_e)$, such that $s \in \text{subs}(t_i)$.

Example 4 (String Event) Let the host language L be the language that contains the string representations of every natural number, and only these strings. Let the event

grammar G_e be a grammar that recognizes palindromes with more than one digit on the language of positive decimal numbers. Tokens, in this case, are single digits. Consider the text over L in which $t_i = "i"$, for $i \in \mathbb{N}^+$, i.e., the text is "1", "2", ..., "10", "11", A string event occurs on $t_{1223} = "1223"$, because "22" \in subs("1223") is a palindrome with more than one digit.

3.1.1 Synthesizing the Grammar for the Host Language

As seen in Definition 1, capturing string events involves detecting occurrences of substrings produced by a context-free grammar G_e within text pertaining to a recursively enumerable language L . We call a grammar G that recognizes L , i.e., $L = \text{lang}(G)$, the *host grammar*. In the context of handling string events from a black-box event generator, as explained in Section 2.1, we cannot assume that the host grammar is known. Thus, it is necessary for L to be discovered while string events are being captured. Moreover, only examples of strings that are part of the language, denoted "positive examples", are available to do so. As demonstrated by Gold [?], this problem is undecidable for most classes of languages, including context-free.

3.1.2 On-Line Grammar Synthesis

The intuition behind Gold's result is simple: since L is being determined by positive examples, whichever grammar has been synthesized up to time m can fail to parse an example $t_n, n > m$. However, up to time m , it is always possible to build a grammar G_m that recognizes t_1, \dots, t_m : in the worst case, G_m contains m production rules, one for each string $t_i, 1 \leq i \leq m$. Therefore, Gold's conclusions indicate that a grammar for L should be recognized by an *on-line* algorithm, which builds successive grammars G_1, \dots, G_m up to time m , such that $\{t_1, \dots, t_m\} \subseteq \text{lang}(G_m), 1 \leq i \leq m$.

The Language Separation Problem. In this work, we assume that the event grammar G_e that encodes string events is known, ZHEFUSCATOR requires it as a parameter¹ and ZHELANG provides operators to define it. Therefore, to capture string events we

¹Section 3.2.1 discusses the approach that we have chosen to let users specify events. Notice that users do not need to provide G_e explicitly: they specify events through examples valid in G_e , which is assumed to be already known by the language synthesis system.

```

1 # An infinite sequence of strings:
2 val text: String stream
3
4 # Parameters of the implementation
5 val TOKENIZE: String -> Token list
6
7 fun add_example
  (tokens: Token list, current_grmr: Grammar) =
8   if successfull_parse(current_grmr, tokens)
9   then current_grmr # Success!
10  else
11    let
12      val new_grammar = fill_holes(tokens)
13    in
14      merge(current_grmr, new_grammar)
15    end
16
17 fun build_grammar((example::text): String stream, grammar: Grammar) =
18   let
19     val new_grammar = add_example(TOKENIZE example, grammar)
20   in
21     build_grammar(text, new_grammar)
22   end
23
24 # Start language separation with the simplest sketch grammar:
25 val grammar:Grammar = build_grammar(text, R1 ::= ε)

```

Figure 3.1: The language separation procedure.

must be able to distinguish occurrences of strings from $\text{lang}(G_e)$ within the input text. From these observations, we define the *language separation* problem as follows:

Definition 2 (Language Separation Problem) *Let $T = \{t_1, \dots, t_m\}$ be a set of strings pertaining to an unknown host language L . Given $G_e = \langle S_e, N_e, T_e, P_e \rangle$, find grammars $G_m = \langle S_m, N_m \cup N_e, T_m \cup T_e, P_m \cup P_e \rangle$, such that $\{t_1, \dots, t_m\} \subseteq \text{lang}(G_m)$, $1 \leq i \leq m$.*

ZHEFUSCATOR language separation algorithm is outlined in Figure 3.1 as a program written in ML syntax. The entry point of this program is function `build_grammar`, which receives `text`, the infinite sequence of strings t_1, t_2, t_3, \dots corresponding to the language to be recognized. The function `build_grammar` operates in a classic *counterexample-guided inductive synthesis* (CEGIS) [?, ?] loop, in which a learner proposes solutions and a verifier checks them, providing counterexamples for failures. In our context the learner produces grammars that recognize the examples seen so far and the verifier checks whether they can generate the subsequent examples.

For each string `example` in the `text` stream, `build_grammar` refines a `grammar` that recognizes `example`. Thus, the `grammar` variable at line 25 of Figure 3.1 refers to the grammar that recognizes `text` on the limit, that is, after an infinite number of examples have been produced. Notice, nevertheless, that even though `build_grammar` never halts, it produces a new grammar each time it is recursively invoked (Line 21). Function `build_grammar` uses an auxiliary routine `add_example`. This procedure checks if the current grammar can parse a string in `text` (Line 8). If it can, nothing else happens (Line 9). However, if parsing fails, then `add_example` refines the current grammar (Lines 10-15). The next section describes this refinement.

On the TOKENIZE Function ZHEFUSCATOR does not focus on synthesis of lexers. Instead, it relies on a predefined lexer, the `TOKENIZE` function, which transforms examples into sequences of tokens. Said function is invoked at Line 19 of Figure 3.1. Examples of

tokens are $int = \{\dots, -2, -1, 0, 1, 2, \dots\}$ and $time = int : int$. Our solution to language separation (Fig. 3.1) is parameterized by this function. The tokenizer might bear an impact on the number of examples necessary to synthesize a definitive grammar for the host language. It can also modify the speed of the algorithms that we shall discuss in the next section. In Section 3.3.2.4 we analyze these two facts empirically.

3.1.3 Grammar Synthesis from Examples

Whenever `build_grammar` fails for a new example t_i , ZHEFUSCATOR uses the function `fill_holes` to produce a grammar G_i that recognizes it. This function is invoked at Line 12 of Figure 3.1, and its implementation is given in Figure 3.2. We shall be explaining this code in the rest of this section. Notice that the auxiliary function `build_hcnf` contains comments mentioning two “Rules”. These rules will be explained shortly.

```

1 # Build a grammar in Heap-CNF that recognizes “tokens”
2 fun build_hcnf(n:int, [token]: Token list): Grammar =
3   Rn ::= token      # Rule 1
4   | build_hcnf(n:int, token::Rest: Token list): Grammar =
5     Rn ::= R2nR2n+1 # Rule 2
6     R2n ::= token    # Rule 1
7     build_hcnf(2×n+1, Rest) # R2n+1 ::= ...
8
9 fun fill_holes(tokens: Token list): Grammar = build_hcnf(1, tokens)

```

Figure 3.2: The grammar synthesizer.

To build a parser for the host language L , thus solving the Language Separation Problem, ZHEFUSCATOR applies a programming-by-examples [?] approach. For each example t_i ZHEFUSCATOR synthesizes a grammar G_i that generates it. Then we merge this grammar into a previously synthesized grammar G that generates the previous examples, thus obtaining a new grammar G such that $\{t_1, \dots, t_i\} \subseteq G$. Each grammar G_i synthesized for generating the given example t_i is in Heap-CNF, a restrictive form of CNF defined as follows:

Definition 3 (Heap-CNF) *A Heap-CNF grammar has restrictions on the non-terminals and the production rules. Non-terminals are $R_1, R_2, R_3, \dots, R_{2^n-2}, R_{2^n-1}$, for some arbitrary n . The allowed production rules are*

- $R_{2^{k+1}-2} ::= a,$
- $R_{2^k-1} ::= R_{2^{k+1}-2}R_{2^{k+1}-1},$ and

- $R_{2^k-1} ::= \mathbf{a}$

in which \mathbf{a} is a terminal and $k \in \{1, \dots, n\}$. Since non-terminals are numbered in the same way as data in the heap data structure, we call this restricted version of Chomsky Normal Form, *Heap-CNF*.

ZHEFUSCATOR restrict itself to Heap-CNF grammars for three reasons. First, given two grammars in this format, it is possible to merge them in linear time on the number of non-terminals, thus producing a new Heap-CNF grammar, as we will see in Section 3.1.3.1. Second, they are not ambiguous (Theorem 5). Finally, they admit LL(1) parsing (Theorem 7). We shall leverage the two latter properties to demonstrate that ZHEFUSCATOR’s solution to the language separation problem is correct. The two latter properties are a consequence of Heap-CNF grammars encoding regular languages. Indeed, a Heap-CNF language can be described by a regular automaton. Nevertheless, we shall call them grammars, as ZHEFUSCATOR is using them to synthesize parsers.

The grammar G_i is built by successively increasing its vocabulary and by “filling holes”, i.e. adding production rules to a partial grammar while $t_i = t_i^1 t_i^2 \dots t_i^n$ is traversed. Initially the partial grammar contains only the starting non-terminal R_1 and no terminals or production rules. At each iteration, the grammar is augmented to generate the first token in the sequence, which is then removed from it. The grammar is also expanded so that it can be further augmented to generate the remaining tokens. This is represented by the application of the following two expansions, which add production rules to G_i :

$$\begin{aligned}
 R_k ::= ? & \stackrel{\text{(Rule 1)}}{\Rightarrow} R_k ::= t_i^j, \text{ if } t_i^j \text{ is the last token of } t_i \\
 & \text{ or } G_i \text{ contains } R_{k+1} \\
 R_k ::= ? & \stackrel{\text{(Rule 2)}}{\Rightarrow} R_k ::= R_{2k} R_{2k+1}, R_{2k} ::= ?, R_{2k+1} ::= ? \\
 & \text{ otherwise}
 \end{aligned}$$

in which R_k is a non-terminal not yet associated with a production rule.

The first rule allows the consumption of the first remaining token in t_i . It can be applied when the respective non-terminal is not the last one in G , except if there is only one token left to be consumed. Otherwise, the second rule is applied, which introduces two new non-terminals in the grammar: one for consuming the first remaining token, via Rule 1, and another to continue the process for generating the subsequent tokens in t_i . This process continues until the grammar that parses t_i is obtained. Function `fill_holes` (Figure 3.2), which implements this procedure, takes as input a sequence of tokens t_i and yields a grammar G_i in Heap-CNF that consumes said sequence, as stated below.

Theorem 1 (Correctness) *Function `fill_holes` (Fig. 3.2) produces a grammar G_i that recognizes an example $t_i = t_i^1 \dots t_i^n$ in n steps with $2n - 1$ non-terminals.*

Proof Sketch 1 *The proof is by induction on the size $|t_i|$ of the example.*

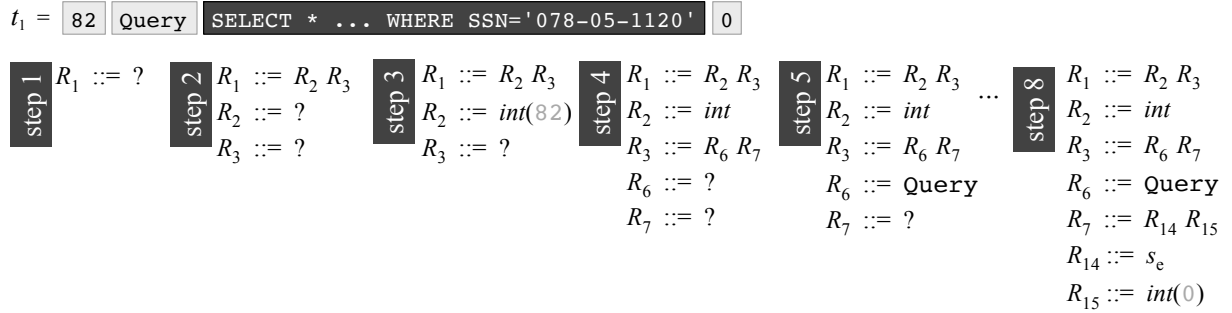


Figure 3.3: Grammar inference via `fill_holes`.

Example 5 *Figure 3.3 illustrates `fill_holes` “stepwisely” building a grammar that generates the first example from Figure 2.1. Note that the characters ‘82’ and ‘0’ were tokenized to `int` and the SQL query to `se`. At Step 1 the partial grammar consists only of the starting non-terminal R_1 . Given that the list of tokens to generate contains more than one element, `fill_holes` applies rule 2 (line 5 of Figure 2.1), producing the partial grammar in Step 2 with two new undefined non-terminals. Rule 1 is then applied to generate the first token in the list (line 6), producing the grammar in Step 3. The `fill_holes` algorithm proceeds to recursively build a grammar to generate the remaining tokens, applying rules 2 and 1 in sequence, until it reaches the case when there is only one token to be generated. This triggers a final application of rule 1 (line 3), yielding the grammar in Step 8.*

The grammar synthesis has the following properties:

Lemma 1 (fill_holes yields Heap-CNF) *Given an example t_i , the resulting grammar G_i produced by `fill_holes` that generates t_i is in Heap-CNF.*

Theorem 2 *Given an example $t_i = t_i^1 t_i^2 \dots t_i^n$, the resulting grammar G_i produced by `fill_holes` is such that $R_{2^{n-1}} ::= t_i^n$.*

Theorem 2 and Lemma 1 perfectly define the structure of grammars produced by `fill_holes`, as stated below:

Corollary 1 *Given an example $t_i = t_i^1 t_i^2 \dots t_i^n$, the resulting grammar G_i produced by `fill_holes` is such that*

$$R_{2^{k+1}-2} ::= t_i^k, k \in \{1, \dots, n-1\}$$

$$R_{2^k-2} ::= \begin{cases} R_{2^{k+1}-2} R_{2^{k+1}-1} & k \in \{1, \dots, n-1\} \\ t_i^n & k = n \end{cases}$$

Figure 3.4 illustrates the structure of a derivation of a string of 5 tokens from the Heap-CNF grammar that would be produced by `fill_holes`.

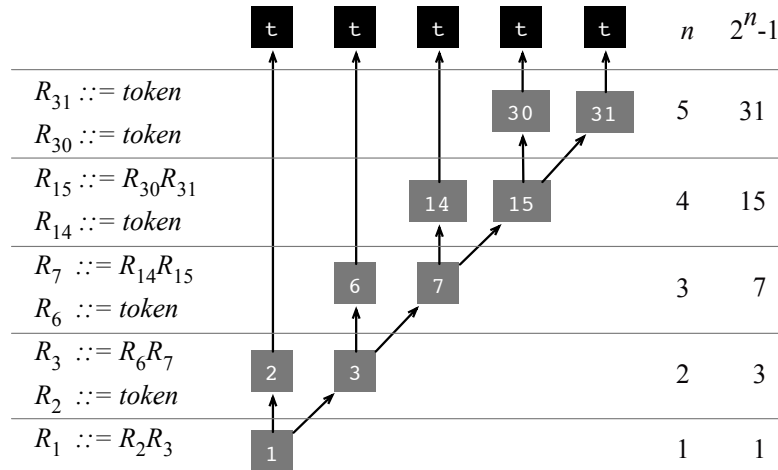


Figure 3.4: The format of the leftmost derivation tree of a Heap-CNF grammar to produce a string with 5 tokens.

The s_e Token. Throughout this work, ZHEFUSCATOR will treat the string event as a single token. As an example, in Figure 3.3, we represent it as s_e , the starting symbol of the event grammar G_e . However, the string event is not a single token; rather, it is a complex sentence pertaining to $lang(G_e)$. Consequently, the sentence represented by s_e does not even need to be formed by the same tokens as the host language. In other words, the tokens in s_e do not, necessarily, need to be recognizable by the `TOKENIZE` function adopted in our implementation of `build_grammar`. That function recognizes tokens from the host language, not from the event language.

Recognizing s_e is necessary when augmenting the current grammar via the `fill_holes` routine. To perform this task, we resort to a brute-force approach: it tries to recognize the largest subsentence $s \in subs(t_i)$ within the active example t_i so that $s \in lang(G_e)$ (See Definition 1). This heuristic is $O(|G_e||t_i|)$, because we can imagine a scenario in which every prefix of t_i is also a prefix of some—incomplete—sentence in $lang(G_e)$.

Nevertheless, the brute-force search tends to fail already in the first token, at least in the setting in which ZHEFUSCATOR was built for: redaction of SQL queries embedded in an unknown language. For instance, consider that the event language is a subset of SQL performing the so called *CRUD* operations, i.e., *SELECT*, *UPDATE*, *CREATE* and *DELETE*. Only sentences that start with one of these four tokens can be part of the event language. Therefore, as soon as the brute-force algorithm stumbles on a different token, it can stop searching immediately. Typical data-representation languages, such as YAML, XML or JSON bear similar properties, meaning that valid sentences in these languages start with a limited number of token combinations.

Furthermore, it is important to consider that the brute force approach is only necessary to augment the current grammar. Whenever line 8 of Figure 3.1 succeeds, no brute-force heuristics are used. As we shall see in Section 3.3.1, in a typical SQL or PostgreSQL log, four to nine samples—among an unbounded number of examples—are

enough to give us a definitive grammar to solve the language separation problem.

3.1.3.1 Merging grammars

Once `fill_holes` produces a grammar G_i for a new example t_i , this grammar is merged into the current grammar G , as it can be seen at Line 26 of Figure 3.1. We define the merging of two Heap-CNF grammars as follows:

Definition 4 (Grammar Merging) *Let $G = \langle R_1, N, T, P \rangle$ and $G_i = \langle R_1, N^i, T^i, P^i \rangle$ be two Heap-CNF grammars. $G' = \langle R_1, N \cup N^i, T \cup T^i, P \cup P^i \rangle$ is the grammar that merges G and G_i .*

Our goal is that G' be also in Heap-CNF and still generates $\text{lang}(G)$ and $\text{lang}(G_i)$. This is achieved by combining the production rules of G and G_i , i.e. by defining the production rules of G' as $P \cup P^i$. Since G and G_i are in Heap-CNF they have the same non-terminals up to $R_{2^{k-1}}$, in which $R_{2^{k-1}}$ is the maximum non-terminal in either G or G_i . So each non-terminal in G' up to $R_{2^{k-1}}$ will generate the combined tokens of G and G_i , while every non-terminal beyond $R_{2^{k-1}}$ has the same production rules of the grammars it comes from, thus generating the same strings that grammar generates.

Example 6 *Figure 3.5 shows the grammars produced for the first, third and fifth lines of Figure 2.1. The tokenization applied to the characters is illustrated in the derivation trees. Each grammar but the first is formed by the merging of a current grammar plus the grammar newly built to match the latest example in the log. Since grammars share non-terminals up to a given index, the union of the production rules has the effect of adding tokens as alternatives to a given non-terminal. For example, in the grammar that generates the first example, R_{15} generates the terminal `int`, while in the second grammar it generates the non terminals $R_{30}R_{31}$, so merging these two grammars results in a grammar in which $R_{15} ::= R_{30}R_{31} \mid \text{int}$. This observation ensures that the merged grammar will be able to generate both the examples generated by the two original grammars.*

Notice that the final grammar that results from merging multiple grammars recognizes a language that is larger than the union of all the examples seen thus far. For instance, the final grammar in Example 6 recognizes the string “`int Query se Query se`”, which encodes two SQL queries.

Lemma 2 (Merging) *If G is the grammar that results from merging two Heap-CNF grammars G' and G'' , then G is Heap-CNF, and $\text{lang}(G') \cup \text{lang}(G'') \subseteq \text{lang}(G)$*

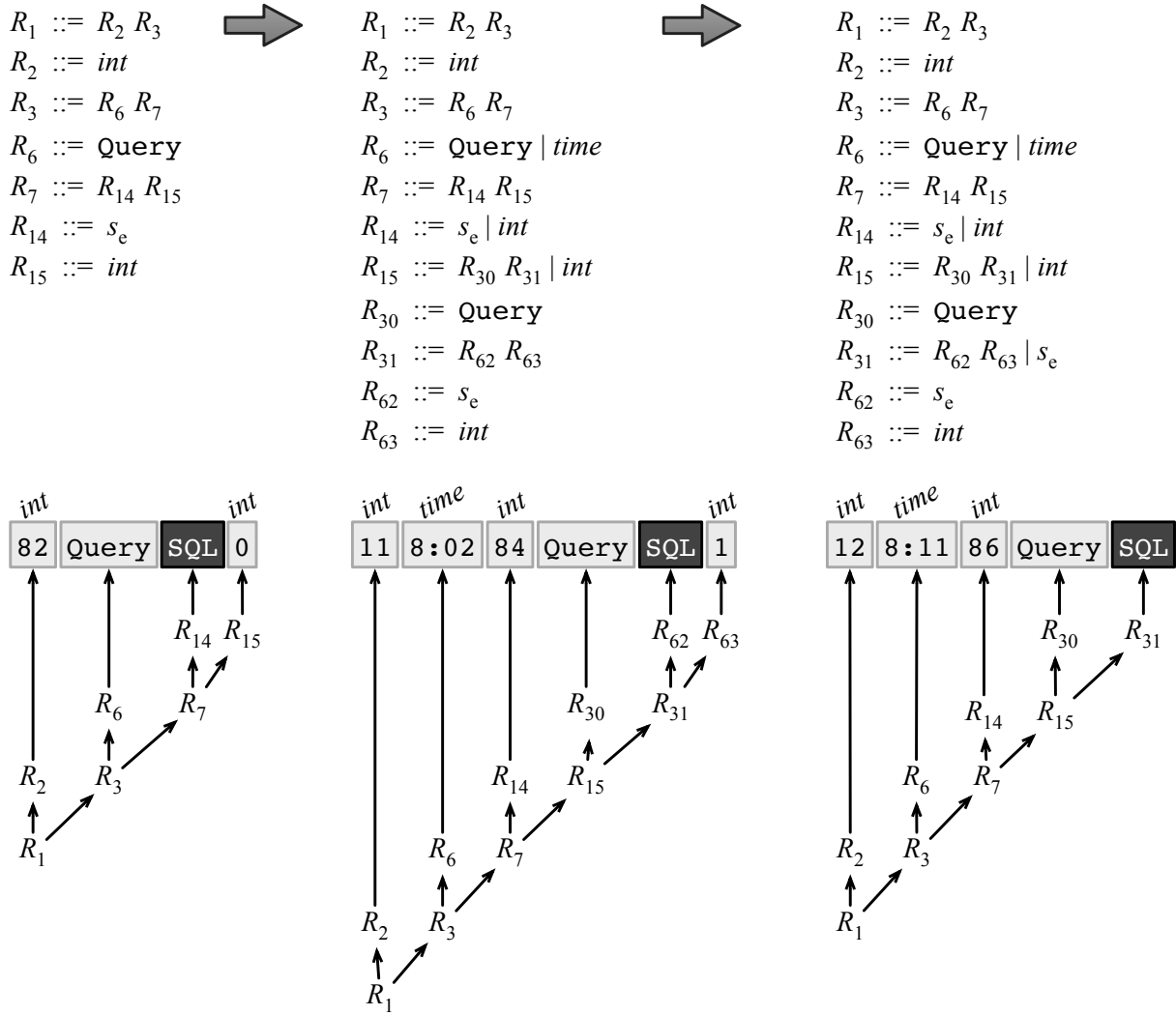


Figure 3.5: Grammars produced from the examples in Figure 2.1.

Proof Sketch 2 We demonstrate the lemma analyzing each one of the four cases involved in the process of merging two Heap-CNF grammars.

Theorem 3 The procedure *build_grammar* (Fig. 3.1) constructs grammars in Heap-CNF.

Proof Sketch 3 The proof of Theorem 3 is the junction of two facts: (i) function *fill_holes* (Fig. 3.2) builds only grammars in Heap-CNF; and (ii) the merging of grammars (Def. 4) yields Heap-CNF grammars.

Theorem 4 (Semantics) Let G_1, G_2, \dots, G_n be the grammars constructed by function *build_grammar* (Fig. 3.1) for input strings t_1, t_2, \dots, t_n . Grammar $G_i, 1 \leq i \leq n$ recognizes every input $t_i, 1 \leq i \leq n$.

Proof Sketch 4 The proof works by induction on the number of examples t_i .

Lemma 3 (Size Complexity) *Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . The size of G_n is $O(N)$, where N is the number of tokens in t_1, t_2, \dots, t_n .*

Proof Sketch 5 *The `fill_holes` procedure only augments the rightmost node of a derivation tree.*

Theorem 5 (Determinacy) *Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . G_n is not ambiguous.*

Proof Sketch 6 *As a consequence of Lemma 3, the rightmost derivation tree of a Heap-CNF grammar always has height $n - 1$ and $O(N)$ nodes.*

Corollary 2 (Time Complexity) *Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . G_n recognizes $t_i, 1 \leq i \leq n$ with $O(N)$ derivations, where N is the number of tokens in t_i .*

Proof Sketch 7 *This corollary follows from Lemma 3, plus the fact, already mentioned in the proof of Theorem 5, that only one rightmost derivation tree is possible.*

3.1.3.2 Limitations: False Positives

The procedure `build_grammar` synthesizes a grammar G that over-approximates the host language L . By over-approximation, we mean that there might exist strings that belong to $\text{lang}(G)$, but that do not belong to L . This observation leads to the notion of *false positives*, which we define as follows:

Definition 5 (False Positive) *Let G_n be the grammar synthesized by `build_grammar` after observing n examples from the host language L . We call a false positive an example t_{fp} such that $t_{fp} \notin L$, $t_{fp} \in \text{lang}(G_n)$ and t_{fp} contains a string event (Definition 1).*

Example 7 (False Positive) *Consider the third grammar seen in Figure 3.5. This grammar recognizes four strings with four tokens: (i) `int query se int`; (ii) `int time se int`; (iii) `int query int int`; (iv) `int time int int`. Only sentences in the format (i) fit the examples seen in Figure 3.5. Sentence (ii) does not correspond to any example and contains a string event (marked by s_e).*

A false positive will lead to the treatment of a string event that, in principle, should be ignored. In the context of this work, ZHEFUSCATOR that will be described in

the next sections, will redact information that is not sensitive. Such action is innocuous in the settings where said system is deployed. Furthermore, the logs that we evaluate in Section 3.3 never lead to false positives. Therefore, we have decided to take no account of false positives in this work. There are two more reasons that led us to ignore them. First, the number of events is unbounded; hence, strings that are false positives up to a certain instant in time might become true positives later. Second, we follow Parsimony’s approach [?] when specifying the event grammar, as we discuss in Section 3.2.1. Parsimony does not support negative examples—a potential way to avoid some false positives.

3.2 Case Study: the ZHEFUSCATOR

We have used the grammar inference techniques discussed in Section 3.1 to implement a system that redacts sensitive information present in program logs. This system is called the ZHEFUSCATOR. ZHEFUSCATOR receives as input an *event language*, given as a grammar G_e , and a running instance of the Java Virtual Machine (JVM). Notice that the ZHEFUSCATOR does not need the source code of the program under execution in the JVM—this program is treated as a black box. Figure 3.6 provides an overview of this tool. In the rest of this section, we discuss particular details of its implementation.

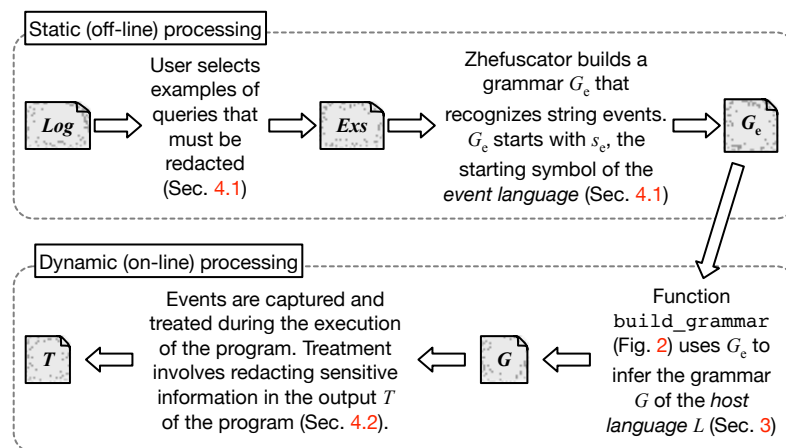


Figure 3.6: ZHEFUSCATOR: event handler for the JVM.

3.2.1 User Interface

ZHEFUSCATOR is meant to be used by professionals who are not necessarily programmers. Therefore, to simplify the task of specifying string events, we provide users with an example-based interface, in which users select substrings from log entries, and a base grammar, G_b , which will be used as a basis for building the event grammar G_e . As a helper to the user, once a substring l is marked for being redacted, all other substrings in the log entries that are recognized by the same rule from G_b which recognizes l are highlighted. That rule is then added to G_e . Through this iterative procedure, the event grammar G_e is built from the basic grammar G_b according to the example substrings selected by the user. Note that users do not deal directly neither with the basic nor with the event grammar: they only deal with textual examples, from which they must choose samples. Currently, we use the SQL grammar as the base grammar, but our implementation is not specific to any grammar. Just keep in mind that if the base grammar does not recognize the example substring selected to be redacted, this example will be ignored in the final event grammar. To determine which data must be redacted, users follow the procedure `markup`, in Figure 3.7.

Procedure `markup`(G_b : Base Grammar, L_f : Log Example)

1. Let G_e be an empty grammar.
2. The user selects a literal l to be redacted, which occurs in a given example from L_f .
3. ZHEFUSCATOR uses the event grammar to extract the largest string $s \in \text{lang}(G_b)$ that contains l .
4. A grammar G'_e , formed with the production rules of G_b necessary to recognize s is constructed.
5. The terminal symbol T that recognizes l is marked to be redacted, T must occur within the rule that recognizes s .
6. G_e is augmented with the rules in G'_e .
7. Every sentence $s' \in L_f$ that G_e recognizes is highlighted.
8. If there are more literals in L_f that must still be obfuscated, the user goes back to step 2.

Figure 3.7: The `markup` procedure that determines which literals must be redacted.

Theorem 6 (Markup) Grammar G_e produced by `markup` (Fig. 3.7) recognizes a subset of $\text{lang}(G_b)$ or the empty language.

Proof Sketch 8 *The proof works by induction on the number of times Step 2 in Procedure `markup` runs.*

As seen in the proof of Theorem 6, grammars G_e and G_b start with the same initial symbol s_e . This symbol is used to compose the instance of the language separation problem (Definition 2) that routine `build_grammar` solves.

3.2.2 Engineering

This section describes details concerning the engineering of the ZHEFUSCATOR—a language-specific problem. For reasons related to the business model in which the authors of this work are involved, ZHEFUSCATOR has been implemented in Java, and deployed onto the Java Virtual Machine. Therefore, it intercepts and treats string events produced by programs written in any programming language that runs on the JVM, including Java, Scala, Kotlin, Clojure and many others. In what follows, we discuss particular aspects of the implementation of this tool.

Parsing ZHEFUSCATOR uses the theory seen in Section 3.1 to build parsers incrementally. These parsers are constructed via the ANTLR [?] parser generation tool. This tool takes as input a grammar that specifies a language and generates as output source code for a recognizer of that language. Procedure `build_grammar` gives ANTLR a new grammar whenever it fails to parse the current text example. ANTLR produces LL(*) parsers, which suits the needs of `build_grammar`, because Heap-CNF grammars are always Left-to-right, Leftmost derivation and can be parsed with one token of lookahead, as the Theorem 7 states. In terms of implementation, we update the grammar by relying on the JVM’s ability to load classes dynamically. The JVM does not need to be restarted in this process. The new grammar is compiled into Java bytecodes by a separate thread, and, as we will see in Section 3.3, such updates take negligible time.

Theorem 7 (LL) *Any Heap-CNF grammar is LL(1).*

Proof Sketch 9 *This fact follows from the observation that Heap-CNF grammars are not recursive.*

Corollary 3 *There are languages whose grammars cannot be synthesized by ZHEFUSCATOR.*

Proof Sketch 10 *This fact follows from the observation that Heap-CNF grammars are not recursive.*

Proof Sketch 11 *A formal language is called an $LL(k)$ language if it has an $LL(k)$ grammar. The set of $LL(k)$ languages is properly contained in that of $LL(k+1)$ languages, for each k greater than or equal to zero [?]. Therefore, there exist context-free languages that are not $LL(1)$.*

The proof of Theorem 7 mentions that Heap-CNF grammars recognize languages with a finite number of possible derivation trees. In fact, strictly speaking, a Heap-CNF language is finite, as the grammar is not recursive. However, in practice, ZHEFUSCATOR deals with infinite languages. Infiniteness comes from the lexer. The procedure `build_grammar` is parameterized by a string tokenizer. In the context of ZHEFUSCATOR’s implementation, this tokenizer is given by ANTLR. The regular language used to recognize tokens can accept an unbounded number of strings. In Section 3.3.2.4 we evaluate the impact of the tokenizer on the performance of ZHEFUSCATOR.

Method interception. ZHEFUSCATOR uses Java Agents to intercept calls to the `System.out.*` singleton object. The Java Agent API [?] provides support to the dynamic instrumentation of JVM applications. Intercepted strings are first fed to `build_grammar`, and then redacted. The first action might result in an expansion of the host language’s grammar. The second might lead to modifications in the output of the program. Literals that must be redacted are specified using the technique discussed in Section 3.2.1.

String Obfuscation. ZHEFUSCATOR performs the redaction of sensitive information via asymmetric cryptography. A sensitive literal l is replaced with a new string l_s , which can be later used as a key to retrieve the true value of l from a classified table. Currently, we use *Advanced Encryption Standard* (AES) to ensure safe redaction of values.

3.2.3 Discussion

The developments explained in this section are necessary to make the ideas introduced in Section 3.1 practical. We do not claim them as contributions, given that the interface and implementation that we adopted have been already discussed in previous work. Our choice for these aspects of our work are pragmatical. On the one hand, the interface discussed in Section 5.4 and the implementation discussed in Section 3.2.2 were effective enough to realize the ideas discussed in this work. However, this choice comes with limitations, which we discuss in the rest of this section.

3.2.3.1 Lack of Negative Examples

The main limitation of our example-based approach is a lack of negative examples. This limitation is also present in Parsimony [?]; hence, it has naturally persisted in our implementation of it. We opted to avoid negative examples because it is our understanding that in most of the cases where ZHEFUSCATOR is useful, negative examples are unnecessary. In other words, database logs tend to follow simple formats, with a small set of sentences of interest. Nevertheless, if necessary to handle more complex formats, then ZHEFUSCATOR might produce false positives. In the context of this work, as explained in Section 3.1.3.2, false positives might cause the redaction of sentences that do not contain sensitive information.

3.2.3.2 Expressiveness

Additionally, an example-based interface lacks resources that would be promptly available in a domain-specific language, such as the ability to specify logical combinations of events. For instance, users could be interested in enabling certain events only after particular events of interest have been detected. ZHEFUSCATOR current interface lacks such sequencing operations. Users interested in such ability are encouraged to use ZHELANG, which a DSL that we have defined for the treatment of string events. ZHELANG will be discussed in Chapter 4.

3.3 Evaluation

We have implemented the techniques discussed in this work onto an actual on-line obfuscator, which we call the ZHEFUSCATOR. ZHEFUSCATOR is open source and can be used to redact queries produced by database logs. This section investigates the following research questions related to this implementation, as well as the techniques that support it:

- **RQ1—Convergence:** how many examples are necessary to produce grammars for languages typically used by SQL logging systems?

- **RQ2—Effectiveness:** are the parsers derived from the synthesized grammars effective?
- **RQ3—Practicality:** what is the runtime overhead of ZHEFUSCATOR when deployed onto a database system dealing with a heavy workload?

We chose these three particular questions to demonstrate that the theory developed in Section 3.1, and its implementation described in Section 3.2, once combined into a concrete tool, lead to a system that is not only novel, but also practical.

Runtime Setup. Every result reported in this section has been produced on an 8-core Intel(R) Core(TM) i7-3770 at 3.40GHz, with 16GB of RAM running Ubuntu 16.04.

3.3.1 RQ1—Convergence

Methodology. To answer RQ1 we measure how many times the predicate `successfull_parse`, invoked at Line 8 of Figure 3.1, fails before we produce a definitive grammar for a certain log generator. We perform this analysis on logs from two database systems and from the OSX operating system. Logs are given as a *text* of examples t_i , as defined in Section 3.1. Each t_i is the entire output produced by the generator, be it a database, be it the operating system, at time unit i . To determine the parts of the log that should be obfuscated, we chose, from each one, four examples, following the steps enumerated in Figure 3.7. We chose the first four sentences that did not fit into the same SQL production rule. However, this choice bears no impact on the results reported in this work. Convergence does not depend on it, and the time to redact strings (running time will be evaluated in the next section) is the same for the different approaches that we compare.

3.3.1.1 Logs from Database

On this experiment, we have generated logs from two different SQL Databases: MySQL version 14.14 Distribution 5.7.27 and PostgreSQL version 9.2.24. Workloads for these two databases were produced by the 9 real-world web applications emulated by OLTP-Bench [?], which include systems such as Wikipedia, Twitter and an ordinary seats system.

Discussion. Figure 3.8 shows the average prefix necessary to synthesize a grammar in different database systems. ZHEFUSCATOR requires approximately eight examples to infer a grammar for the logs produced by MySQL, and five for those produced by PostgreSQL. In the former collection, logs contain an average of 662K lines; in the latter, 1,867K. This experiment indicates that, for typical database logs, the grammar inference procedure of Section 3.1 tends to converge to a definitive parser after five to eight examples. Furthermore, these examples are a very small portion of the entire log: in every case, we had a definitive grammar after observing less than 0.01% of the whole log file.

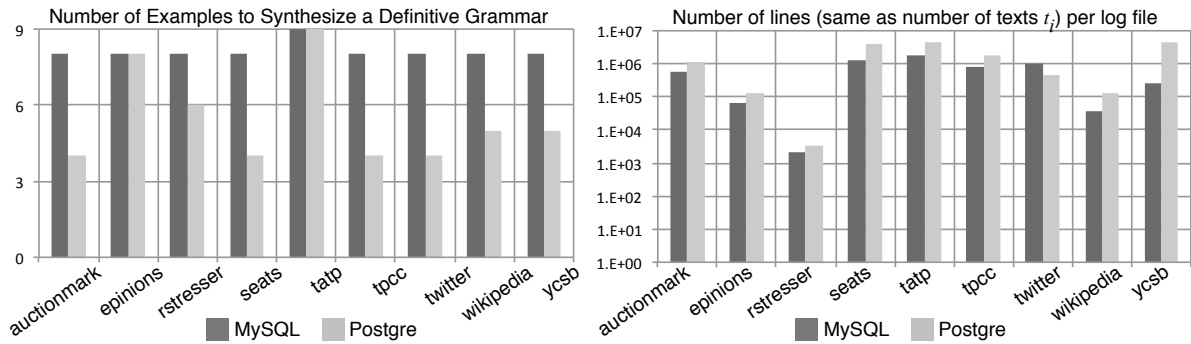


Figure 3.8: Average prefix size necessary to synthesize a grammar for different log files produced by either MySQL or PostgreSQL.

3.3.1.2 Logs from the Operating System

This experiment uses the logs produced by default by MacOS version 10.14.6 in the `/var/log` directory. Contrary to the examples that use the databases, these logs are very different one from the other (the format of sentences is not shared across them). This fact will be made clear once we analyze how many examples are necessary for synthesizing a definitive grammar—this number varies substantially across the logs. We gathered four logs from five distinct OS users, whose usage pattern corresponds to the profile of professional programmers. The logs used in this experiment are:

- `corecaptured.log`: logs operations of the network hardware. On average, these logs have 174K lines.
- `wifi.log`: logs network traffic. On average, they contain 9K lines.
- `system.log`: logs the operations executed in the whole system. On average, they contain 4K lines.
- `fsck_apfs.log`: logs file system operations, and contain 4K lines on average.

Discussion. Figure 3.9 shows the average prefix necessary to synthesize grammars for the OSX logs. The number of required examples is higher than what has been observed in Section 3.3.1.1. The ratio of examples per log size is also higher. In one case (`user3:system`) we had a log with only five lines, whose grammar demanded three examples. This case is an anomaly, due to the small log size. The largest prefix consisted in 170 examples, for a log with 6,579 samples (`user1:system`). In general, the ratio of examples per sample is still very low. For instance, our largest logs (`corecapture`) have almost 200K lines on average, and yet our on-line grammar inference engine finds a grammar that recognizes all these samples after observing 57 to 64 examples.

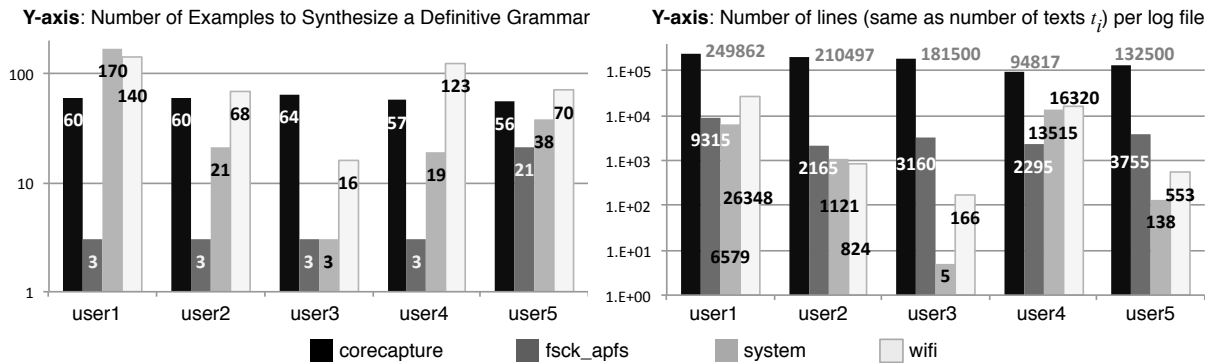


Figure 3.9: Average prefix size necessary to synthesize a grammar for different MacOS logs.

3.3.2 RQ2—Effectiveness

This section evaluates the practicality of the grammars synthesized by ZHEFUSCATOR. To this effect, we shall answer the five questions below. BF refers to the *Brute Force* approach, which searches event patterns exhaustively within text examples:

1. Section 3.3.2.1: how does ZHEFUSCATOR compare to BF to parse one individual example for which a parser has not already been synthesized.
2. Section 3.3.2.2: how does ZHEFUSCATOR compare to BF to parse 1,000 examples in an actual log file produced by a MySQL database.
3. Section 3.3.2.3: how does ZHEFUSCATOR compare to BF to parse 1,000 examples in artificially generated logs of different sizes.
4. Section 3.3.2.4: how does the tokenizer change the runtime of ZHEFUSCATOR.

3.3.2.1 Parsing Effectiveness

There exists a trivial approach to solve the Language Separation Problem introduced in Definition 2: given an example t_i in the host language, we start a search for an event s , an SQL query in our context, at every token of t_i . If two events can start at the same token, we choose the longest one. This solution is called the *brute-force* approach. The developments in Section 3.1 are attractive inasmuch as they lead to a faster solution to language separation than the brute-force technique. In this section, we compare the parsing speed of both approaches.

Before we discuss our methodology, two observations are in order. First, when ZHEFUSCATOR’s current grammar is not able to recognize the active example, it behaves in a similar manner as the brute force approach: it must scan the SQL query, assuming that it can start at any token. In addition to this, it must augment the current grammar using the techniques discussed in Section 3.1. Second, when ZHEFUSCATOR’s parser is able to recognize the active example, parsing happens via $O(N)$ productions, where N is the number of tokens. Yet, the number of characters per token varies, and the lexer’s runtime must be taken into consideration. Thus, the overall runtime is $O(M)$, where M is the number of characters in the active example. The brute force approach might expand $O(N^2)$ productions; however, such worst case seldom happens. Most of the tokens in a valid example cannot be the prefix of any SQL query. Therefore, although naïve, the brute force approach is still likely to outperform Zhefusctor for examples with few characters.

Methodology. The brute-force approach becomes less practical as the number of characters in the examples t_i of the host language L increases. To investigate at which point the grammars synthesized by `build_grammar` become more efficient, we have used the logs seen in Section 3.3.1.1. To obtain examples of varying sizes, we either split or concatenate lines from these logs; hence, producing strings of different lengths.

Discussion. Figure 3.10 compares the brute-force with our synthetic grammars. Our grammars are more asymptotically efficient than the brute-force approach. After multiple merging operations, a Heap-CNF grammar still recognizes a sentence in $O(N)$ derivation steps, where N is the size of the sentence. The brute force approach, in turn, will always require $O(N^2)$ steps. Figure 3.10 shows that for examples between 128 and 256 characters (about 16 tokens) our approach becomes consistently better than the trivial brute-force parsing. In Section 3.3.2.2 we observe the effect of this improvement when applied onto an actual log.

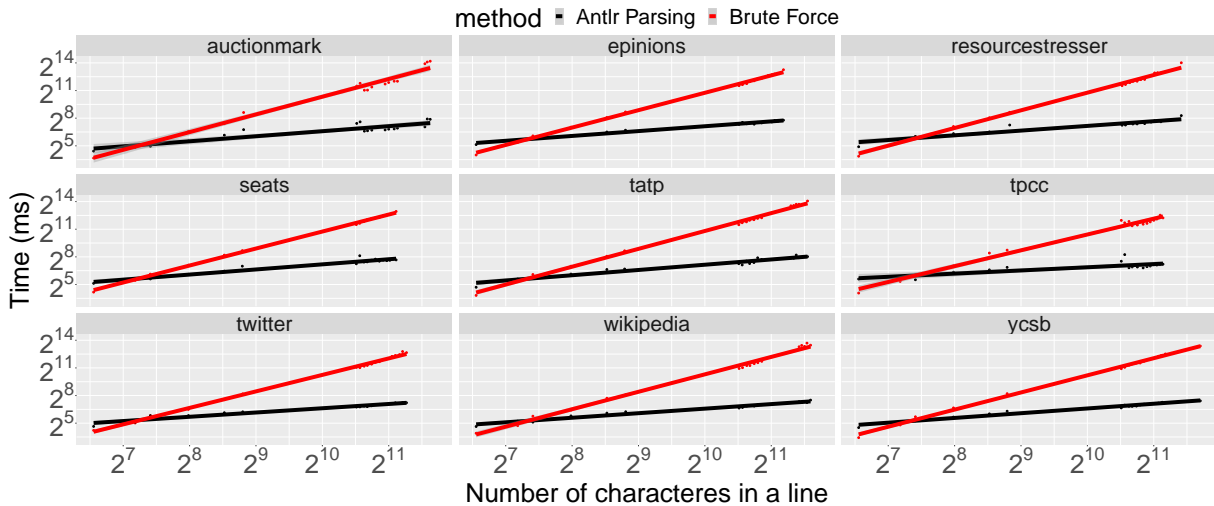


Figure 3.10: Time comparison of brute force approach and the ANTLR parser.

3.3.2.2 Effectiveness on an Actual Log File

In Section 3.3.2.1, we compared the average time taken by the ZHEFUSCATOR and the brute force approach to parse one example. However, the benefit of our parser synthesizer becomes more evident once we analyze its effect when amortized onto a long chain of examples. In this section, we analyze this effect via *skyline* charts. These charts show the time taken per individual example in the log. For this experiment, we chose the log produced by the MySQL implementation for the AUCTIONMARK application. We emphasize that the choice of log, for this experiment, is immaterial: all the logs produced by MySQL follow the same format, and ZHEFUSCATOR’s parser needs to be augmented only 8 times for all of them. AUCTIONMARK has been chosen simply because it is the first benchmark in OLTPBench.

Methodology. We compare both the approaches, ZHEFUSCATOR and the brute force, when given the first 1,000 examples in the log that MySQL produces for AUCTIONMARK. For each example, we count only the time to recognize strings—redaction is not accounted for, because it applies the same algorithm, the same number of times, in both the approaches. Notice that choosing more than 1,000 examples will not change the results reported in this section, because ZHEFUSCATOR builds a definitive parser after observing 19 entries in the log file.

Discussion. Figure 3.11 shows the result of this experiment, juxtaposing the skyline produced by the brute force approach and by ZHEFUSCATOR. The log file contains two distinct parts. The first 250 examples are system configuration commands, and have 1,021 characters, on average. The last 750 examples are various SQL queries, and contain 106 characters on the average. Using the C tokenizer, we obtain 105 tokens, on the average,

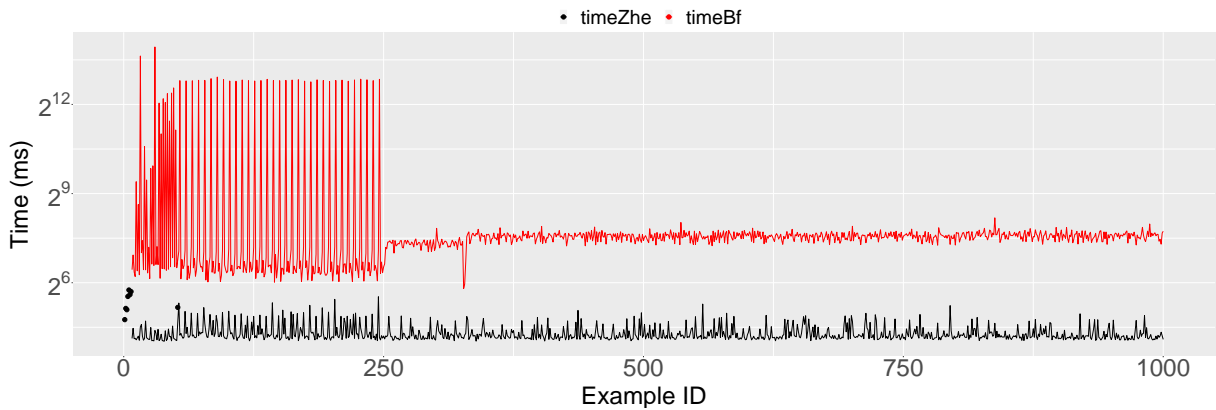


Figure 3.11: Skilene comparison between ZHEFUSCATOR and the brute force approach to parse 1,000 examples from the AUCTIONMARK log.

considering the 1,000 examples in the log file. Under this circumstance, the performance gap between ZHEFUSCATOR and the brute force approach is noticeable.

ZHEFUSCATOR spends, on the average, 26.15 milliseconds per example, with a standard deviation of 17.86 ms. This number includes the extra time ZHEFUSCATOR needs to augment the current parser—an action that happened 8 times in this experiment. The brute force approach spends 586.63 milliseconds per example, with a standard deviation of 1,830.33 ms. ZHEFUSCATOR is 22.5x faster, per example, than the brute force approach. However, this experiment uses an ideal scenario for ZHEFUSCATOR: a long stream of homogeneous textual examples. In the next section, we shall analyze the behavior of ZHEFUSCATOR under more unfavorable conditions.

3.3.2.3 Increased Effectiveness via Amortized Cost

The logs produced by MySQL and PostgreSQL are formed by long individual examples (more than 100 tokens on average). However, these examples are all similar; hence, as already observed in Figure 3.8, ZHEFUSCATOR synthesizes a definitive parser after observing a very short subset of them. To stress out the performance of ZHEFUSCATOR, in this section we analyze its behavior when dealing with more complex logs, which we have produced artificially.

Methodology. To produce the logs, we use six different types of tokens: booleans, integers, doubles, strings, dates and sets of comma-separated integers within curly brackets, e.g., $\{2, 3, 5, 7\}$. We generate four types of logs. Each log contains a random number of tokens between 0 and $R \in \{4, 8, 16, 40\}$, before and after an SQL query. We use always the query “SELECT *string* FROM *string* WHERE *id* = *int*”. With $R = 4$, we have

$4^6 + 4^6 = 8,192$ possible example formats; with $R = 8$, we have $8^6 + 8^6 = 524,288$, and so on. Therefore, `fill_holes` will be invoked a much larger number of times than in the setup used in the previous section.

Discussion. Figure 3.12 shows the result of this experiment. Whereas BF shows homogeneous behavior—its runtime per example varying only slightly—ZHEFUSCATOR has two types of responses. Such responses depend on the current parser recognizing or not the active example. When recognition is possible, parsing is fast; otherwise, the parser must be augmented with new productions, and we observe a runtime spike, which is marked in Figure 3.12 with a black dot. Said spikes are compulsory for the initial examples. However, as the current grammar increases, sentence recognition becomes more common, and spikes tend to disappear. As a consequence, the more events are observed, the larger is the performance improvement of ZHEFUSCATOR over the brute force approach.

Figure 3.13 shows average time per example, plus standard deviations observed for ZHEFUSCATOR and for the brute force approach. The figure shows two results for ZHEFUSCATOR: the first considers only the time when parsing succeeds; the second considers, in addition, the time taken by `fill_holes`, when ZHEFUSCATOR fails. In the former scenario, ZHEFUSCATOR always outperforms the brute force approach. In the latter, it always loses. The conclusion is that, once it reaches a steady state, ZHEFUSCATOR’s $O(N)$ parser is consistently a better option than BF’s $O(N^2)$ algorithm. However, if necessary to augment the current parser too often, our technique loses its attractiveness. In this particular experiment, `fill_holes` performs worse than in Section 3.3.2.2, because the host language is much more complex.

3.3.2.4 Impact of the Tokenizer on Runtime

The `add_example` routine, which is invoked by `build_grammar` (Figure 3.1, Lines 7-15) is parameterized by a tokenizer. The tokenizer is a function that converts the input text into tokens. The tokenizer is just an artifact of our implementation: users of our system will never have to deal with it. The implementation of ZHEFUSCATOR can use any tokenizer that ANTLR supports. As we have hinted in Section 3.1.2, the tokenizer impacts both the number of examples as well as the runtime of ZHEFUSCATOR. In this section, we analyze this impact by verifying the behavior of ZHEFUSCATOR when parameterized by two different lexers.

Methodology. We have tried ZHEFUSCATOR with two different lexers. Both were taken from public projects that use ANTLR—they have not been implemented as part of this research.

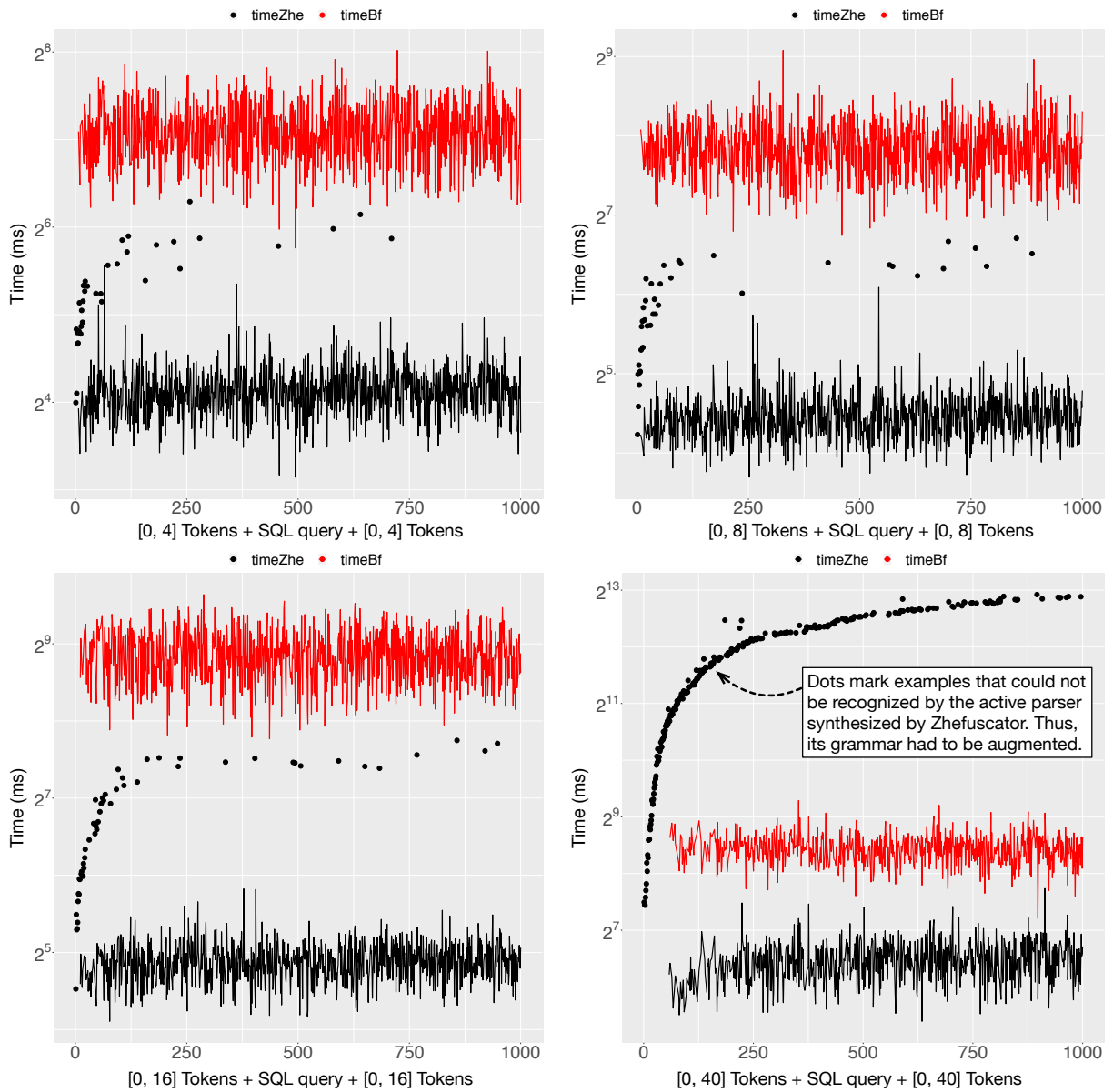


Figure 3.12: Runtime comparison between ZHEFUSCATOR and the brute force approach to parse 1,000 examples of artificially generated logs. Black dots mark invocations of `fill_holes`.

Discussion. Although the choice of tokenizer might modify the number of examples necessary to reach a definitive grammar, the two tokenizers that we have used led to the same prefix size in Figures 3.8 and 3.9. This happens because C and SQL have many similar tokens, including identifiers—the most common in the examples. However, the impact on runtime is different. Using the C tokenizer, ZHEFUSCATOR takes 26.15 milliseconds, on average (STD = 17.85ms), per example from the AUCTIONMARK log (Figure 3.11), including the eventual time taken to augment the grammar. Using the SQL version, this time drops down to 18.81 milliseconds, with a standard deviation of 3.33ms. The latter is faster because the SQL lexer uses a smaller automaton than the C lexer.

Format	[0,4]+SQL+[0,4]		[0,8]+SQL+[0,8]		[0,16]+SQL+[0,16]		[0,40]+SQL+[0,40]	
Tk/Ex	29.99		38.26		54.10		298.43	
	parsing only	fill_holes: 35	parsing only	fill_holes: 39	parsing only	fill_holes: 51	parsing only	fill_holes: 264
avg (Zhe)	17.53	18.40	22.08	23.67	29.47	34.05	88.88	1,085.01
std (Zhe)	3.45	6.32	4.51	10.09	5.00	23.70	21.15	2,004.49
avg (BF)	137.15		232.14		464.65		347.20	
std (BF)	31.41		55.38		103.24		60.48	

Figure 3.13: Average time and standard deviation (per example, in milliseconds) that ZHEFUSCATOR (Zhe) and the brute force approach (BF) take to analyze the artificial logs. “Parsing only” reports runtimes for examples in which ZHEFUSCATOR’s current parser succeeds without having to synthesize a new grammar. “fill_holes: XX” includes the time of “parsing only”, plus the time to augment the current parser. XX reports the number of times ZHEFUSCATOR had to augment the current parser (via the `fill_holes` routine).

3.3.3 RQ3—Practicality

The techniques described in Section 3.1 have a computational cost. The goal of this section is to measure such cost. This empirical evaluation shall allow us to claim that the overhead of ZHEFUSCATOR, when deployed onto typical Java applications, is low enough to be practical.

Methodology. It is difficult to measure the overhead of ZHEFUSCATOR in our experimental setup involving actual deployments of MySQL and PostgreSQL. This difficulty comes from the fact that logging, at least in that particular setting, is a rare event. Log entries are produced only when users enter queries in the database. In this scenario, the overhead of ZHEFUSCATOR is negligible. Thus, to probe this overhead in a more heavily loaded scenario, we shall proceed with two experiments. In Section 3.3.3.1 we measure the runtime overhead that event handling imposes onto a single invocation of the `System.out.println` routine used to output log information in a database server. This evaluation provides some insight into the absolute overhead of event handling; however, it does not give us much information about how ZHEFUSCATOR would impact user experience, for the time of handling one single string event is very fast. To circumvent this limitation, in Section 3.3.3.2 we measure the overhead that ZHEFUSCATOR imposes onto batch computations, i.e., that perform a fixed number of steps. In this case, we focus on the Java Dacapo benchmark suite [?].

3.3.3.1 Overhead of Treating one String Event

To measure the overhead of treating one string event, we have built a system that reads a log file and outputs it line by line using the *System.out.println* method from the Java Standard Library. For maximum stress, we assume that every SQL literal must be redacted. In this experiment, we adopt the same logs from the MySQL databases used in Section 3.3.1.1.

Discussion. Figure 3.14 presents the results of this evaluation. Each log was evaluated ten times; hence, each box plot contains ten samples. The figure makes it clear that ZHEFUSCATOR’s event handler has an overhead over individual method invocations. This overhead can be as high as two orders of magnitude, as observed in *resourcestresser*. However, this cost accounts for a very small proportion of the runtime of a typical database system. In the case of *resourcestresser*, the average time to redact every literal in the log is 0.03sec per invocation of *System.out.println*. This time includes the invocation of *build_grammar* (Fig. 3.1) and the obfuscation of literals. Obfuscation includes the time to encrypt literals using AES.

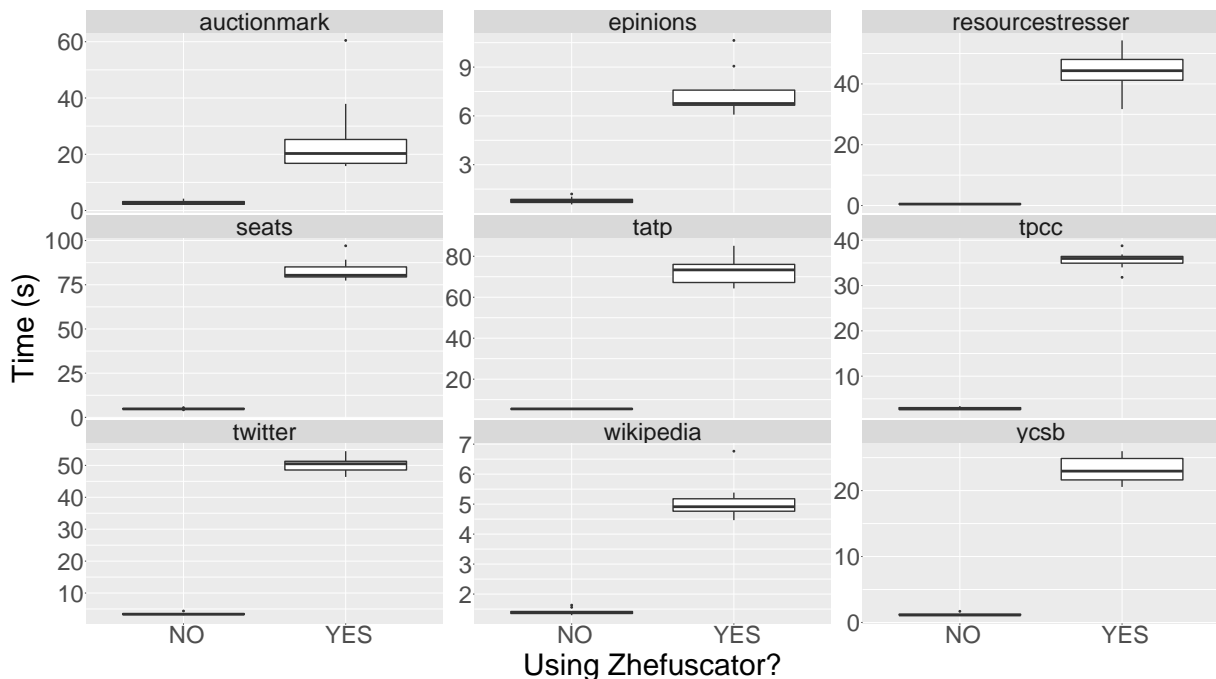


Figure 3.14: Overhead of ZHEFUSCATOR on an extreme case: a system that only outputs different database logs.

3.3.3.2 Deploying on Java Dacapo

In this experiment, we measure the overhead of building a grammar for every output produced by the programs in the DaCapo Benchmark Suite. DaCapo’s logs do not contain SQL queries; hence, in this section, we are measuring the time to build grammars, but not the time to redact queries.

Discussion. Figure 3.15 compares the runtime of DaCapo without and with interventions from ZHEFUSCATOR. Figure 3.16 shows accompanying data: p-values, number of log events and number of production rules in the final grammar that we synthesize. The p-value provides us with some notion of statistically significant runtime difference: the lower the p-value, the more noticeable is the gap in runtime between the two versions of each DaCapo program. Typically, p-values below 0.05 are considered statistically significant. These p-values have been obtained via a T-Test applied on the same data used to produce Figure 3.15. The T-Test provides us with an idea on how different are a “control” and a “test” groups. In our setting, the control group is formed (in Figure 3.15) by applications that do not run the ZHEFUSCATOR. The test group, in contrast, is formed by the same applications using the ZHEFUSCATOR. The lower the p-value returned by the T-Test, the more statistically significant is the different between these two groups.

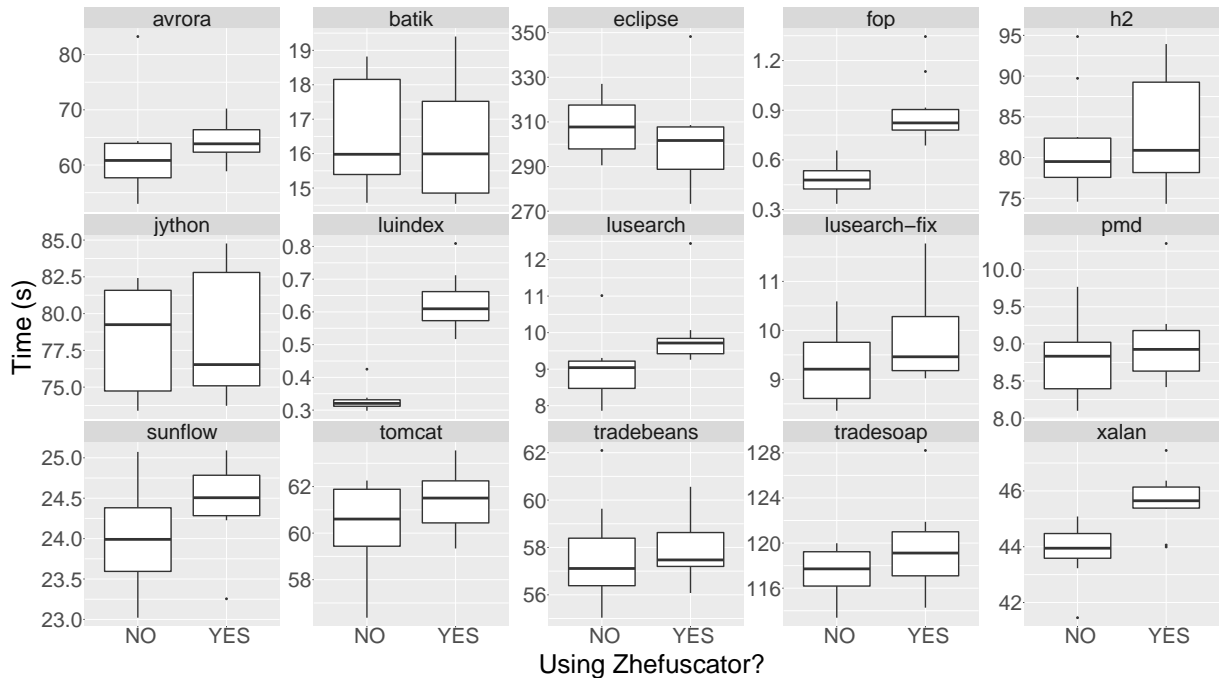


Figure 3.15: The overhead of ZHEFUSCATOR on Dacapo.

The runtime overhead of ZHEFUSCATOR, even when deployed onto a batch system, tends to be small. In 11, out of 15 cases, we could not perceive any statistically significant runtime difference. The largest runtime gap that we have observed was in `fop`; however,

	avro	batik	eclipse	fop	h2	jython	luindex	lusearch	lusrch-fix	pmd	sunflow	tomcat	tradebeans	tradesoap	xalan
P-values	0.47	0.68	0.38	0	0.54	0.99	0	0.03	0.17	0.52	0.07	0.13	0.66	0.16	0
Log Lines	13	22	25	13	24	94	13	45	45	13	13	19511	31	31	14
Productions	30	60	10	16	68	30	16	46	46	30	46	56	54	54	46

Figure 3.16: The overhead of ZHEFUSCATOR on Dacapo. The lower the p-value, the more statistically significant the overhead.

this is the benchmark that runs for the shortest time. Thus, this overhead, due to the initialization of ZHEFUSCATOR's agent, tends to be amortized in systems that run for more time. The largest absolute overhead was observed in `xalan`: 1.7 seconds on average, over a system that runs for 44 seconds on average.

Chapter 4

ZHELANG

4.1 Language Specification

This section defines the syntax (Sec.4.1.2) and the semantics of ZHELANG (Sec.4.1.3), a reactive domain specific language to threat string events. However, before we dive into any formalism, Section 4.1.1 provides some rationale on the choice of operators that constitute ZHELANG.

4.1.1 Core Operators

ZHELANG programs are made of the combination of five different operators: `atomOp`, `combOp`, `orOp`, `optOp` and `searchUntilOp`. The first operator, `atomOp`, determines *patterns of interest*, i.e., regular expressions that trigger events once they occur in the input stream. The other operators are *combinators*, meaning that they let the user specify combinations of events. The first four operators, `atomOp`, `combOp`, `orOp` and `optOp`, let users specify *context free grammars*. The next example illustrates this possibility.

Example 8 *Figure 4.1 (a) shows a context-free grammar. Figure 4.1 (b) shows a ZHELANG program that recognizes the same grammar. We shall postpone the discussion of the semantics of this program to Section 4.1.3.*

The fifth operator, `searchUntilOp`, exists for two reasons. First, to enable the repetition of multiple patterns as triggers of events. Second, because ZHELANG programs need a way to skip text that is not of interest. In other words, a ZHELANG program encodes context-free grammars that recognize specific patterns within an infinite text stream. The full language that generates this infinite text does not need to be recognized by the program.

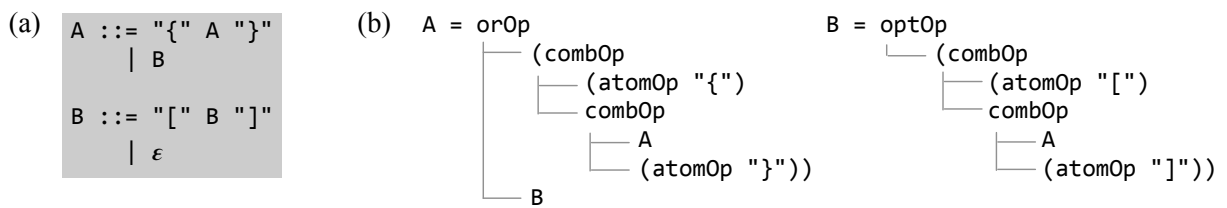


Figure 4.1: (a) A context-free grammar that recognizes balanced sequences of braces and brackets. (b) The same grammar implemented as a combination of four ZHELANG operators.

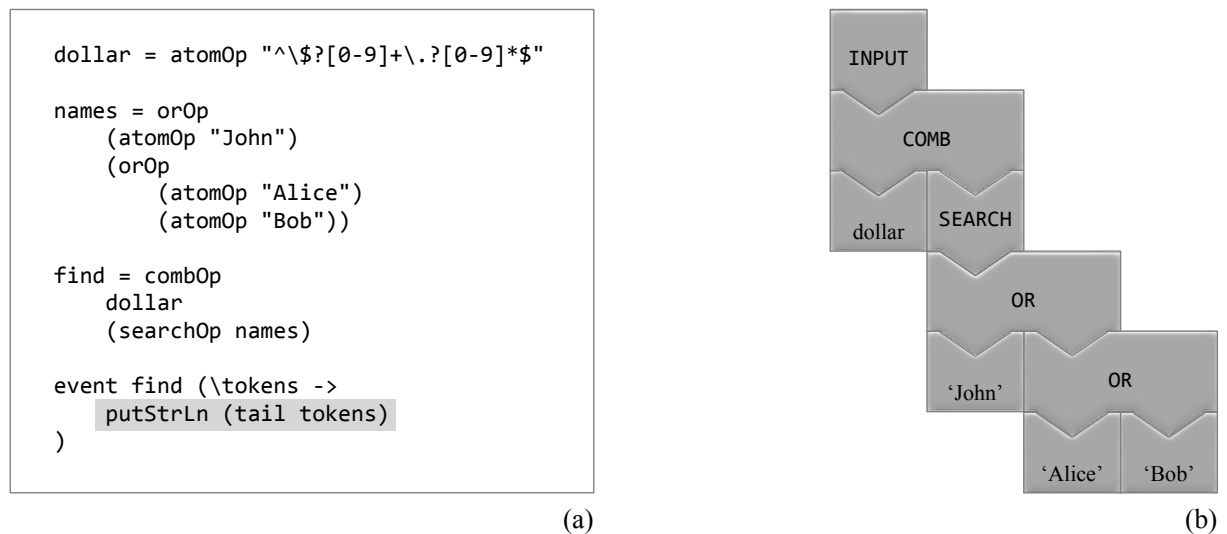


Figure 4.2: (a) A ZHELANG program that detects all the occurrences of either one of three particular names after observing a currency pattern. The action language, embedded in the guard, is Kotlin. We have highlighted the only Kotlin command (`println`) used. (b) The block representation of the program.

Example 9 Figure 4.2 (a) shows an example of a ZHELANG program. This program prints all the occurrences of any of three different names (John, Alice or Bob) after having observed the occurrence of a pattern that represents an amount of dollars, e.g., \$ 2, \$ 2.33, \$ 2.333. Thus, the underlined strings in the sentences below would be printed:

- John owns \$ 2.30 to Alice Tokien and John Wayne.
- Amount: \$ 2.30; Name: Alice Alba; Name: Alice Page.

How this semantics is produced by the syntax seen in Figure 4.2 (a) is yet to be discussed in Section 4.1.3. ZHELANG supports a graphic representation of programs, which Figure 4.2 (b) shows. Section 4.1.2.1 provides more detail about this visual program representation.

4.1.2 Syntax

A ZHELANG program is written according to the grammar in Figure 4.3. A program is a sequence of *guards* and *events*. Guards represent bindings between names and *observers*. Events represent bindings between observers and programs written in an *action language*. The latter—any programming language that can be linked with ZHELANG—lets users specify actions to be executed once an event is detected. As Figure 4.3 shows, actions are invoked as lambda expressions. In this work, we assume, for simplicity, a curried representation.

```

observer ::= atomOp regex                ; Atomic Operators
           | combOp observer observer   ; Combinator Operator
           | orOp observer observer     ; Optional Operators
           | searchUntilOp observer observer ; Search Until Operator
           | optOp observer             ; Optional Operators
           | Fail                         ; A pattern matcher that always fails
program  ::= event observer action     ; Event Definition
           | cutEvent observer action   ; Cut Operator
action   ::=  $\lambda v \rightarrow E(v)$  ;  $E$  is any expression in the action language

```

Figure 4.3: The syntax of ZHELANG.

Observers are operators to recognize patterns in the input stream. Currently, ZHELANG provides users with five types of observers. Their semantics is explained in Section 4.1.3. For now, it suffices to know that they can be combined, as a recursive data type. Example 9 shows a program that results from this combination.

4.1.2.1 Visual Representation

To ease program understanding, we have equipped ZHELANG’s observers with a visual representation. As previously mentioned, observers form a recursive data type. Visualization relies heavily on this recursive nature. Figure 4.4 (a) shows the five blocks of observers used in ZHELANG. The horizontal length of two of these blocks, denoting *combOp* and *orOp*, is variable. This length can be extended whenever necessary to accommodate multiple observers in a contiguous block, i.e., with successive applications of the binary observers, simulating an n -ary application. The next example illustrates this fact.

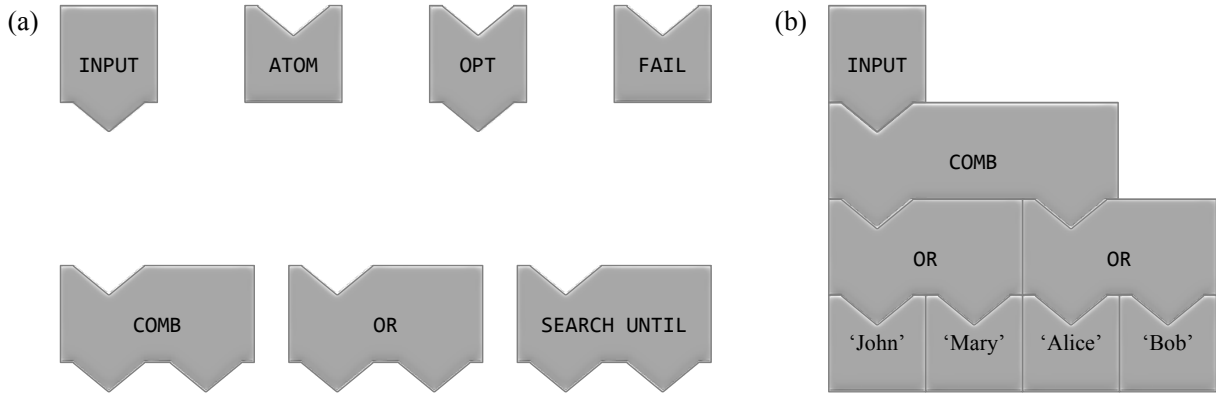


Figure 4.4: (a) The seven observers currently defined in ZHELANG (we use a symbolic representation for the input stream—`input`; however, this symbol does not denote an observer). (b) Program that recognized the pattern `('John'|'Mary')('Alice'|'Bob')` in the input stream.

Example 10 *Figure 4.2 (b) shows a pictorial representation of the program earlier seen in Example 9. Notice that this program, contrary to the one seen below in Example 11, recognizes multiple (indeed infinitely many) occurrences of the pattern (“John” | “Alice” | “Bob”). This fact is due to the use of the `loopOp` observer.*

Example 11 *Figure 4.4 (b) shows a program that recognizes events formed by two tokens, in sequence. The first token is either the string John or the string Mary. The second token is either Alice or Bob.*

4.1.3 Semantics

Figure 4.5 shows the semantics of the different observers that constitute ZHELANG. Each observer returns a triple BL_tL_s , in which $B \in \{\text{Satisfied}, \text{Unsatisfied}\}$; L_t is a list of tokens that have been recognized by the observer; and L_s is the rest of the infinite input stream. We say “rest” because observers might modify the input stream after processing it. For example, the `atomOp` observer removes the head token of the input stream if a match happens.

Although the return type of each observer is the same, their input types might differ. The operator `optOp` has two curried arguments: The first is another observer and the second is the input stream. The operator `searchUntilOp` takes three arguments: another observer, a pattern that when matched ends the execution of the operator, and the input stream. The operator `atomOp` also has two arguments. The second is also the input stream, but the first is a pattern to be matched. The other two observers, `combOp`

$$\begin{array}{c}
\frac{match(t, h) = \text{True}}{\text{atomOp } t (h : s) \text{--} > \text{Satisfied } [t] s} \{At1\} \quad \frac{match(t, h) = \text{False}}{\text{atomOp } t (h : s) \text{--} > \text{UnSatisfied } s} \{At2\} \\
\\
\frac{c_1 \text{ input--} > \text{UnSatisfied } s}{\text{combOp } c_1 c_2 \text{ input--} > \text{UnSatisfied } s} \{Co1\} \\
\\
\frac{c_1 \text{ input--} > \text{Satisfied } t_1 s_1 \quad c_2 s_1 \text{--} > \text{UnSatisfied } _}{\text{combOp } c_1 c_2 \text{ input--} > \text{UnSatisfied } (\text{tail input})} \{Co2\} \\
\\
\frac{c_1 \text{ input--} > \text{Satisfied } t_1 s_1 \quad c_2 s_1 \text{--} > \text{Satisfied } t_2 s_2}{\text{combOp } c_1 c_2 \text{ input--} > \text{Satisfied } (t_1 ++ t_2) s_2} \{Co3\} \\
\\
\frac{c_1 \text{ input--} > \text{UnSatisfied } _}{\text{orOp } c_1 c_2 \text{ input--} > c_2 \text{ input}} \{Oo1\} \\
\\
\frac{c_1 \text{ input--} > \text{Satisfied } t_1 s_1}{\text{orOp } c_1 c_2 \text{ input--} > \text{Satisfied } t_1 s_1} \{Oo2\} \quad \frac{}{\text{Fail input--} > \text{UnSatisfied } (\text{tail input})} \{Fl1\} \\
\\
\frac{c_1 \text{ input--} > \text{Satisfied } t_1 s_1}{\text{searchUntilOp } c_1 c_2 \text{ input--} > \text{Satisfied } t_1 s_1} \{Su1\} \\
\\
\frac{c_1 \text{ input--} > \text{UnSatisfied } _ \quad c_2 \text{ input--} > \text{Satisfied } _}{\text{searchUntilOp } c_1 c_2 \text{ input--} > \text{UnSatisfied } (\text{tail input})} \{Su2\} \\
\\
\frac{c_1 \text{ input--} > \text{UnSatisfied } s \quad c_2 \text{ input--} > \text{UnSatisfied } s}{\text{searchUntilOp } c_1 c_2 \text{ input--} > \text{searchUntilOp } c_1 c_2 s} \{Su3\} \\
\\
\frac{c_1 \text{ input--} > \text{Satisfied } t s}{\text{optOp } c_1 \text{ input--} > \text{Satisfied } t s} \{Op1\} \quad \frac{c_1 \text{ input--} > \text{UnSatisfied } _}{\text{optOp } c_1 \text{ input--} > \text{Satisfied } [] \text{ input}} \{Op2\} \\
\\
\frac{\text{searchUntilOp observer Fail stream--} > \text{Satisfied } t s}{\text{event observer action stream--} > \text{action}(t); \text{event observer action } s} \{Eg1\} \\
\\
\frac{\text{searchUntilOp observer Fail stream--} > \text{Satisfied } t s}{\text{cutEvent observer action stream--} > \text{action}(t);} \{Ce1\}
\end{array}$$

Figure 4.5: The semantics of ZHELANG. We use s to denote the input stream, and c to denote combinations of observers.

and `orOp` receive three arguments. The first two are observers, and the third is the input stream.

The `atomOp` observer is the simplest of the five: it matches a token, its first argument, with the head of the input stream. There are two operators that combine observers:

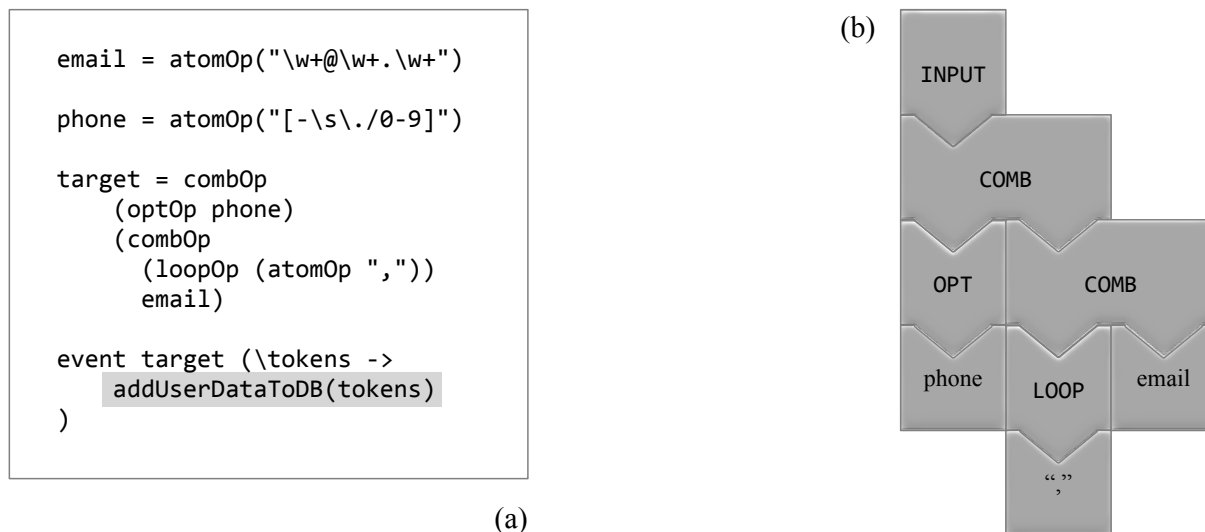


Figure 4.6: (a) A ZHELANG program that recognizes a phone and an email separated by a comma. (b) The block representation of the program.

`combOp` and `orOp`. The former applies observers in sequence: it reports an event when both its input observers return successful events. The latter, `orOp` reports an event in case any of its input observers does it.

The `searchUntilOp` operator is used to implement iterations over the input stream. The construction `searchUntilOp c1 c2 s` applies the observer `c1` on the stream `s` until `c1` triggers an event for some token in `s` or until `c2` is successfully applied to `s`. Finally, the `optOp` operator is present in ZHELANG to allow the existence of combinators that do not modify the input stream.

Example 12 Figure 4.6 shows a program that extracts a phone number from the input stream, in case one exists, plus an email. These tokens must have a comma between them, as in a CSV file or as in a semi-structured address format. We use the `optOp` operator to indicate that the phone number is optional. Thus, the program is able to recognize the regular expression formed by zero or one occurrences of a phone number, plus a comma and then an email. The last two tokens are mandatory.

Figure 4.5 uses a special atom, `Fail`. This atom denotes a pattern that never matches any other pattern. `Fail` is defined by Rule *Fl1* and is used in rules *Eg1* and *Ce1* to indicate that events are searched *ad infinitum*. While writing programs in ZHELANG, we have found this pattern—consuming tokens until an event takes place—so common, that the current language specification defines an operator for it. Figure 4.7 shows how this new operator, `searchOp`, is defined. Example 13 explains its behavior.

Example 13 The program in Figure 4.7 (b) searches for an SQL query containing a `SELECT` block plus a `WHERE` clause. Once search starts, the guard will be deactivated when a line breaker is read. Figure 4.7 (c) shows a similar program; however, once this event becomes active, it will only stop running if the query is found.

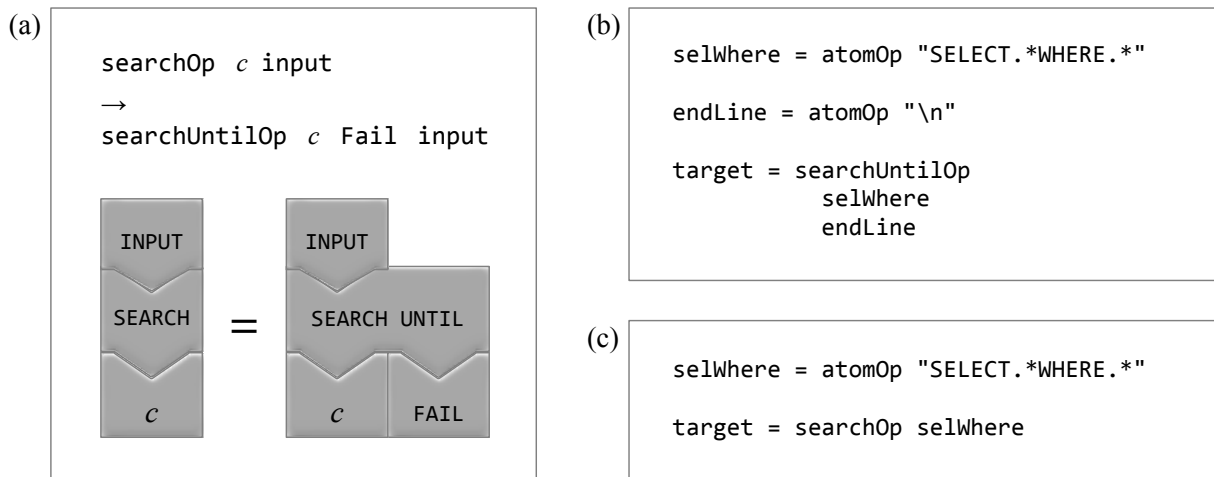


Figure 4.7: (a) The definition of the operator `searchOp`. (b) Program that searches for an SQL query until a line break is found. (c) Program that searches for the same SQL query, until one is found.

$$\begin{aligned}
O(\text{Fail}) &\approx O(1) \\
O(\text{atomOp } t) &\approx O(t) \\
O(\text{optOp } c) &\approx O(c) \\
O(\text{combOp } c_1 \ c_2) &\approx \max(O(c_1), O(c_2)) \\
O(\text{orOp } c_1 \ c_2) &\approx O(c_1) + O(c_2) \\
O(\text{searchUntilOp } c_1 \ c_2) &\approx T * (O(c_2) + O(c_1))
\end{aligned}$$

Figure 4.8: The asymptotic complexity of the different operators of ZHELANG. We let n be proportional to the size of the input stream observed until when the event takes place.

4.1.4 Asymptotic Complexity

Figure 4.5 gives us the means to estimate the computational complexity of each one of the six operators used in ZHELANG. This complexity is given in Figure 4.8. The estimates provided in Figure 4.8 can be inferred directly from the semantics of each operator. This process is used in the proof of Theorem 8.

Theorem 8 (Complexity) *The asymptotic complexity of ZHELANG observers is given by the relations in Figure 4.8.*

The proof is by case analysis on the rules in Figure 4.5. We show a few cases:

At1: The cost of solving `atomOp t (h : s)` is the cost of matching the regular expression t with the token h .

At2: Same as the case for At1. Therefore, that limits the complexity of `atomOp t (h : s)` into the cost of performing one pattern matching.

Co1: The cost of solving `combOp c1 c2 input` is the cost of running `c1 input`.

Co2: The cost of solving `combOp` c_1 c_2 `input` is the cost of running c_1 `input` and c_2 s , where s is also an input stream. Therefore, the complexity of this case is the sum of the complexity of the subcases.

Co3: Same as the case for Co2.

Su1: The cost of solving `searchUntilOp` c_1 c_2 `input` is the cost of running c_1 `input`.

Su2: The cost of solving `searchUntilOp` c_1 c_2 `input` is the cost of running c_1 `input` plus the cost of running c_2 `input`.

Su3: The cost of solving `searchUntilOp` c_1 c_2 `input` is as in the case Co2, plus the recurrent cost of invoking `searchUntilOp` on the smaller input, which can be reduced by one token. This case can be repeated n times, where n is the number of tokens inspected until either case Su2 or Su1 happen.

The other cases are similar, and we shall omit them. Notice that the only difference between `opOp` and `combOp`, in what regards their computational complexity, is the fact that the second observer in `combOp` runs only when the first observer fails. Therefore, we replace the summation by the upper bound (max) on the asymptotic costs of these operators. \square

Most of the programs written in ZHELANG will run in time linear on the number of tokens observed times the cost of matching these tokens against regular expressions. This is the case, for instance, of the three case studies in Section 4.3. Linear complexity is the most common behavior, because to go beyond this bound, observers must contain either `orOp` or recursive combinations involving `searchUntilOp`. Nevertheless, the relations in Figure 4.8 indicate that it is possible to derive programs with exponential complexity using a combination of `orOp` and `searchUntilOp` observers. This result holds even when the input stream is bounded. Example 14 illustrates this possibility.

Example 14 *Figure 4.9 shows four different programs, and recurrence relations defining their asymptotic complexities. The solution for each one of these four recurrence relations appears in Figure 4.9 (c). Notice that the last observer, `expnLoop`, has exponential complexity.*

4.2 Implementation

The reference implementation of ZHELANG is written in Kotlin and runs on the Java Virtual Machine. This section presents the reader with some details of this implementation.

<pre>find1Dot = searchOp (atomOp ".")</pre>	$T_{\text{find1Dot}}(n) = O(\text{atomOp "."}) + T_{\text{findDot}}(n-1)$	$O(\text{find1Dot}) = O(n)$
<pre>find2Dot = orOp find1Dot find1Dot</pre>	$T_{\text{find2Dot}}(n) = O(\text{find1Dot}) + O(\text{find1Dot})$	$O(\text{find2Dot}) = O(n)$
<pre>quadLoop = searchOp find2Dot</pre>	$T_{\text{quadLoop}}(n) = O(\text{find2Dot}) + T_{\text{quadLoop}}(n-1)$	$O(\text{quadLoop}) = O(n^2)$
<pre>expnLoop = searchOp compOp (atomOp ".") orOp expnLoop expnLoop</pre>	$T_{\text{expnLoop}}(n) = O(\text{atomOp "."})$ $\quad + T_{\text{expnLoop}}(n-1)$ $\quad + T_{\text{expnLoop}}(n-1)$	$O(\text{expnLoop}) = O(2^n)$

(a)
(b)
(c)

Figure 4.9: (a) Examples of ZHELANG programs. (b) The recurrence relation that denotes the complexity of each program. (c) The solution of each recurrence relation, giving the asymptotic complexity of the examples.

4.2.1 Overview of the Implementation

The implementation of ZHELANG follows closely the observer design pattern. Thus, we have a number of objects representing *observers* registered onto a central notifier, the *subject*. The latter encodes a text stream and a panel onto which observers can be plugged. Therefore, the runtime environment of a ZHELANG program is composed by the components listed below. Figure 4.10 shows the relation between them.

- *Observers*: objects that encode any of the six operators described in Section 4.1;
- ZHESTREAM: a *facade* representing the infinite text that must be processed by the program;
- ZHEPANEL: the subject of the observer pattern, that is, the object in which observers are registered to wait for events.

4.2.1.1 The Implementation of Observers

Structurally, observers are *parser combinators* [?]. In other words, an observer is a high-order function that takes parsers (other observers) as input and returns a parser (it-

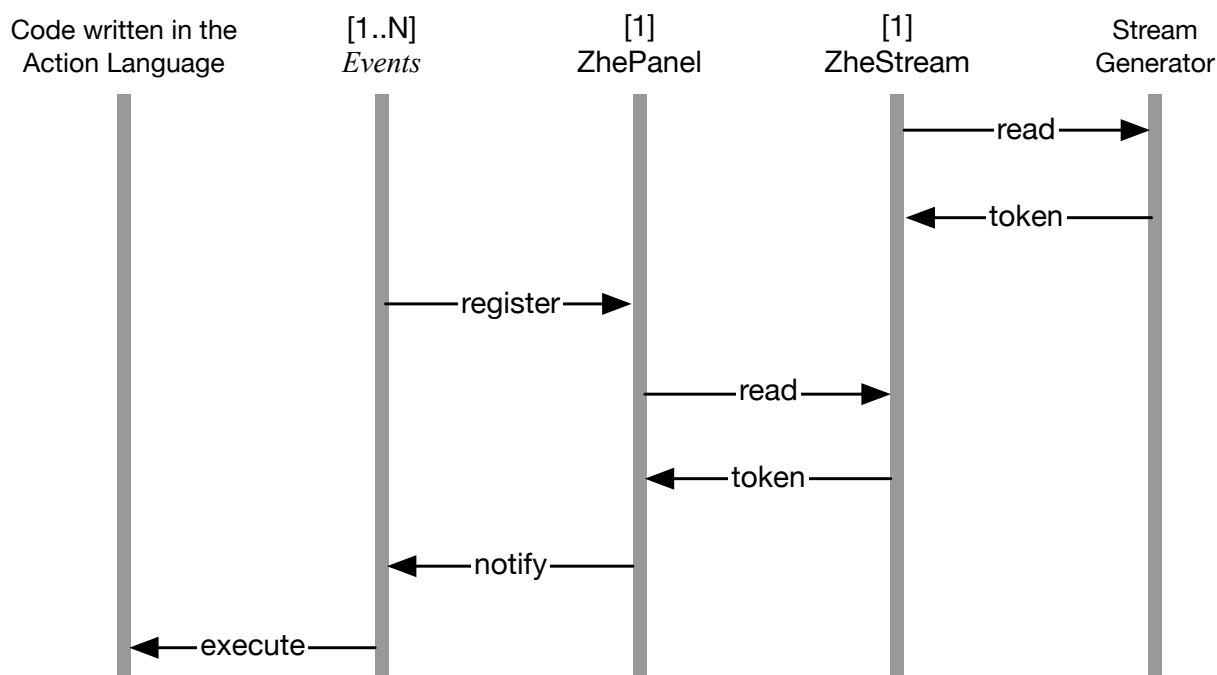


Figure 4.10: The relation between the core components of ZHELANG. The stream generator and the code that implements actions are not written in ZHELANG.

self). These parsers read tokens provided by the ZHEPANEL and output a parser tree. An observer has the following type: $\text{String} \rightarrow (\text{Boolean}, [\text{String}], \text{String}, \text{Maybe Error})$. If o is a ZHELANG observer, and $(status, tokens, rest, error)$ its output, then we have the following definitions:

status: This boolean tells if the matching was successful or not. This information needs to be explicitly given, rather than otherwise inferred, because the `optOp` operator, and any observer based on it, might be successful but recognize no tokens.

tokens: A list with the recognized tokens. Items are added to the list by a successful application of the `atomOp` operator, or by combining the output of other observers. For example, the `combOp` operator concatenates the lists of two successful observers and the `orOp` operator returns the list of tokens of the first successful observer.

rest: The third element of the tuple is the remaining input string to be parsed after the application of the observer. When the observer fails, it consumes the first token from the input, the only exception being again `optOp`, which returns the input unchanged if it fails. When the observer succeeds, the remaining string will be the result of removing the recognized tokens from the input string.

error: The optional error indicator is a tuple containing a number and a string. The number is an error code and the string is an error message explaining why the observer failed. An observer might fail for different reasons. A non-exhaustive list includes:

1. `AtomOp` fails if the input string does not match the provided token;
2. `SearchUntilOp` fails if the termination condition was matched;
3. `combOp` fails if any of the provided observers fail.

4.2.1.2 The ZHSTREAM

We use the expression “infinite string” throughout this text. However, this is rather an abuse of language. In practice, an infinite string is represented by a generator. Therefore, an infinite string is an unbounded set of finite strings produced at different moments. The ZHEPANEL receives and records each element of this set. These elements are received from a stream abstraction, the ZHSTREAM.

ZHSTREAM is a Stream of text that keeps track of the string chunks arriving at different time intervals. This abstraction breaks the incoming strings into chunks of the same size. Said chunks are broadcasted to any registered observer by the ZHEPANEL. ZHSTREAM splits strings into equal-sized pieces for performance reasons. ZHELANG’s implementation uses Java `StringBuilder` to manipulate strings. Tokenization of an immutable string is linear on the size of the string. Thus, splitting happens to avoid having to copy large chunks of data to separate them into tokens.

However, splitting can compromise ZHELANG’s correctness. An Error may happen if an observer requires more than one chunk of string to match its pattern. To deal with this problem, the ZHEPANEL relies on the error returned by the observer. If the observer fails because it consumed all the string chunks, then the ZHEPANEL asks ZHSTREAM to provide another chunk. This new piece is concatenated to the current chunk and matched against the observer once again. This process loops until the match is successful or if it fails due to a reason other than the length of the current buffer. Notice that this modus operandi allows the creation of observers that never terminate, as Example 15 illustrates.

Example 15 Consider the event $Ob = (\text{atomOp } ". * ")$. Ob will consume all the string chunks received thus far until failure ensues. If every chunk has N characters and M chunks have been observed, then up to that point, Ob will have performed $N * (\frac{M^2}{2} + \frac{M}{2})$ matching operations.

4.2.1.3 The ZHEPANEL

This data structure works as a bridge between data (the input stream) and the events. It handles the following responsibilities:

- Tracks which parts of the input string were matched against each observer;
- Binds observers to actions, thus creating events;
- Activates events, spawning them into new threads;
- Broadcast string chunks to observers.

In terms of implementation, the ZHEPANEL contains a method `addEvent(Observer, Action)`. The first argument is defined in Section 4.2.1.1. The second is an interface with the action language. This interface is implemented in Kotlin, and links with a library with hooks to invoke external processes. The implementation of ZHELANG discussed in Section 4.2 invokes actions implemented in either Java or Kotlin. Therefore, interfacing is simple: actions are implemented as classes compiled into JVM bytecodes. The `Action` interface contains a method `execute([String])`, which receives a list with the tokens matched by the enclosing observer.

4.2.2 Core Property: Isolation

Each event in ZHELANG executes within a separate thread. Actions run as separate processes. The observer does not block other actions while executing it's own. Thus, multiple actions might run concurrently, having been created by the same observer, albeit to handle different string events. Events request string chunks independently from each other, by invoking the proper method upon the ZHEPANEL.

Events run in isolation, meaning that they do not share state. The ZHESTREAM is only read by the events; therefore, it is not replicated. However, different events, given the nature of the observers that guard them, can read different parts of the input stream. We emphasize that isolation is maintained for correctness, not performance reasons, as Example 16 illustrates.

Example 16 *Figure 4.11 shows two events. The first waits for the pattern "AB"; the second, for the event "BC". If these events do not run independently, the string "ABC"*

could be matched in two different ways. In both, an event would be missed. For instance, if the program matches *Ob1* before *Ob2*, then it will successfully match *Ob1* consuming tokens "AB". The string "C" will be returned to the shared state. After that, it will try to match *Ob2*. Failure would ensue, as part of the necessary input had already been consumed. This semantics would be in disaccord with the behavior expected in Figure 4.5

<pre>Ob1 = combOp (atomOp "A") (atomOp "B")</pre>	<pre>Ob2 = combOp (atomOp "B") (atomOp "C")</pre>	<pre>event Ob1 (\tokens -> callAction(tokens))</pre>	<pre>event Ob2 (\tokens -> callAction(tokens))</pre>
-------------------------------------------------------	-------------------------------------------------------	------------------------------------------------------------	------------------------------------------------------------

Figure 4.11: Program that demonstrates the need of isolation for correctness.

4.3 Use Cases

This section presents three use cases of ZHELANG. Section 4.3.1 shows how to build a log obfuscator. Section 4.3.2, in turn, presents a crawler built on top of our DSL. Finally, Section 4.3.3 shows how to run simple refactorings using ZHELANG. The goals of these use cases are twofold. First, we want to provide the reader with examples of useful programs that can be easily expressed in ZHELANG. Second, we want to compare the efficiency of the implementation discussed in Section 4.2 with other, better-known systems.

Hardware Experiments discussed in this section were performed on an Intel Core i5 with a clock of 1.8GHz, and 8 GB of RAM (DDR3) at 1,600MHz.

Software ZHELANG was implemented in Kotlin release 1.4.20. For comparison purposes, we have also implemented every case study in Comby [?], which is publicly available at <https://github.com/comby-tools/comby>. Comby is a parser-parser combinator, that factors out commonalities among different programming languages into a single abstract syntax tree. Additionally, Sections 4.3.1 and 4.3.2 compare ZHELANG with tools that are specific to those case studies. In Section 4.3.1, this tool is ZHEFUSCATOR, a log redaction system used by Cyral Inc., a data security company (<https://github.com/lac-dcc/Zhe>). Section 4.3.2 uses BEAUTIFULSOUP4, version 4.9.3. This is a python library for the analysis of HTML files. Every experiment runs on OSX 64-bits, version 10.14.16.

Methodology When reporting running times, we show the averages of five executions. Samples for the same program are not collected in sequence. Instead, programs run in random order. Differences in running time between averages collected for distinct systems are considered statistically significant within a confidence interval of 99%. All the applications employed in this evaluation have a timeout of 600s. The timeout was necessary to ensure a fair comparison between tools, as it is hardcoded in Comby. Incidentally, only case studies implemented in this tool have reached the time barrier during our evaluation.

4.3.1 Log Obfuscation

Log obfuscation consists in the identification and redaction of sensitive sentences produced by log generators. A core motivation for this kind of process are recent data protection laws [?, ?, ?], such as the *General Data Protection Regulation (GDPR)*¹, valid in the European Economic Area since 2016. The *log obfuscator* is an online system that intercepts the output of a *log generator*, e.g., a database server, an operating system, etc, and recognizes string patterns denoting classified information. Once such patterns are identified, it replaces them with encrypted text.

Experimental Setup. In this section, we compare log obfuscators implemented with ZHELANG and with Comby against a tool that has been built with this purpose: Cyral Inc.’s ZHEFUSCATOR. ZHEFUSCATOR is implemented in Java, and runs as an extension of the Java Virtual Machine, analyzing any output recorded in file descriptors. In this experiment, we have replaced ZHEFUSCATOR’s grammar synthesis tool with programs written in COMBY and ZHELANG.

In this experiment, we have generated logs from MySQL SQL Database version 14.14 Distribution 5.7.27. Workloads were produced by the nine real-world web applications emulated by OLTP-Bench [?]. OLTP uses data from Wikipedia, Twitter and other large databases. The three competing approaches employed in this experiment, ZHELANG, COMBY and ZHEFUSCATOR were configured to redact every literal in every SELECT query in the log. Figure 4.12 shows the ZHELANG program that performs this action. The three systems use the same action, namely, routine `redactQuery`, which has been implemented in Kotlin, and compiled to Java bytecodes.

Discussion. Figures 4.13, 4.14 and 4.15 show the result of the comparison between the running time of the three tools considered in this section. OLTP gives us nine workloads.

¹<https://eugdpr.org/>

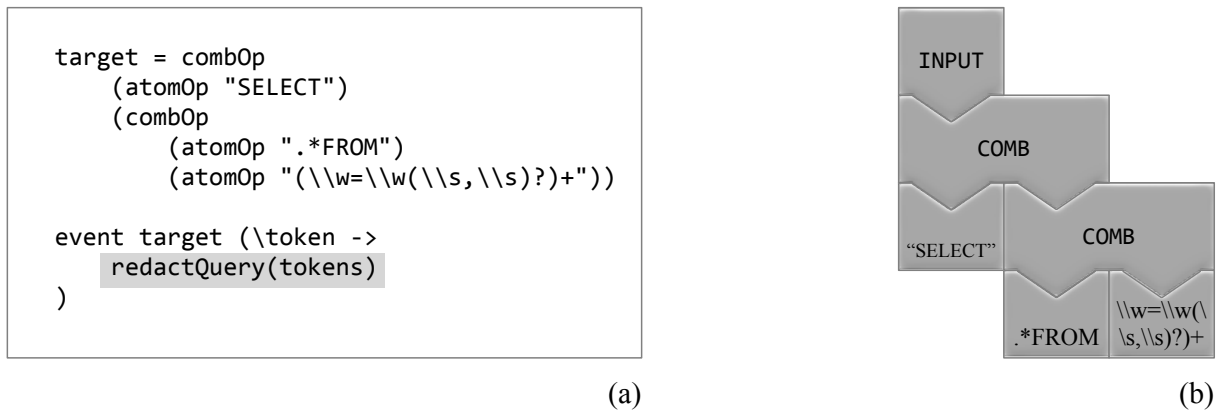


Figure 4.12: (a) A ZHELANG program that recognizes **SELECT** queries in SQL programs. (b) The corresponding block representation of the program.

We have run each tool on each workload five times. Figure 4.13 contains one point for each one of these experiments (45 points per tool), except when considering COMBY. This parser-parser combinator timed out in the majority of the experiments. The next figure, 4.14, organizes these samples into boxplots, per workload, so that the reader can check averages and outliers. Finally, Figure 4.15 summarizes this experiment presenting boxplots per tool, for all the workloads together.

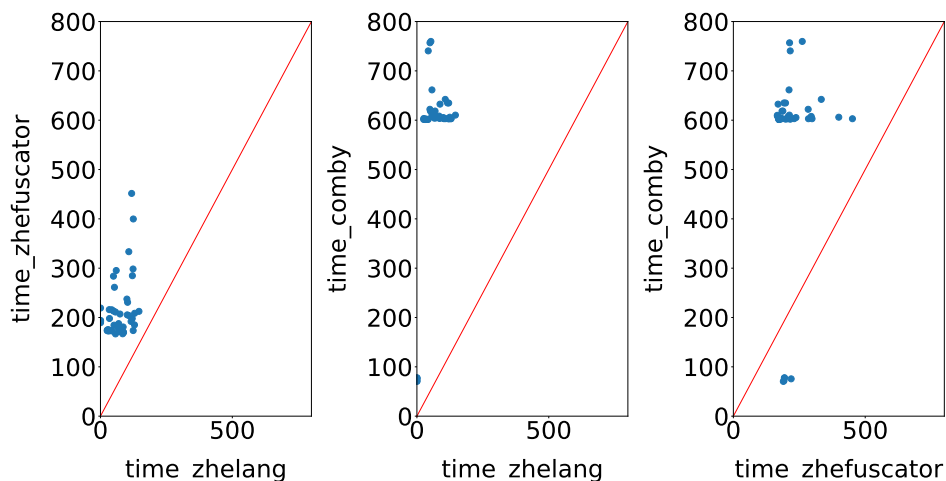


Figure 4.13: Time comparison between different log redaction tools. Each point represents a workload. Axes show time, in seconds. The red line is the line that crosses the origin and point (1.0, 1.0). The less points a tool has within its side of the diagonal, the faster it is.

ZHELANG is the fastest player in this evaluation. However, the reader must understand that these tools do different things; hence, they cannot be compared directly. Programmers must encode the same pattern of interest (SQL queries) in these three technologies. What distinguishes them is how they identify this pattern in the input stream, which contains tokens that are not of any interest. ZHEFUSATOR synthesizes a grammar for an unknown language—in this case, the language that recognizes the log stream.

Therefore, it requires a few examples until converging to a definitive grammar. In this experiment, ZHEFUSCATOR converged after eight to nine examples. The ZHELANG program, in turn, does not need to synthesize a grammar to recognize the full log language. It discards any token that cannot be used to build the pattern of interest, i.e., `SELECT` queries. COMBY also does not need to recognize the full log language; however, it expects the input string to be encodable in a Dyck-Extended language. Such is not the case of SQL, and COMBY falls back into a generic parser that performs declarative matching on the input—a rather slow process.

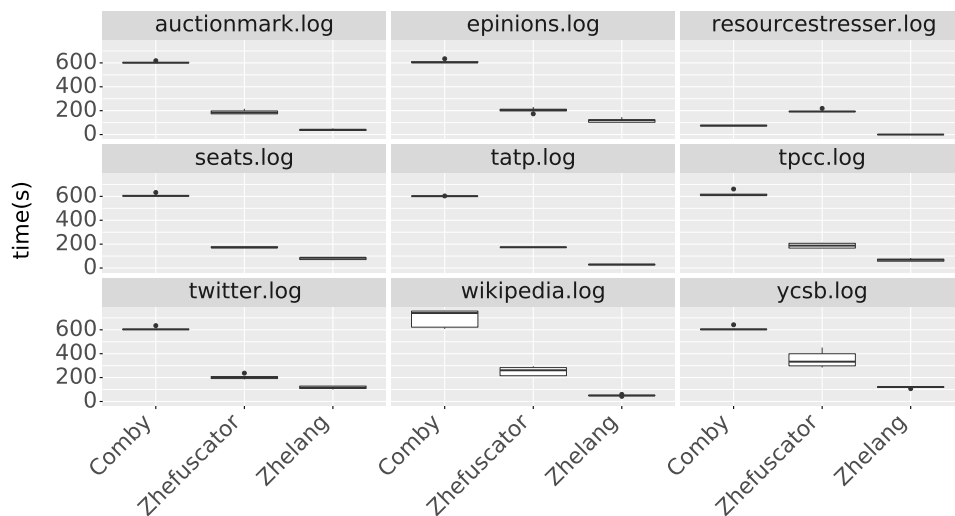


Figure 4.14: Time distribution for ZHELANG, ZHEFUSCATOR and COMBY, when analyzing the OLTP logs.

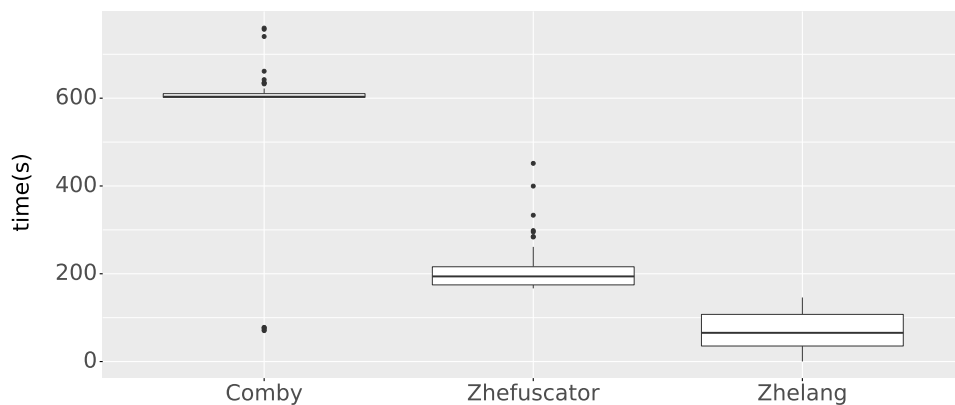


Figure 4.15: Summary of the running times observed in Figure 4.14.

4.3.2 Web Crawling

Web crawling is a scenario that can be naturally modeled as a problem of processing an infinite stream of text. The set of webpages that form the World-Wide Web is, obviously, finite. However, for any practical purpose, it can be treated as unbounded, given that its sheer volume prevents typical string processing tools from storing its contents before textual analysis starts. In this experiment, we use ZHELANG to build the parsing component of a web crawler. This crawler builds a catalogue with the links found in webpages.

Experimental Setup. The workload for this experiment consists in the 50 pages listed in *Alexa Top Sites*². The goal of this evaluation is to measure the time to parse these pages. Thus, to avoid timing fluctuations due to the network latency, we have downloaded all the pages using GNU `wget` Version 1.21. In this experiment, we compare ZHELANG with BEAUTIFULSOUP4's parser, which is implemented in Python, and is customized for parsing webpages. We have also implemented this parser in COMBY. The goal of this experiment is to extract every link tag from HTML files and rewrite them into a JSON format. Figure 4.16 shows the ZHELANG implementation of this task.

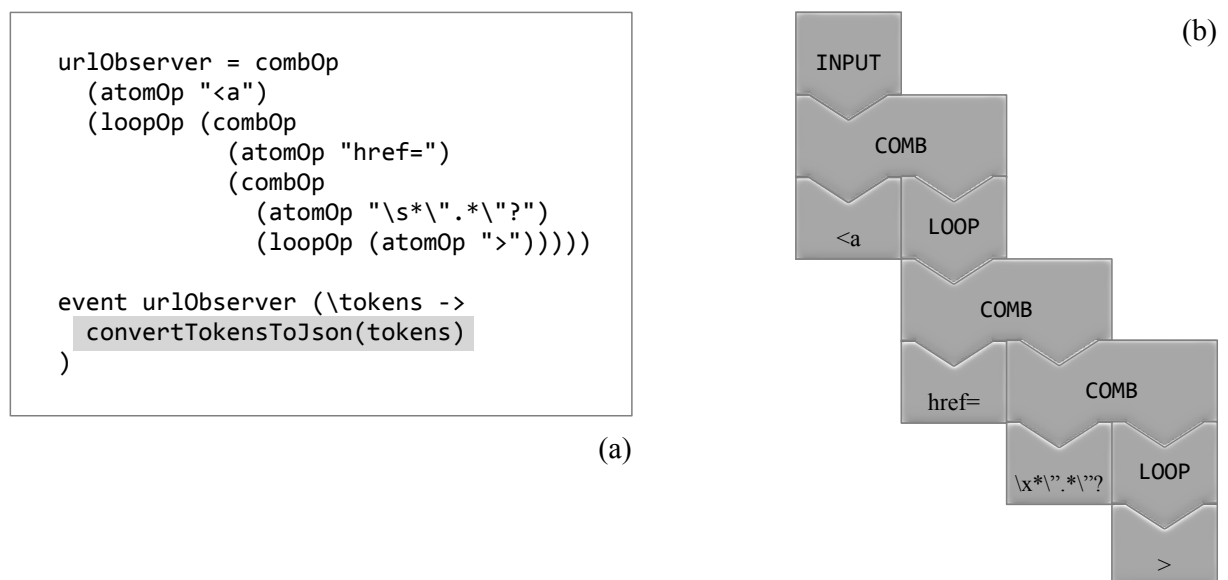


Figure 4.16: (a) A ZHELANG program that finds links in webpages. (b) The block representation of the program.

Discussion. The parser built using ZHELANG outperforms the parsers built with COMBY and with BEAUTIFULSOUP4. In the latter case, the difference is small, albeit non-

²Available at <https://www.alexa.com/topsites> on February 1st, 2021

negligible. Figure 4.17 shows absolute running times for every experiment (5×50 samples). This difference in running times is uniform across the 50 webpages that constitute our workload. Figure 4.18 provides absolute times per webpage. Each bar in Figure 4.18 is the average of five runs. That figures makes it clear that COMBY’s poor performance is not due to outliers. Our experience shows that “declarative matching”, the technique used by this tool, although flexible, yields slow parsers when ported to environments that lay outside the domain of Dyck-Extended languages.

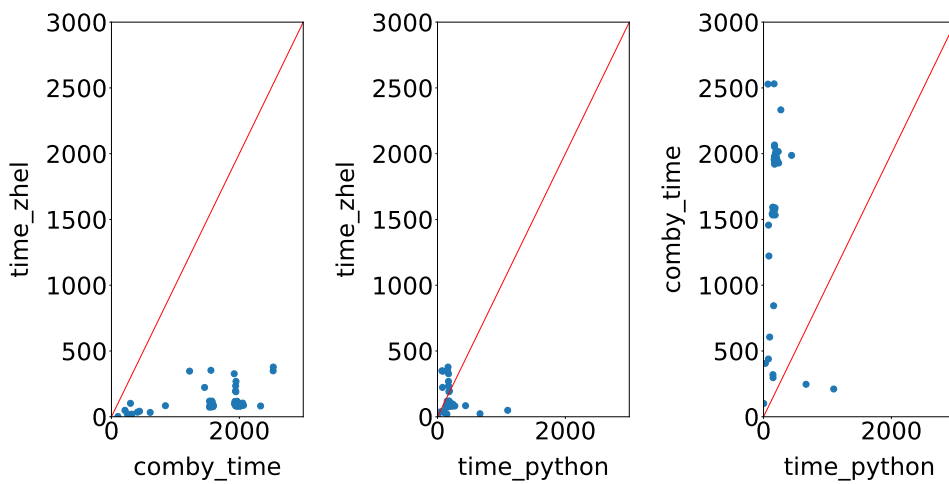


Figure 4.17: Time comparison between different web crawlers. Each point represents a workload (one of five executions for each one of fifty webpages). We show the main diagonal, to ease comparison. The less points a tool has on its side of the diagonal, the faster it is.

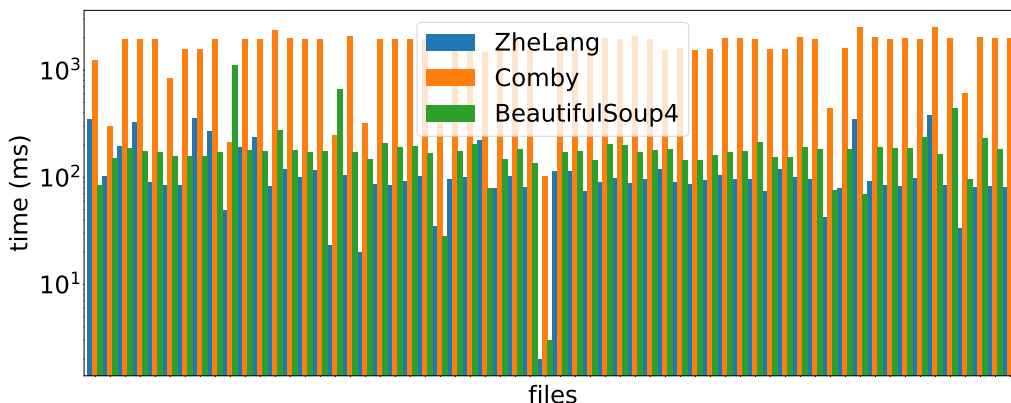


Figure 4.18: Running time necessary to process each webpage in the workload. Each tick on the X-axis represents a different webpage. Bars show averages of five samples.

The running time difference observed between ZHELANG and BEAUTIFULSOUP4 is an artifact of the underlying programming language used to implement these two systems. ZHELANG is implemented in Kotlin, and runs on the JVM. The BEAUTIFULSOUP4

parser is implemented in Python, and runs on CPython. Visual inspection of both programs shows that the number of parsing actions (reductions and shifts) necessary in either implementation is approximately the same. In this experiment, the bulk of the running time of each web crawler is spent on text processing. The action triggered once a link is found takes minimal time. This action consists in saving the link into persistent storage. We shall provide data to back up this observation in Section 4.3.4 (Figure 4.24, Page 65).

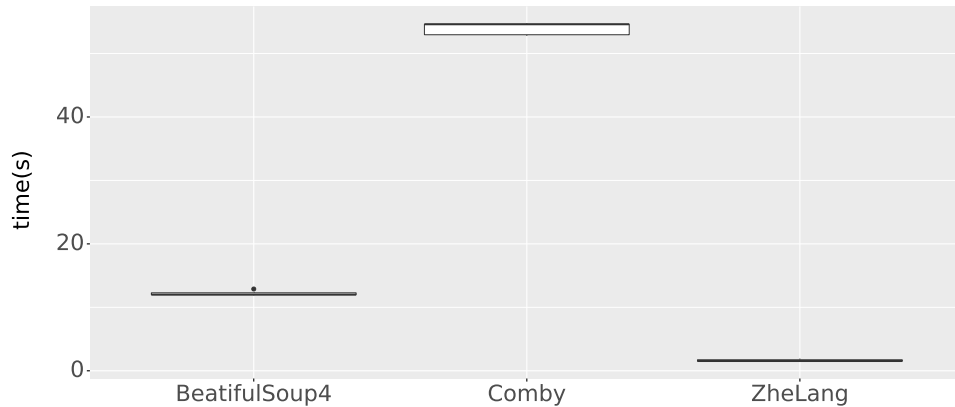


Figure 4.19: Summary of running times of the different web crawlers built using either ZHELANG, COMBY or BEAUTIFULSOUP4.

4.3.3 Program Refactoring

Refactoring is a family of automatic program transformations that developers use to improve the quality of the code that they write. Refactoring is not the main domain for which ZHELANG has been conceived: the program text is finite. Nevertheless, we use this case study to compare ZHELANG with COMBY in the scenario for which the latter has been designed.

Experimental Setup. We have written a refactoring tool for C, that transforms “for” loops into “while” loops. In terms of parsing, this is a trivial transformation: it consists in recognizing terms in the format “for(*; * ;*)”, where * is any number of characteres. Figure 4.20 shows the implementation of this tool in ZHELANG.

As benchmarks, we have used the 10,000 largest C programs available in the ANGHABENCH collection of benchmarks [?]³. We chose the “single-function” category of benchmarks, meaning that each C file contains one single function, plus all the definitions (types, external functions, macros) necessary to ensure its compilation.

³Available at <http://cuda.dcc.ufmg.br/angha/home> on February 1st.

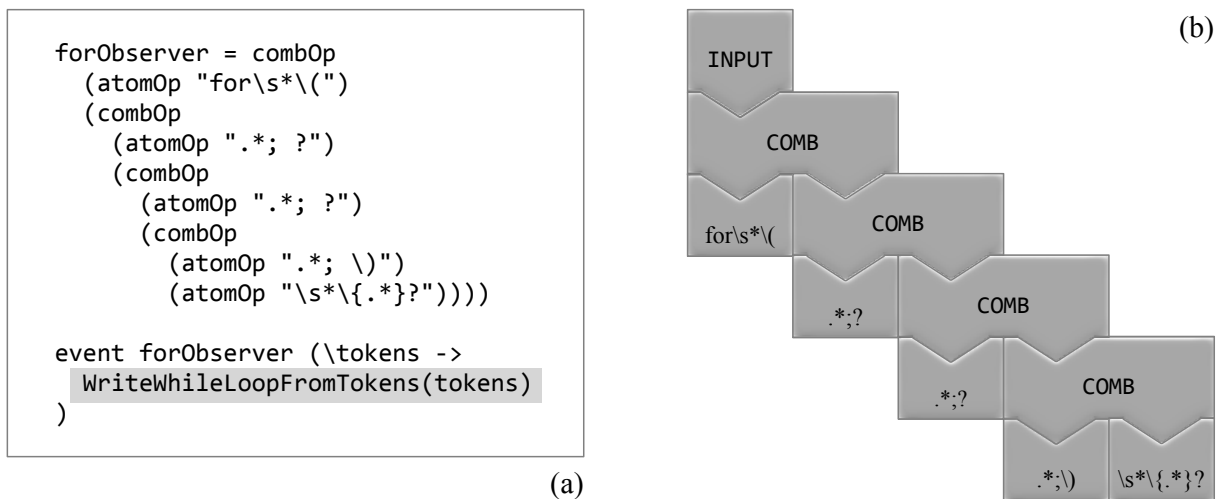


Figure 4.20: (a) A ZHELANG program that recognizes `for` loops in programs. (b) The block representation of the program.

Discussion Figures 4.21 and 4.22 indicate that the ZHELANG program is faster than its counterpart implemented in COMBY. On the average, refactorization based on ZHELANG takes 35 milliseconds per file. Using COMBY, this time jumps up to 68 milliseconds on average. To give the reader some perspective on these numbers, notice that the C files have 7853 characters on average. Furthermore, the output produced by both parsers is identical.

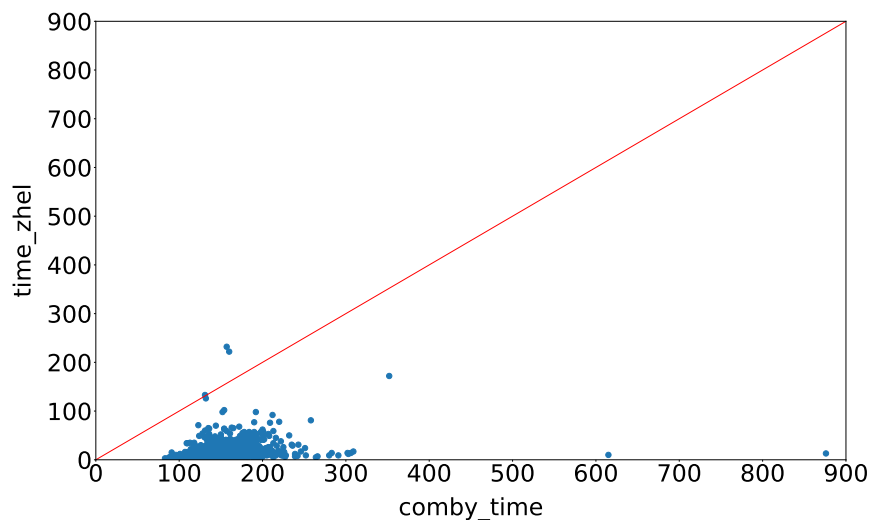


Figure 4.21: Time comparison between ZHELANG and COMBY with respect to the task of transforming “`for`” loops into “`while`” loops. Each point represents a workload (5×10^4) samples. The red line is the main diagonal. The less points a tool has on its side of this line, the faster it is.

Figure 4.23 plots the time that ZHELANG and COMBY take to process the 1024 largest files in our benchmark collection. The ZHELANG implementation has outperformed COMBY’s in every one of the 1024 samples in Figure 4.23. The maximum time

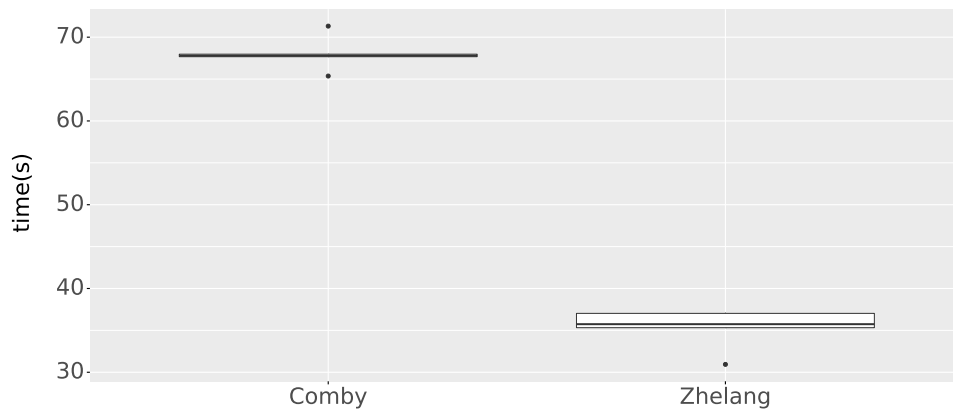


Figure 4.22: Summary of results presented in Figure 4.21.

taken by ZHELANG was 232 seconds, for a program with 6211 characters. The maximum time taken by COMBY was 876 seconds. In this case, the input file has 11166 characters. Nevertheless, there are programs, not among the 1024 largest files, in which COMBY's implementation has been able to outperform ZHELANG. Figure 4.21 shows two occurrences (out of 10,000 experiments). Because the running time of both implementations is so short, we believe that this result is caused more by standard fluctuations in time measurement than to some inability of ZHELANG's implementation to handle particular inputs.

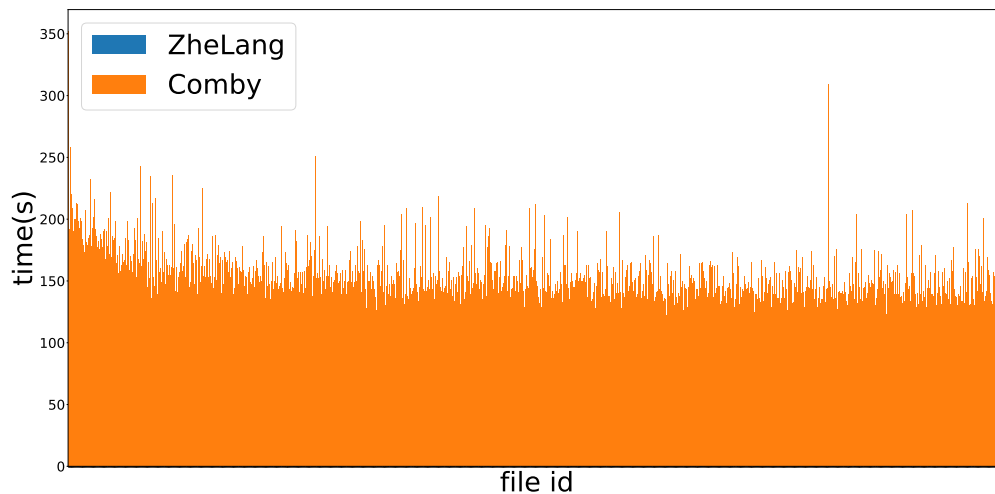


Figure 4.23: Running times collected for the 128 largest files in the ANGHABENCH collection of benchmarks. Each bar is the arithmetic average of five executions. Programs are sorted by size, measured as Kilobytes.

4.3.4 Analysis of the Runtime Behavior of ZHELANG

The running time of ZHELANG programs can roughly be divided into the three parts below:

Tokenization: the separation of the visible part of the input stream into tokens. The implementation of ZHELANG discussed in Section 4.2 uses a very conservative definition of tokens: any single character read in the input is a token.

Matching: the verification of which atoms can be recognized in the input. To this end, sequence of tokens are given to the `atomOp` combinators in the event guard. Each one of these operators define a string pattern, i.e., a regular expression. Once a sequence of tokens is successfully matched with a regular expression, a string event takes place.

Actuation: the execution of an action, once an event is detected by an observer.

In this section we analyze how much time our different case studies spend on each one of these parts. However, before we start this analysis, notice that the time spent on actuation is not directly dependent on metric properties of the input stream (the quantity and average length of tokens observed). In other words, actuation consists in the invocation of an external process to handle the occurrence of an event—a pattern of text on the input stream. The case studies in Sections 4.3.1–4.3.3 invoke very short actions, for our goal when evaluating them was to compare the speed of different parsing techniques. Therefore, the bulk of time reported in this section is spent in either the tokenization or in the matching phases of a ZHELANG program.

Methodology In this section, we profile the time taken by two of our case studies: the one in Section 4.3.2 and the one in Section 4.3.3. These two case studies have the benefit of receiving individual inputs (web pages or C files), which can be profiled independently from each other.

Discussion Figure 4.24 shows how the running time of the program in Figure 4.16 is split into tokenization, matching and actuation. Figure 4.25 provides similar information for the program in Figure 4.20. Time spent on parsing is proportionally higher in Figure 4.25 than in Figure 4.24. One could believe, in principle, that this difference happens because in the latter case, texts manipulated by ZHELANG are larger. This intuition is motivated by the observation that time spent on tokenization varies approximately linearly with program size, whereas parsing has higher complexity. However, this assumption is false. Each C file contains, on average, 7853 tokens, whereas each webpage

contains 281621. Profiling reveals that the culprit for the difference between Figures 4.24 and 4.25 is the time necessary to evaluate regular expressions. The parser used for the code refactoring tool (Figure 4.20), contains more regex operations than the parser used for web crawling (Figure 4.16). As already observed in Figure 4.8, the computational cost of running a regex-based guard is given by the cost of checking the regular expression against every token in the input stream.

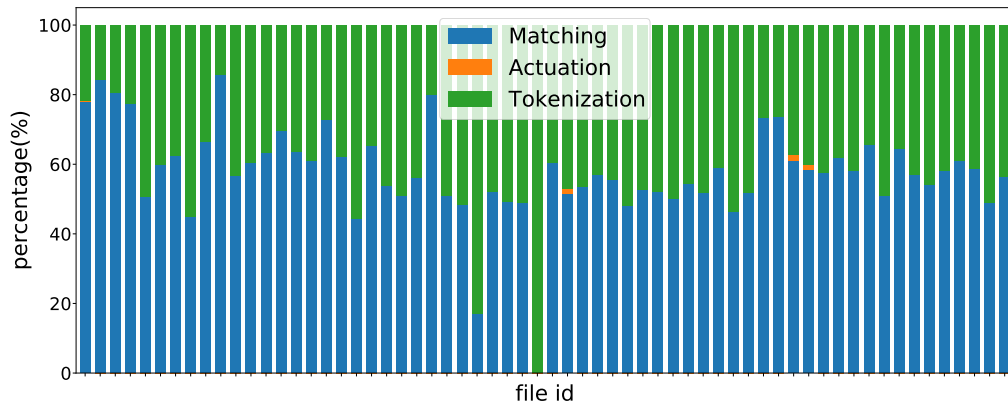


Figure 4.24: Division of the time that the program in Figure 4.16 takes to process the text of different webpages. Each tick on the X-axis represents one of the webpages analyzed in Section 4.3.2.

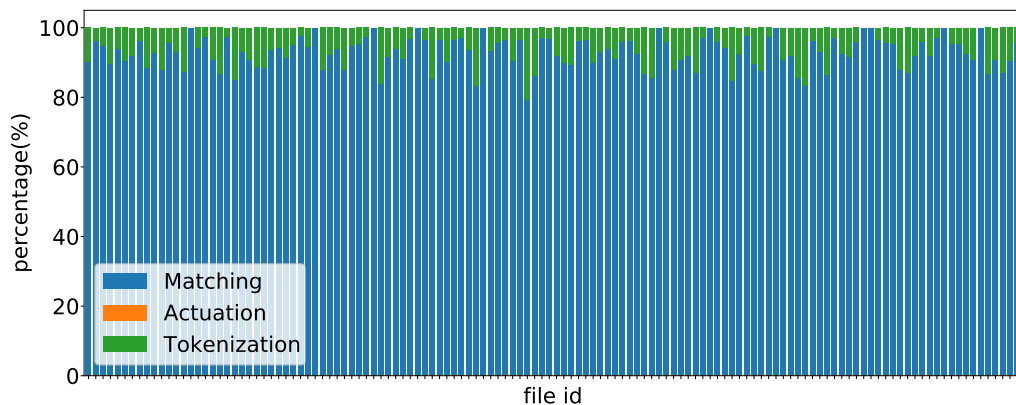


Figure 4.25: Time that the program in Figure 4.20 spends on different parts of its implementation. Notice that the time spent on actuation is present, although not visible, given that it is short, compared to the other running times. Each tick on the X-axis represents one of the 128 largest programs analyzed in Section 4.3.3.

Chapter 5

Literature Review

The treatment and processing of strings has been extensively discussed in the literature. In Section 5.1 we discuss the foundations and general concepts of parsing, providing some perspective on the context in which this work is inserted. Section 5.2 discusses this theory, to give the reader some perspective on the foundations of the present work. We also notice that much of the developments in this work bear resemblances to programming fuzzing. Yet, whereas fuzzing is concerned with recognizing a language that describes the input of a program, this dissertation deals with the inverse problem: we recognize a language that describes the output of the program. Section 5.3 discusses work related to fuzzing. Additionally, there exists a vast body of literature concerned with the synthesis of grammars from examples. This is the approach that we use in Section 3.2.1 to equip the ZHEFUSCATOR with a user interface. In Section 5.4 we discuss work related to the synthesis of grammars from examples. Finally, as our theory, once implemented into an actual tool, yields a reactive system, we cover those in Section 5.5.

5.1 Parsing

ZHELANG solves a parsing problem. However, contrary to traditional parsers, a ZHELANG program defines a grammar that does not necessarily recognize the whole of the input it receives. Instead, ZHELANG grammars recognize patterns within this input. In this sense, programs written in ZHELANG solve a *language separation problem* [?], i.e., how to recognize a target language within a host language. However, in our context, we care only for the target language: tokens in the host language are skipped through the `searchUntilOp` observer. Carrying this analogy further, we could say that `searchUntilOp` observer “recognizes” the host language, whereas the other operators recognize the target language. Notice that this way of looking into the language separation problem differs from its original conception, where grammars must be synthesized for the host and for the target languages.

The `orOp` observer is analogous to the *choice* operator used in Parsing Expression Grammars (PEGs) [?]. Therefore, except for the `searchUntilOp` observer, the other operators used in ZHELANG could have been defined as PEG operations. Nevertheless, PEGs and ZHELANG are concepts that fit into different categories: the former is a formalism to define top-down parsers; whereas the latter is a domain specific programming language. Although indeed many tools have been built as parsing expression grammars [?, ?, ?], we found it more natural to define ZHELANG’s observers as parser combinators [?].

There is much literature around the notion of parser combinators. For an overview, we recommend the Related Work sections of `okasaki1998functional` or `brown2016build`. Nevertheless, we emphasize that while ZHELANG programs can be seen as ensembles of parser combinators, its reactive nature does not appear in tandem with said literature—inasmuch as our knowledge goes.

Still in what parsing is concerned, while carrying out our experiments, we realized that several ZHELANG programs could be easily reimplemented as instances of the COMBY framework [?], as long as the input to be parsed was finite. Nevertheless, notice that whereas COMBY is a tool, ZHELANG is a domain specific language that must be linked with a general purpose programming language. COMBY is much more specific: it searches and changes code structures. In this task, COMBY and ZHELANG share the same declarative flavour. Yet, given how specific COMBY is, it differs from ZHELANG in two important ways. First, none of ZHELANG’s reactive components is present in COMBY. Second, COMBY’s actions are rewriting rules, whereas ZHELANG’s can be anything, for they are invocations of separate processes.

5.2 Inductive Grammar Synthesis

The notion of *language identification in the limit*, which we have used as a motivation for our on-line grammar inference algorithm, was introduced by Edward Gold in the mid sixties [?]. Much research evolved from Gold’s initial problem formulation. The main developments in the field are due to Angluin and her collaborators [?, ?, ?, ?]. Nevertheless, several research groups have formalized grammar inference for specific types of languages [?, ?, ?, ?, ?]. Since the nineties, decidability for inference of grammars for several classes of languages is already known [?]. Usually, the language thus produced is deterministic, although Eman *et al.* have shown how to derive probabilistic automata on the limit [?]. The identification of string events fits into the framework of language inference in the limit; however, in this work, we do not try to guess the right host language L that contains said events. Instead, we try to infer a grammar G that recognizes

string events in any prefix of this language. Notice that G might also recognize strings that do not belong into L . This possibility has no practical implications in the context of this work: we are interested in finding string events, not in recognizing exactly the host language that contains it.

Recent progress in the field of machine-learning has imbued Gold’s original program with renewed attractiveness. For an overview of how machine-learning techniques are used to solve language recognition in the limit, we recommend Bennaceur *et al.*’s [?]. The literature contains several examples of how statistical inference techniques are used to learn a language in the limit, such as the work of Li *et al.* [?], who employ a genetic-based algorithm to learn the structure of XML documents. Or, along a different direction, the work of Graben *et al.* [?], who have developed an interactive system to gradually learn a simple language of English numerals. We contend that such techniques, although effective in their contexts, are not ideal fits to our problem—online language recognition—because they require slow, exploratory-based algorithms, which would be too heavy for our needs.

5.3 Program Fuzzing

In this work, we are interested in approximating a grammar that characterizes the output of a program. The inverse problem has received more attention in the programming language community: to infer a grammar that describes the input of a program. This kind of inference is useful in testing via software fuzzing, as demonstrated by Bastani *et al.* [?] and Blazytko *et al.* [?], for instance. The many approaches described in the literature [?, ?, ?, ?] differ from our work in many ways. First, there is the obvious difference in direction: we infer grammars for program outputs, not inputs. Second, these techniques typically rely on negative examples to refine the inferred grammar, whereas negative examples play no role in our formulation. Finally, there is a difference in purpose: we are not interested in testing a program; rather, our intention is to intervene in the program already in production.

5.4 Interactive Grammar Inference

There exists prior work about the construction of parsers for programming languages based on examples [?, ?, ?, ?, ?]. Such systems synthesize and refine grammars, one example at a time. Much of the inspiration behind our approach to select which literals must be redacted (see Section 3.2.1) came from Parsimony [?], an IDE for example-guided synthesis of lexers and parsers. This line of work is an instance of a much broader field known as *programming-by-examples* (PBE) [?]. ZHEFUSCATOR is not a framework to support programming by example. It infers grammars on-the-fly that recognize examples produced automatically by a machine, not a person. Therefore, the speed to synthesize a parser is an essential requirement of our work—more than clarity, or the efficiency of the parser itself. That is the reason why we have opted to produce Heap-CNF grammars: it is fast to generate and merge them.

5.5 String Events

This work is not the first work to deal with the on-line detection of string events. Research along this direction was mostly concerned with security. String events have been handled, for instance, in the context of intrusion detection [?, ?], dynamic taint analysis [?, ?], recognize regular patterns on a stream of text [?] and on-the-fly spam identification [?]. Nevertheless, if we do not claim primacy, we claim generality. All these previous works would identify string events in very specific situations, e.g., as particular patterns embedded in an SQL query, in the case of tainted flow analysis [?], or as a combination of specific tokens within a network package, in the case of intrusion detection [?]. This work is the first work to provide a general framework that, in a way, “learns” a language, and recognizes string events embedded into it.

Chapter 6

Conclusion

This work has presented two reactive methodologies to detect string events. Said events are described by a language whose grammar is known. They occur within a potentially infinite text, defined by a host language, whose grammar is unknown.

ZHEFUSCATOR shows how to synthesize a grammar G that recognizes any prefix of the infinite text stream. By defining a specific restriction of Chomsky Normal Form, the Heap-CNF, we guarantee that G is non-ambiguous (Theorem 5) and admits LL(1) parsing (Theorem 7). We have shown, empirically, that this theory can be implemented into an efficient log anonymization system, which redacts sensitive information from the output of programs, while treating these programs as black-box software. We have tested the ZHEFUSCATOR onto logs from databases (MySQL and PostgreSQL), operating systems (OSX) and Java benchmarks (DaCapo). In every case, the performance overhead of this system is very small.

With ZHELANG we have defined a general reactive methodology for handling infinite text processing. ZHELANG expands on ZHEFUSCATOR solutions by giving the user the ability to define the event grammar G . This can be done by using simple operators to define observers, which can be further binded to an action to create an event. We also have shown other applications for infinite text processing, such as implementing a web crawler and performing automatic code refactoring.

Future work. We speculate that recent developments in the programming languages community can be used to strengthen the theory and the practice discussed in this work. First, concerning formalization, our theorems are not mechanically verified. This shortcoming is due to the lack of a general framework to reason about properties of LL(1) parsers. However, Edelmann *et al.* [?] have showed how to build LL(1) parsers with derivatives and zippers that are correct by construction. Second, ZHEFUSCATOR is parameterized by a tokenizer, which our current implementation borrows from ANTLR. The fact that users have no way to specify a lexer in our system can be considered a limitation of our current implementation. Thus, it would be desirable to give users the possibility to define their own tokenizers without exposing them to minutia related to automata theory. Recent work by Chen *et al.* [?] has provided a clear interface for this purpose, which is based on examples supported by a natural language (NL) description

of regular expressions. We believe that NL-based specifications will be able to improve purely example-based approaches that have recently been shown to be effective to specify regular expressions [?, ?].

ZHELANG could also benefit on the usage of NL-based interfaces. In order to define an event the users need to provide a grammar G by writing an observer and need to define an action using a general purpose programming language. This requirements make it hard to non-programmers to use ZHELANG. The usage of template based matching, such as the one used in COMBY [?], can allow non programmers to define observers using natural language which enhances the usage of ZHELANG. However, in order to make it fully accessible to users without a programming background it is also required to synthesize actions. The programming-by-examples approach, such as the work done in flash fill in [?], has already shown that it's possible to use examples to synthesize simple string processing algorithms. This research direction is even more promising once we consider the availability of efficient string solvers such as CVC4Sy [?], which supports a wide range of logical theories, including strings and regular expressions.

bibfile

Appendix A

Proofs of Lemmas and Theorems

This appendix contains proofs of Lemmas, Theorems and Corollaries present in the work “REACTIVE METHODOLOGIES TO INFINITE TEXT PROCESSING”.

Theorem 1 Function `fill_holes` (Fig. 3.2) produces a grammar G_i that recognizes an example $t_i = t_i^1 \cdots t_i^n$ in n steps with $2n - 1$ non-terminals.

Proof 1 *The proof is by induction on the size $|t_i|$ of the example. On the **Base Case**, we have that $t_i = \text{token}$; hence, $|t_i| = 1$. `fill_holes` produces $R_1 ::= \text{token}$, which recognizes t_i trivially. On the **Inductive Case**, we assume that $t_i = \text{token} \bullet \text{Rest}$. By induction, we have that `fill_holes` generates a grammar with starting symbol R_{2i+1} that recognizes **Rest** in $n - 1$ steps (Line 7 of Figure 3.2). The extended grammar recognizes t_i :*

$$\begin{aligned} R_n & ::= R_{2n}R_{2n+1} \\ R_{2n} & ::= \text{token} \\ R_{2n+1} & ::= \dots \end{aligned}$$

By induction, we know that R_{2n+1} starts production rules with $2(n-1) - 1$ non-terminals. Adding R_n and R_{2n} , we have that the resulting grammar contains $2n - 1$ non-terminals.

Lemma 2. If G is the grammar that results from merging two Heap-CNF grammars G' and G'' , then G is Heap-CNF, and $\text{lang}(G') \cup \text{lang}(G'') \subseteq \text{lang}(G)$

Proof 2 *We demonstrate the lemma analyzing each one of the four cases involved in the process of merging two Heap-CNF grammars. We let $R_i^' ::= P^'$ be the production rule that corresponds to R_i in G_i . Similarly, we let $R_i^{''} ::= P^{''}$ be the production rule that corresponds to R_i in $G_i^{''}$. We let tk be a token:*

- $P^' = \text{tk}'_1 \mid \dots \mid \text{tk}'_n$ and $P^{''} = \text{tk}_1'' \mid \dots \mid \text{tk}_n''$. In this case, we have that $R_i ::= \text{tk}'_1 \mid \dots \mid \text{tk}'_n \mid \text{tk}_1'' \mid \dots \mid \text{tk}_n''$, which is still Heap-CNF.
- $P^' = R_{2i}R_{2i+1} \mid \text{tk}'_1 \mid \dots \mid \text{tk}'_n$ and $P^{''} = \text{tk}_1'' \mid \dots \mid \text{tk}_n''$. In this case, we have that $R_i ::= R_{2i}R_{2i+1} \mid \text{tk}'_1 \mid \dots \mid \text{tk}'_n \mid \text{tk}_1'' \mid \dots \mid \text{tk}_n''$, which is still Heap-CNF.
- $P^' = \text{tk}'_1 \mid \dots \mid \text{tk}'_n$ and $P^{''} = R_{2i}R_{2i+1} \mid \text{tk}_1'' \mid \dots \mid \text{tk}_n''$. In this case, we have that $R_i ::= R_{2i}R_{2i+1} \mid \text{tk}'_1 \mid \dots \mid \text{tk}'_n \mid \text{tk}_1'' \mid \dots \mid \text{tk}_n''$, which is still Heap-CNF.

- $P' = R_{2i}R_{2i+1} \mid tk'_1 \mid \dots \mid tk'_n$ and $P'' = R_{2i}R_{2i+1} \mid tk_1'' \mid \dots \mid tk_n''$. In this case, we have that $R_i ::= R_{2i}R_{2i+1} \mid tk'_1 \mid \dots \mid tk'_n \mid tk_1'' \dots \mid tk_n''$, which is still Heap-CNF.

Notice that if we have a token tk_x that appears in both lists: $tk'_1 \mid \dots \mid tk'_n$ and $tk_1'' \mid \dots \mid tk_n''$, then this token will appear only once—by definition—in the corresponding list of the merged grammar.

Theorem 3. The procedure `build_grammar` (Fig. 3.1) constructs grammars in Heap-CNF.

Proof 3 The proof of Theorem 3 is the junction of two facts: (i) function `fill_holes` (Fig. 3.2) builds only grammars in Heap-CNF; and (ii) the merging of grammars (Def. 4) yields Heap-CNF grammars. To demonstrate Fact-i, notice that `fill_holes` only produces rules in the format $R_i ::= \text{token}$, or $R_i ::= R_{2i}R_{2i+1}$; hence, the grammar is in Heap-CNF. Fact-ii follows from Lemma 2.

Theorem 4. Let G_1, G_2, \dots, G_n be the grammars constructed by function `build_grammar` (Fig. 3.1) for input strings t_1, t_2, \dots, t_n . Grammar $G_i, 1 \leq i \leq n$ recognizes every input $t_i, 1 \leq i \leq n$.

Proof 4 The proof works by induction on the number of examples t_i . In the **base case**, `build_grammar` fails compulsorily in the attempt to parse t_1 , because its current grammar recognizes only the empty string, i.e.: $R_1 ::= \epsilon$. Failure happens in the conditional at Line 8 of Figure 3.1. A new grammar G_1 will be constructed for t_1 by routine `expand_grammar`, via function `fill_holes`. By Theorem 1, G_1 recognizes t_1 . In the **inductive step**, we have a grammar G_k , that recognizes every example t_1, \dots, t_k . When `build_grammar` is given a new example t_{k+1} , two scenarios are possible:

- G_k recognizes t_{k+1} ; hence, the conditional at Line 19 of Figure 3.1 is true.
- G_k fails to recognize t_{k+1} . In this case, a new grammar G' will be constructed by `fill_holes`, and the resulting grammar $G_{k+1} = \text{merge}(G_k, G')$ recognizes t_1, \dots, t_{k+1} , by Lemma 2.

We let $\text{merge}(G_k, G')$ above be the grammar that results from merging G_k and G' .

Lemma 3. Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . The size of G_n is $O(N)$, where N is the number of tokens in t_1, t_2, \dots, t_n .

Proof 5 The `fill_holes` procedure only augments the rightmost node of a derivation tree. In other words, given a sentence of n tokens, `fill_holes` produces a grammar with¹:

- $2n - 1$ non-terminal symbols;

¹We treat s_e , the starting symbol of the event grammar, as a single token.

- $2n - 1$ production rules;
- n terminal symbols;

The *merge* routine never adds new terminals or non-terminals to a grammar; hence, it maintains its asymptotic size complexity.

Theorem 5. Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . G_n is not ambiguous.

Proof 6 As a consequence of Lemma 3, the rightmost derivation tree of a Heap-CNF grammar always has height $n - 1$ and $O(N)$ nodes. Only one rightmost derivation tree is possible, which Figure 3.4 illustrates. The rightmost token is always recognized by a production from non-terminal R_{2^n-1} .

Corollary 2. Let G_n be the grammar constructed by function `build_grammar` (Fig. 3.1) after observing inputs t_1, t_2, \dots, t_n . G_n recognizes $t_i, 1 \leq i \leq n$ with $O(N)$ derivations, where N is the number of tokens in t_i .

Proof 7 This corollary follows from Lemma 3, plus the fact, already mentioned in the proof of Theorem 5, that only one rightmost derivation tree is possible. Thus, the grammar built by `fill_holes` recognizes a sentence with n tokens with $2n - 1$ derivations.

Theorem 6. Grammar G'_e produced by `markup` (Fig. 3.7) recognizes a subset of $\text{lang}(G_e)$ or the empty language.

Proof 8 The proof works by induction on the number of times Step 2 in Procedure `markup` runs. In the **base case** (Step 1), we have that G'_e recognizes the empty language. In the **inductive step**, we assume that G'_e recognizes a subset of $\text{lang}(G_e)$ after n iterations of Step 2. In the next iteration, Steps 3 and 4 ensure that G_e'' recognizes a subset of $\text{lang}(G_e)$. The junction of G'_e and G_e'' uses only production rules of G_e ; hence, it must recognize a subset of the language that G_e recognizes. Furthermore, because these two grammars start with s_e , the initial symbol of G_e , the resulting grammar after the junction also starts with s_e .

Theorem 7. Any Heap-CNF grammar is LL(1).

Proof 9 This fact follows from the observation that Heap-CNF grammars are not recursive. Therefore, no left recursion is possible, and the language that these grammars recognize has a finite number of possible derivation trees. The one token of lookahead follows from Definition 3 and Corollary 1, because the position of a token in the derivation tree is uniquely determined by the position of that token in the input string.

Corollary 3. There are languages whose grammars cannot be synthesized by ZHEFUSCATOR.

Proof 10 *A formal language is called an $LL(k)$ language if it has an $LL(k)$ grammar. The set of $LL(k)$ languages is properly contained in that of $LL(k+1)$ languages, for each k greater than or equal to zero [?]. Therefore, there exist context-free languages that are not $LL(1)$. This restriction mean that even on the limit, ZHEFUSCATOR would not be able to synthesize perfect grammars for some languages. However, up to any number n of events, ZHEFUSCATOR will synthesize a grammar G_n that recognizes every $t_1, \dots t_n$, and potentially other strings, as discussed in Section 3.1.3.2.*