



SCHOOL OF COMPUTATION, INFORMATION
AND TECHNOLOGY — INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

Inferring String Properties from Code Property Graphs

Severin Schmidmeier





SCHOOL OF COMPUTATION, INFORMATION
AND TECHNOLOGY — INFORMATICS

TECHNISCHE UNIVERSITÄT MÜNCHEN

Bachelor's Thesis in Informatics

Inferring String Properties from Code Property Graphs

Herleitung von Eigenschaften von Strings aus Code Property Graphen

Author: Severin Schmidmeier

Supervisor: Prof. Dr. Claudia Eckert

Advisor(s): Alexander Küchler, Florian Wendland

Submission: 15.03.2023



I confirm that this bachelor's thesis is my own work and I have documented all sources and material used.

Ich versichere, dass ich diese Bachelorarbeit selbstständig verfasst und nur die angegebenen Quellen und Hilfsmittel verwendet habe.

München, 15.03.2023

(Severin Schmidmeier)

Acknowledgments

I want to express my gratitude to my friends for their constructive feedback, moral support, and encouragement.

Furthermore, I would like to thank my advisors for their feedback and help and Prof. Eckert for providing the topic of this thesis.

Lastly, I need to thank my cat for his help and the valuable in-depth discussions about this thesis we had.

Abstract

As more and more aspects of our daily lives move online, the security of online data becomes increasingly important. Injections – i.e., the insertion of unwanted strings – are among the most common threats to online security. Identifying potential injection vulnerabilities in source code is crucial to preventing this type of attack. Static analysis tools that try to gain insight into the behavior of programs from their source code and are used for this detection can benefit from precise information about the values of potentially vulnerable strings. Currently, most of the available static analysis tools have no way of obtaining such information.

In this thesis, we adapt an approach to extract formal grammars that describe strings from a code property graph. We then approximate these grammars to be strongly regular and generate human-readable regular expressions from them. To accomplish this, we transform them into automata and employ the state elimination strategy.

The obtained regular expressions describe properties of the values the analyzed strings can take on and match all such values. We provide a proof of concept implementation and show that it is able to accurately describe complex strings. We believe that, with further refinement, the information our approach provides can be used to successfully detect, for example, SQL injection vulnerabilities. This can increase the capabilities of static analysis tools and helps in preventing security issues.

Contents

Acknowledgments	i
Abstract	ii
1. Introduction	1
1.1. Motivation	1
1.2. Contribution	2
1.3. Outline	2
2. Problem Description	3
3. Background	5
3.1. Formal Grammars	5
3.2. Strongly Regular Grammars	6
3.3. Automata	6
3.4. Regular Expressions	7
3.5. Code Property Graph	7
4. Related Work	9
4.1. String Property Analysis	9
4.2. String Property Analysis Used for SQL Injection Vulnerability Detection .	10
4.3. Other Approaches for SQL Injection Vulnerability Detection	11
4.4. Overview of the Different Approaches	12
5. Approach and Implementation	14
5.1. General Approach	14
5.2. Grammar Creation	15
5.3. Regular Approximation	17
5.3.1. Character Set Approximation	17
5.3.2. Mohri-Nederhof Approximation	20
5.4. Strongly Regular Grammar to Automaton	22
5.4.1. Algorithm	22
5.4.2. Operation Productions	26
5.5. Automaton to Regular Expression	29
5.5.1. State Elimination	29
5.5.2. Delgado Heuristic	30
5.6. Hotspot Collection	31

6. Evaluation and Discussion	32
6.1. Evaluation and Benchmarking	32
6.1.1. Quality	32
6.1.2. Performance	37
6.2. Limitations of the Evaluation	41
6.3. Future Work	42
6.3.1. Assertions	42
6.3.2. Polyvariance	43
6.3.3. More Extensive Implementation	44
6.3.4. State Elimination Heuristics	44
6.3.5. Automata Centric Approach	45
7. Conclusion	46
A. Evaluation Results	47
Bibliography	49
List of Figures	52
List of Tables	53
List of Listings	54
List of Acronyms	55

1. Introduction

We first provide the motivation for our work in Section 1.1, then describe our contribution to the field in Section 1.2, and further establish the outline for the rest of the thesis in Section 1.3.

1.1. Motivation

The increasing reliance on software applications in various aspects of modern life has led to a growing concern for the security of these applications. Among the many security threats that can affect software, injection vulnerabilities are among the most dangerous and prevalent. The Open Web Application Security Project (OWASP) consistently lists injection attacks, which include SQL injections, LDAP injections, and command injections, as one of the top ten web application security risks [23].

Injection vulnerabilities occur when an attacker is able to insert malicious code or strings into an application, often through input fields that accept user input, such as search boxes or login forms. This can result in the attacker gaining unauthorized access to sensitive data, executing arbitrary code, or even taking control of the entire system.

To assist developers and security researchers in spotting injection vulnerabilities in code, a variety of tools and techniques, including static analysis tools for string values exist. These tools analyze the source code of an application to identify potential vulnerabilities, including injection vulnerabilities. Static analysis tools are particularly useful because they can detect vulnerabilities that may not be apparent during testing or manual code review. Making developers aware of such issues during development enables them to fix the vulnerability early. One such static analysis tool is Codyze¹, which uses the implementation we extend in this thesis for its analyses.

In order to detect injection vulnerabilities, these tools may try to analyze the possible values a string that is passed as a query to e.g. a database can take on. From these inferred properties, a tool can then assess whether the analyzed program contains any potential injection vulnerabilities and alert the programmer.

Since the analyzed strings often contain unknown user inputs, enumerating all possible values is not feasible. Therefore, we use regular expressions to describe the

¹<https://www.codyze.io>

1. Introduction

properties of all possible values. For example, consider a string variable, that is used as an SQL query and is determined to be described by the regular expression `DELETE * FROM myTable WHERE id='.*'`. This information can be used to issue a warning during a static analysis, because the analyzed program allows for an arbitrary unchecked string to be inserted into the SQL query, which poses a severe security vulnerability.

1.2. Contribution

In this thesis, we extend a Code Property Graph (CPG) implementation [29] to increase its capabilities in analyzing string values. We adapt the theoretical approach by Christensen et al. [3], which creates regular languages describing the values of a string, to the present CPG implementation. Note that Christensen et al. work on a different representation of the analyzed code and use the results in a different way. For example, they use a different definition of the data flow graph (DFG) and a novel data structure representing the complete analyzed code, from which they extract deterministic finite automata (DFAs). We only analyze parts of the code that are required for the specific query. Therefore, using their approach requires adaptation of the used techniques and solving of new problems like the handling of string operations for our use case. We also want to provide the information to an analyst in a human-readable format, in our case in the form of regular expressions. To achieve this, we first combine the mentioned approach with the Nederhof algorithm [22] to transform the obtained results to automata. Further, we use state elimination in combination with a heuristic by Delgado and Morais [5] to convert them to regular expressions for users. We also provide a proof of concept implementation covering a subset of the Java standard library.

1.3. Outline

We first describe the problems and posed research questions in Chapter 2. After providing some theoretical background in Chapter 3, we then present selected related work in Chapter 4 and describe the different steps of our approach in Chapter 5. In Chapter 6, we evaluate the results and benchmark our implementation. There, we also highlight some limitations of our approach and evaluation and include ideas for future continuation and improvement of the presented design before concluding the outcome of this thesis in Chapter 7.

2. Problem Description

We extend a CPG implementation by the Fraunhofer AISEC research institute [29]. It currently has no means of providing information about the structure and contents of string variables that go beyond propagating literals if they are not changed.

Describing such strings is not trivial, as often at least part of a given string stems from an unknown source, for example runtime user input. We solve this issue by first describing a given string with a formal grammar, which conservatively approximates the values the string can take. This means that for a string s , the language generated by the grammar we obtain to describe s always contains all possible values s can have. We want to use regular expressions as a final representation because they are human-readable and therefore allow users to manually evaluate our results for a security analysis. Regular expressions also allow automated tools to match given queries.

The obtained grammars generate context free languages, which are a superset of the regular languages accepted by regular expressions. Therefore we can't directly convert the obtained grammars to regular expressions without any loss of information, but rather need to approximate them with grammars that generate regular languages first.

This approximation poses the challenge of deciding, which information to retain and which parts to change. We also need to account for the effects of operations like a `replace` function on the analyzed strings. Additionally, it is necessary that our approximation stays conservative.

Furthermore, since the grammars our approximation creates are not textbook regular grammars, but rather strongly regular grammars, we need to use algorithms suited to this type of grammar for the conversion.

Because we use automata as an intermediary step in the conversion from grammar to regular expression, we use the state elimination algorithm [1] to convert a nondeterministic finite automaton (NFA) to a regular expression. As we want our results to be human-readable, we need to reduce the length of the resulting regular expression by optimizing the state elimination algorithm.

2. Problem Description

To summarize, this section poses the following research questions.

1. Which approaches for obtaining information about string values exist?
2. Can such an approach be adapted to the characteristics of the given CPG implementation?
3. How can we approximate the obtained grammar?
4. How can we transform the approximated grammar into a human-readable format?
5. How can we resolve the effect of operations on strings during this transformation?
6. How can we improve the quality of our result and the performance of our approach?

3. Background

We first provide some background information and notation for formalisms used in the following chapters. This includes formal grammars in Section 3.1, strongly regular grammars in Section 3.2, finite automata in Section 3.3, and regular expressions in Section 3.4. In Section 3.5, we then describe the code property graph we use.

3.1. Formal Grammars

A formal grammar consists of a set of nonterminal symbols N , an alphabet Σ of terminal symbols, a set of production rules $(\Sigma \cup N)^*N(\Sigma \cup N)^* \rightarrow (\Sigma \cup N)^*$, also called just “productions” and a start nonterminal $S \in \Sigma$.

In the context of this thesis, Σ is the set of characters making up the strings of the analyzed programming language. In some places, we use regular expressions as terminals, which is an extension to the usual definition.

To define grammars in examples, we use the notation $X \rightarrow Y$ to specify a production rule to transform a nonterminal X to some other nonterminal Y . We use capital letters for nonterminals and lowercase letters for terminals. We also do not specify the start symbol explicitly, but rather the nonterminal on the left hand side of the first given production is considered to be the start symbol.

Context-free grammars (CFGs) are grammars where the left hand side of all productions consists of a single nonterminal. The set of production rules for CFGs therefore is defined as $N \rightarrow (\Sigma \cup N)^*$ compared to the more general definition above.

Regular grammars are a subset of CFGs where the right hand side of each production consists either of a single terminal $a \in \Sigma$ or of exactly one nonterminal and one terminal. Additionally, the nonterminals on the right hand side are either always the first symbol or always the last symbol on the right hand side.

Therefore, for example, a grammar containing both the productions $A \rightarrow aB$ and $B \rightarrow Ab$ with $A, B \in N \wedge a, b \in \Sigma$ is not regular. Here, the nonterminal is the last symbol in the first production but the first symbol in the latter, which does not comply with the above rule.

3.2. Strongly Regular Grammars

\mathcal{R} is the equivalence relation defined on the set of nonterminals N of some grammar:

$$A\mathcal{R}B \Leftrightarrow (\exists \alpha, \beta \in V^* : A \xrightarrow{*} \alpha B \beta) \wedge (\exists \alpha, \beta \in V^* : B \xrightarrow{*} \alpha A \beta) \quad (3.1)$$

Here, V is $\Sigma \cup N$, i.e. the set of all symbols, terminal and nonterminal. $\xrightarrow{*}$ is the reflexive and transitive closure of the production relation \rightarrow defined by the set of productions in the grammar. In other words, iff $A \xrightarrow{*} \alpha B \beta$, there exists a sequence of productions starting at the nonterminal A to produce a set of symbols that contain B . Therefore, \mathcal{R} groups all nonterminals into disjoint equivalence classes, in which each nonterminal in a class can be produced by each other nonterminal in the class. These nonterminals are called mutually recursive.

A grammar is strongly regular iff the production rules in each such equivalence class are either all right-linear or left-linear. A production rule is right-linear iff it is of the form $A \rightarrow w\alpha$, where w is a sequence of terminal symbols and α is empty or a single nonterminal symbol. Left-linear productions are defined accordingly but the nonterminal is on the left side of the production result.

Strongly regular grammars are guaranteed to generate regular languages [20]. Despite the potentially misleading name, strongly regular grammars are not a subclass of regular grammars but rather an extension to them. Every regular grammar is strongly regular, but not vice versa.

3.3. Automata

A deterministic finite automaton (DFA) consists of a set of states Q , an alphabet of input symbols Σ , a transition function $\delta : Q \times \Sigma \rightarrow Q$, an initial state $q_0 \in Q$ and a set $F \subseteq Q$ of accepting states. For a nondeterministic finite automaton (NFA), from a given state multiple states can be reached with the same input, so the transition function is $\delta : Q \times \Sigma \rightarrow 2^Q$, where 2^Q denotes the power set of Q .

We represent automata as graphs, where each state is a node and the transition function is represented by edges labeled with elements of Σ . The start state is marked with an incoming arrow and the accepting states are marked with double circles. An edge in this graph, henceforth also called transition, is denoted as (q_1, a, q_2) , where $q_1 \in Q$ is the origin state, $a \in \Sigma$ is the label and $q_2 \in Q$ is the target state of the edge.

3.4. Regular Expressions

We use a regular expression syntax with metacharacters from the Java regular expression flavor, which can be seen in Table 3.1.

Metacharacter	Name	Description
*	Kleene star	Matches the previous character zero or more times
.	Wildcard	Matches any character
?	Option	Matches the previous character zero or one time
	Choice	Matches either the previous or the following expression
[abc]	Character class	Matches any of the contained characters (here a, b and c)
[^abc]	Negative character class	Matches any character not contained (here anything except a, b and c)

Table (3.1) The regular expression metacharacters we use

Expressions are grouped using round brackets and meta characters are escaped using single backslashes. In a character class consecutive characters can be abbreviated using -, e.g. [0-9], to match any digit.

Languages accepted by regular expressions that use more advanced features like look-ahead assertions or backreferences are not necessarily regular [2]. However, with the presented limited syntax only expressions accepting regular languages can be generated, as all used features can be simulated using the features present in the textbook definition of regular languages.

3.5. Code Property Graph

A Code Property Graph (CPG) is a directed multi graph, where the nodes represent syntactic elements like simple expressions or function declarations and the edges represent the relations between those elements. While the CPG is a general concept, we focus on one the implementation¹ by the Fraunhofer AISEC that extracts a CPG out of source code of a set of different programming languages.

¹<https://github.com/Fraunhofer-AISEC/cpg>

3. Background

Here, nodes and edges have a list of key-value pairs called properties which contain general information for the element. For example, a node representing a statement in a source file contains the location of the underlying code and an edge representing evaluation order may contain whether the target statement is unreachable. The graph is initially created by language frontends, which create partially connected abstract syntax trees (ASTs), which are then enriched by additional information like the mentioned evaluation order and data flow information by multiple passes [29]. A pass iterates the graph in some way to enhance it with, for example, additional edges or other information. Users of the library can extend its functionality by adding additional passes.

While the CPG contains many different types of edges, the most relevant edge type for this thesis is data flow edges, which represent the data flow between different expressions.

```
String s1 = "xyz"; 1
System.out.println(s2); 2
```

Listing (3.1) CPG example code

Consider the short code example in Listing 3.1. Here, among others, the following nodes are part of the CPG:

- **Literal**, representing the string literal "xyz" in line 1
- **VariableDeclaration**, representing the declaration and initialization of **s** in line 1
- **DeclaredReferenceExpression**, representing the reference to **s** in line 2

In this example, the data flows from the **Literal** node to the **VariableDeclaration** and from there to the **DeclaredReferenceExpression**.

The nodes connected by those edges effectively form a subgraph of the CPG, the DFG, from which we then extract the information on string values. Figure 3.1 shows this part of the CPG.

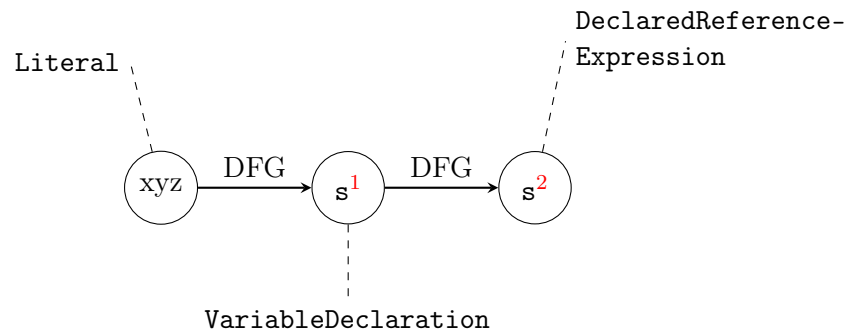


Fig. (3.1) Part of the CPG for the code in Listing 3.1

4. Related Work

In this section we present other work related to this thesis. We first discuss work related to the analysis of string properties in Section 4.1, which is where our implementation fits in. We also include other applications of such analyses. Further, we present different approaches utilizing such methods for the specific application of detecting SQL injection vulnerabilities in Section 4.2. Additionally, in Section 4.3 we introduce other techniques not closely related to the string properties for detecting such vulnerabilities.

4.1. String Property Analysis

Our approach partially follows the approach by Christensen et al. [3]. The authors construct a context-free grammar from a flow graph, but instead of creating it on-demand, starting at the chosen hotspot node like we do, they consider the total flow graph for grammar creation. As the total amount of potential hotspots can be far greater than the number of actual points of interest, this total grammar may be potentially significantly larger than needed, thus reducing the performance of the subsequent steps. Christensen et al. use the same approximation methods for obtaining regular languages from the generated context free grammars, but instead of directly transforming these grammars into automata like we do, they produce a novel formalism, the multi-level automaton (MLFA). This automaton allows the extraction of the respective automata for different hotspots. Due to the aforementioned on-demand generation of the grammar, we don't need the extraction capability for single hotspots that the MLFA provides. In contrast to our implementation, Christensen et al. also do not transform the obtained automata to regular expressions, but rather provide query options using automata. The authors provide a feature rich implementation of their approach and show that it efficiently produces useful results.

Tabuchi et al. [25] describe a type system for a minimal functional calculus where strings have a regular expression as their type. Using type inference and reconstruction algorithms, they assign a type to each string variable. Properties of a variable can then be obtained from its type. They show that their proposed type system can produce good results when applied to their minimal calculus. While we considered adopting this approach for the analysis, there are some problems, especially due to our different requirements and

4. Related Work

prerequisites. To use the presented approach in practice an efficient algorithm for type checking and type reconstruction is needed. The given paper does not include these, but rather indicates several problems in constructing such algorithms for the given situation without losing some of the desired preciseness. The authors mention that using standard type reconstruction by constraint solving for the proposed type system is impossible due to limitations of regular languages. Additionally, this approach is tailored to the mentioned calculus and utilizes specific features like pattern matching, which would make adapting it to our use case more difficult. The additional layer of abstraction introduced by the DFG used in the approach we chose eliminates this problem and makes adaptation easier.

Kirkegaard et al. present an XML transformation system, X_{ACT}, that functions as a part of the J_{WIG} framework developed by Christensen et al.’s research group [16]. This system statically verifies the validity of XML transformations by tracking sets of XML values during dataflow analyses. Comparable to the mentioned approach by Tabuchi et al. [25], the XML transformation language XDuce also uses regular expression types to describe string values [15]. It therefore compares to J_{WIG} as the approach by Tabuchi et al. compares to the previously mentioned work by Christensen et al. [3]. Both of these tools use basic analyses of string values for the purpose of validating XML generation. More advanced approaches, like the work by Christensen et al. we adopt, provide more robust and more generic results for string analysis.

There exist frameworks to model different analyses using set constraints, for example the BANSHEE toolkit by Kodumal and Aiken [17]. BANSHEE allows the user to define a specification of the used constraints, from which a constraint resolution engine is built. Using such a generic framework, different analyses based on constraint resolution can be created, for example pointer analysis or context-free language (CFL) reachability. However, Christensen et al. argue that some operations of their – and by extension our – approach cannot be captured by the constraint operators [3].

4.2. String Property Analysis Used for SQL Injection Vulnerability Detection

Gould et al. [8] build on the work of Christensen et al. and use the obtained automata for further analysis of SQL queries. They use a context-free language reachability algorithm to validate the semantic correctness of SQL queries. For this analysis, they also include the grammar of the SQL language for scoping information and the database schema for type checking. The authors show that their tool can precisely detect errors in SQL queries. However, they currently do not detect SQL injection vulnerabilities, but rather only assess the semantic correctness and type safety of the queries.

4. Related Work

Halfond and Orso also leverage the functionality provided by Christensen et al. to detect SQL injection vulnerabilities using the AMNESIA tool [12]. Their approach differs from the other mentioned work, as it is not entirely a static analysis, but rather a static analysis used to enhance the capabilities of runtime monitoring. They use the implementation provided by Christensen et al. to build an SQL query model from the obtained automata. The model represents all of the possible SQL queries that can be generated at the analyzed hotspot. During runtime, the monitor can now check the dynamically generated queries against the existing query model. If a runtime query does not match the model, the monitor can reject the malicious query to prevent an SQL injection. In another publication, the authors show that their implementation of this dynamic approach can successfully prevent SQL injections in real world use cases without imposing considerable execution overhead on the application [11].

As our approach is able to produce DFAs for the analyzed code as well, the previous approaches could theoretically be adapted to our implementation.

Wassermann and Su [28] present an approach comparable to ours, where they also characterize values of string variables using context free grammars. They specifically target SQL injection vulnerabilities by tracking user-modifiable data and using the generated CFGs to check whether this data can change the syntactic structure of a query. While this approach is successful in detecting those vulnerabilities, our approach is more general and not focused on detecting one specific type of problem, but rather on providing general information for unspecified further use.

4.3. Other Approaches for SQL Injection Vulnerability Detection

Livshits and Lam present an approach for vulnerability detection based on taint propagation [19], which involves finding all sinks derivable from sources via some given derivation rules. In the context of vulnerability detection, for example, a function that extracts some user input from a request is a source, while an SQL query execution function is a sink. A sink is flagged as potentially vulnerable if it is derivable from a potentially malicious source via some set of derivation rules. The authors use a precise points-to analysis to find paths that allow unsanitized data to flow into e.g. SQL queries. As the different sinks, sources, and derivation rules for a specific vulnerability need to be provided by the user of their tool, they provide a user-friendly way using Process Query Language (PQL) to specify vulnerability patterns. The authors show that their approach successfully detects previously unknown security vulnerabilities in real world applications.

Another comparable approach by Halfond et al. uses positive tainting [13]. Their implementation tracks trusted strings rather than untrusted ones and allows only tainted

4. Related Work

strings in queries. The authors argue that this approach is especially suited for preventing SQL injection attacks, as due to the positive tainting incompleteness in the analysis leads to false positives, which can be filtered, instead of false negatives. They show that their implementation of this approach successfully and efficiently detects all tested SQL injection vulnerabilities without false positives.

4.4. Overview of the Different Approaches

In Table 4.1 we provide an overview of the presented approaches and their respective differences. Here, JSA is used as an abbreviation for “Java String Analyzer”, the implementation of their approach by Christensen et al., DB for database and PoC for proof of concept. Generally, the presented work can be categorized into two groups: The first includes generic approaches that aim to provide information and capabilities for different analyses building on them ([3, 17, 25]). The publications in the second group try to solve specific problems like preventing SQL injection attacks, sometimes using the aforementioned analysis techniques as a foundation ([8, 12, 13, 15, 16, 19, 28]). Our work could be categorized into the first group.

4. Related Work

Publication	Purpose	Approach	Provides implementation
Christensen et al. [3]	provide general information	Java String Analyzer (JSA)	✓
Tabuchi et al. [25]	provide general information	regular expression type system	
This thesis	provide general information	adaption of JSA, Nederhof Algorithm, state elimination	PoC
Kirkegaard et al. [16]	validate XML transformation	simple static string analysis	✓
Hosoya & Pierce [15]	validate XML transformation	regular expression type system	✓
Kodumal & Aiken [17]	toolkit for building different analyses	set constraint resolution	✓
Gould et al. [8]	statically verify SQL queries	JSA & CFL reachability with DB information	prototype
Halfond & Orso [11, 12]	prevent SQL injection attacks	JSA & runtime monitoring	✓
Wassermann & Su [28]	prevent SQL injection attacks	CFG based, detect changes to structure of SQL query	prototype
Livshits & Lam [19]	prevent SQL injection attacks	static taint analysis, points-to analysis	✓
Halfond et al. [13]	prevent SQL injection attacks	positive taint analysis	✓

Table (4.1) Comparison of the presented related work

5. Approach and Implementation

We first provide a general overview of our approach and describe which steps are needed for the transformation from DFG to regular expression. Subsequently, we provide a detailed description of each step.

5.1. General Approach

The general approach for our implementation is adapted from the one described by Christensen et al. [3]. Conceptually, we first create a context-free grammar (CFG) from the DFG in the process described in Section 5.2. The created CFG is then approximated to a strongly regular grammar (SRG) using the Character Set Approximation described in Section 5.3.1 and the Mohri-Nederhof algorithm described in Section 5.3.2. We then transform this SRG into a nondeterministic finite automaton (NFA) using Nederhof's algorithm described in Section 5.4. Finally, we describe how to create a regular expression from this automaton using the state elimination strategy in Section 5.5. The following diagram visualizes this process and the different steps from a graph to different types of formal grammars to a regular expression.

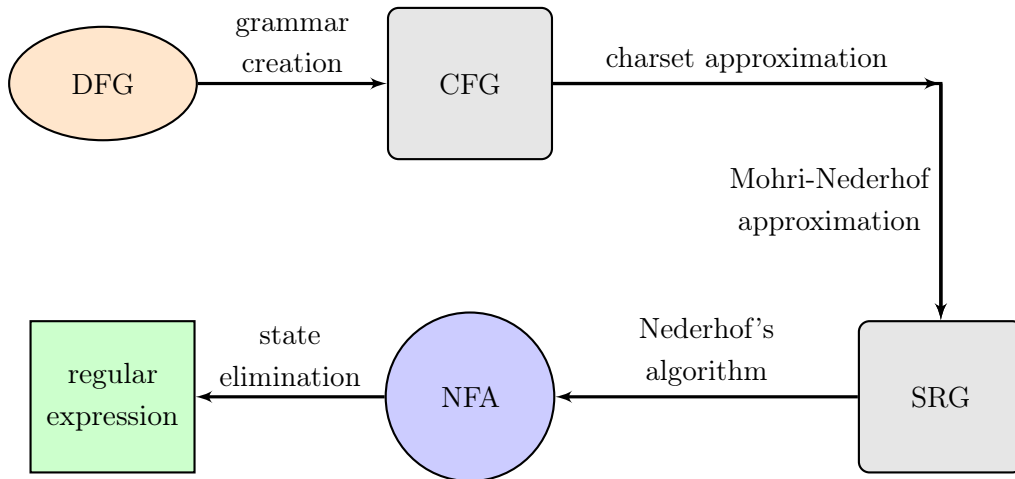


Fig. (5.1) The general approach for obtaining regular expressions

5.2. Grammar Creation

To create the grammar for a given CPG node, we traverse the DFG backwards, starting at the given node. For each visited node, we add a **Nonterminal** and the fitting productions to our grammar.

The different types of productions we use can be seen in Table 5.1, where `<terminal>` represents a terminal symbol containing a regular expression that describes a string value and “*op*” is a placeholder for a string operation that is applied to some arguments.

Name	Production	Usage
UnitProduction	$X \rightarrow Y$	references between nodes
ConcatProduction	$X \rightarrow Y Z$	concatenation of two nodes
TerminalProduction	$X \rightarrow \langle \text{terminal} \rangle$	literal string values and other terminal symbols
OperationProduction	$X \rightarrow \text{op}(Y)$	operations on strings

Table (5.1) Utilized production types

Consider the code example in Listing 5.1 for the following explanations of the different productions.

```
String s1 = "foo";           1
s2 = s3 + "bar";           2
s4 = s5.trim();           3
```

Listing (5.1) Grammar creation example code

UnitProductions mostly represent references between nodes where the underlying string is not changed. In Listing 5.1, this would be the case for the reference from s^3 to the variable declaration in line 1.

ConcatProductions are created for **BinaryOperator** nodes that represent a string concatenation using the + operator. For the example in Listing 5.1, the nonterminal corresponding to the **BinaryOperator** node for the + in line 2 would have a **ConcatProduction** with the right hand side nonterminals corresponding to the nodes for s^3 and the string literal respectively.

TerminalProductions point to a **Terminal** that represents a fixed regular expression. For example, for the **Literal** CPG node representing the "bar" string literal, the corresponding nonterminal has a **TerminalProduction** where the **Terminal** contains a regular expression that matches only the string "bar". **TerminalProductions** also occur at CPG

5. Approach and Implementation

nodes without incoming DFG edges where the value is not known. Those nodes could represent any string value, and therefore, the corresponding **Terminal** contains the regular expression `.*`, matching all strings.

OperationProductions represent function calls or other operators. The CPG for Listing 5.1 contains a **CallExpression** representing the function call to the library function `trim`. We create a **Trim** object representing this operation and the **OperationProduction** $X \rightarrow trim(Y)$, where X is the nonterminal corresponding to the node representing `s`⁴ and Y to the one representing `s`⁵. All operation objects like **Trim** also contain information about possible arguments and implement a character set transformation and an automaton transformation. These transformations describe the effect of the operation on the set of characters making up the words the operation is applied on or automata accepting those words respectively. Examples and how these transformations are used for the approximation are described in Section 5.3. This language agnostic representation of string operations allows developers of the CPG library to add support for functions and operators in other languages with different semantics compared to the corresponding Java functions, without needing to change the grammar approximation. For example, for the Python expression `"abc" * 5`, the `*` operator can be represented using a generic **Repeat** operation object. For all further steps it is not relevant whether this repeat operation is created from the mentioned Python operator `s * n` or from the corresponding Java function `s.repeat(n)`.

Improvements

Unlike Christensen et al. [3], we do not consider the total DFG when extracting the grammar. They parse the whole graph into a grammar describing all nodes, while we create the grammar starting from a single node and ignore all parts of the graph not connected via DFG edges to this node.

Since the majority of a large program is often not relevant for a specific node, this reduces the amount of nodes we need to handle during analysis. Consequently, this reduces the size of the resulting grammar, therefore leading to performance improvements.

Additionally, we can traverse the DFG conditionally, stopping at nodes representing numbers. If the traversal reaches such a node, we use an existing analysis that tries to compute the precise value. For example, for an integer created by the usual arithmetic operations, this analysis can obtain the resulting value. In this case, we can add a **TerminalProduction** with the **Terminal** representing the value literal and otherwise, if the value is not known, the **Terminal** contains a regular expression matching all numbers of the present type, e.g. `"0|(-?[1-9][0-9]*)"` for integers.

5.3. Regular Approximation

To transform the created grammar into a regular expression, we need to approximate the CFG we obtained like described in the previous section. Since basic regular expressions accept regular languages, the result of this approximation has to be a type of grammar that produces only regular languages to allow direct conversion to regular expressions without losing information. Using the two different approximation steps described in the following section, the CFG is approximated into a strongly regular grammar.

5.3.1. Character Set Approximation

To use the Mohri-Nederhof approximation algorithm described in Section 5.3.2, we first need to eliminate all cycles in our grammar that contain operation productions [20].

First, we view the grammar as a graph in which each symbol of the grammar corresponds to one graph node and for each production there are edges from the nonterminal on the left hand side to all symbols on the right hand side. For two nonterminals A and B , there is an edge from the node corresponding to A to the one corresponding to B if and only if there exists a production of form $A \rightarrow \alpha B \beta$ with α and β sequences of arbitrary symbols. This graph allows us to group terminals that are reachable from each other by finding the strongly connected components (SCCs) of the graph.

All nonterminals are assigned a character set, containing all characters that make up the words in the language of the corresponding nonterminal.

$$\begin{aligned} S &\rightarrow \text{replace}[c, x](A) \\ A &\rightarrow BC|CB \\ B &\rightarrow "ba" \\ C &\rightarrow "ca" \end{aligned}$$

Fig. (5.2) Example grammar

We assign a character set to each nonterminal N using a fixed-point iteration inside the graph component of N that constructs a character set $C(N)$ for N from the character sets of the nonterminals on the right hand side of N 's productions.

For productions with terminals on the right hand side, the character set is just the set of all characters occurring in the terminal. To account for the terminals that are regular expressions and therefore may include meta characters, which should not be included in the character set, we create the character set when we construct the regular expression. Consider, for example, the regular expression representing integers mentioned above.

5. Approach and Implementation

There, the corresponding character set contains all digit characters and the minus sign. For concatenation productions like $A \rightarrow BC$ in the example above, we take the union of the character sets of the two nonterminals on the right hand side.

Considering the example grammar in Figure 5.2, B represents the word ba and C represents ca , therefore the corresponding sets of characters are $\{'a', 'b'\}$ and $\{'a', 'c'\}$ respectively. The words that can be generated from A are combinations of B and C and therefore always contain all characters in the character sets of B and C . Thus, the character set for A is $\{'a', 'b'\} \cup \{'a', 'c'\} = \{'a', 'b', 'c'\}$.

Each operation defines a character set transformation - a function $T_{op} : 2^\Sigma \rightarrow 2^\Sigma$ - that approximates how the application of the given operation changes the character set. Here, Σ represents the set of all possible characters. For example the character set transformation for a **replace** operation, where a known character o is replaced by a known character n has the character set transformation described in Formula 5.1.

$$T_{replace[o,n]}(S) = \begin{cases} (S \setminus \{o\}) \cup \{n\}, & \text{if } o \in S \\ S, & \text{if } o \notin S \end{cases} \quad (5.1)$$

In comparison, for a **replace** operation, where the newly inserted character is not known, the transformation is defined as shown in Formula 5.2. Here, if the replaced character is contained in S , the set is transformed to Σ , since the newly inserted character could be any element of Σ .

$$T_{replace[o,?]}(S) = \begin{cases} \Sigma & \text{if } o \in S \\ S, & \text{if } o \notin S \end{cases} \quad (5.2)$$

These approximations are used in the fixed-point computation to assign character sets. As mentioned above, in the example in Figure 5.2 the character set for A is $\{'a', 'b', 'c'\}$. To obtain the character set of S , we apply the transformation defined by the $replace[c, x]$ operation to the character set of A , which gives us $(\{'a', 'b', 'c'\} \setminus \{'c'\}) \cup \{'x'\} = \{'a', 'b', 'x'\}$ as the character set of S .

To determine the SCCs, we use Tarjan's algorithm [26]. The SCCs of a graph form a directed acyclic graph (DAG), because if the graph of SCCs would contain a cycle, all contained components would be strongly connected and therefore joined into one component. This DAG implies that there exists a topological ordering of the SCCs [4]. Tarjan's algorithm topologically sorts the returned components in reverse order as a byproduct, which is necessary for the fixed-point iteration to terminate. During the

5. Approach and Implementation

computation, for a given nonterminal N , its charset is updated using the character sets of its successors. The reverse topological ordering of the components ensures that the first handled component is the root in the graph formed by the SCCs, while leafs in this graph are handled last. This ensures that the successors of each nonterminal are either in the same component or in a component that has already been handled earlier.

To break up the cycles containing operation productions, we replace one operation production $X \rightarrow op(Y)$ in each cycle with a production $X \rightarrow r$, where r is the regular expression that matches the language $C(X)^*$. Here, $C(X)$ again denotes the character set we assigned to the nonterminal X .

To find the cycles in the grammar, we check for each nonterminal A in a given component M , whether it has an operation production, and if yes, whether one of the nonterminals on its right-hand side is also part of M . If this is the case, by definition of SCCs, A is reachable from this nonterminal and therefore the operation production is part of a cycle.

Character Set Implementation

In real world applications, the occurring character sets usually either contain only a few characters, for example the alphanumericals in a string literal, or almost all characters, for example when a single character is removed from an unknown variable represented by Σ .

We also need to efficiently convert both of these types of sets into short regular expressions to keep the results human-readable. A naive implementation like joining all contained characters using the regular expression choice operator would lead to extremely long expressions for very large sets.

To solve this problem and easily represent these two extremes, we have two different implementations, both conforming to a common **CharSet** interface that requires the common set functions **union**, **intersect**, functionality for adding and removing characters and membership checks.

The first, **SetCharSet**, is mostly a simple wrapper around a **Set<Char>** containing the characters. The second, **SigmaCharSet**, is used to easily represent sets like $\Sigma \setminus \{a, b, c\}$ by storing a **Set<Char>** containing the characters *not* contained in the set, while all other characters are assumed to be members.

The behavior of the the set operations **union** and **intersect** can be described using the operations defined in Figure 5.3.

This approach reduces the storage needed to represent the type of character set, where only a few characters are removed from Σ compared to always storing all contained characters. It also simplifies the creation of regular expressions from the character set. As mentioned above, always using all contained characters in the regular expression produces large regular expressions for sets with cardinality close to $|\Sigma|$. Using our approach, we can represent the regular expression created from a **SigmaCharSet** using negated

5. Approach and Implementation

<code>SigmaCharSet union SigmaCharSet</code>	$\hat{=}$	$(\Sigma \setminus A) \cup (\Sigma \setminus B) = \Sigma \setminus (A \cap B)$
<code>SigmaCharSet union SetCharSet</code>	$\hat{=}$	$(\Sigma \setminus A) \cup S = \Sigma \setminus (A \setminus S)$
<code>SetCharSet union SetCharSet</code>	$\hat{=}$	$S_1 \cup S_2$
<code>SigmaCharSet intersect SigmaCharSet</code>	$\hat{=}$	$(\Sigma \setminus A) \cap (\Sigma \setminus B) = \Sigma \setminus (A \cup B)$
<code>SigmaCharSet intersect SetCharSet</code>	$\hat{=}$	$(\Sigma \setminus A) \cap S = S \setminus A$
<code>SetCharSet intersect SetCharSet</code>	$\hat{=}$	$S_1 \cap S_2$

Fig. (5.3) Definitions of `union` and `intersect` using standard set operations.

character classes and those created from a `SetCharSet` using normal character classes. This reduces the average length of the resulting expressions, compared to always using the same approach. For example the `SetCharSet` that represents the set $\{'a', 'b', 'c'\}$ is used to create the regular expression `[abc]*`, while the `SigmaCharSet` representing $\Sigma \setminus \{'0', '1', '2'\}$ corresponds to `[^012]*`.

5.3.2. Mohri-Nederhof Approximation

Mohri and Nederhof [20] describe an algorithm to approximate a given CFG with an SRG.

Recall from the definition of SRGs in Section 3.2, that we can partition the nonterminals of a grammar into equivalence classes based on whether they are reachable from each other. Also recall that a grammar is strongly regular, if for each such equivalence class, all recursive productions of nonterminals contained in it are either left-linear or right-linear.

For determining if a production rule of a given equivalence class is right- or left-linear, all nonterminals that are not part of the class can be considered as terminals. For example a production $A \rightarrow CX$, where A and C are nonterminals in the same equivalence class, and X is a nonterminal in another class is left linear because X can be viewed as a terminal.

To transform a CFG into an SRG, we only need to transform the sets of mutually recursive nonterminals, where not all productions are either left-linear or right-linear.

Transformation

Mohri and Nederhof describe a more general approach for transforming the required equivalence classes, that accounts for productions with an arbitrary number of nonterminals on the right hand side [20]. Since all productions we use have either one or two nonterminals or exactly one terminal on the right hand side, we can reduce this more general approach to the following algorithm described by Christensen et al. [3].

5. Approach and Implementation

The approach to transform a given equivalence class M consists of the following two steps:

First, for each nonterminal A in M add a new nonterminal A' . Intuitively, A' represents a sequence of characters immediately following the sequence recognized by A . If A is the start terminal of our grammar and therefore corresponds to the analyzed hotspot, add a production $A' \rightarrow \epsilon$.

Second, replace all productions of A according to the replacement rules shown in Figure 5.4. Here B and C are nonterminals in M , X and Y are any nonterminals in a different equivalence class and R is a newly created nonterminal.

$$\begin{array}{lll}
A \rightarrow X & \rightsquigarrow & A \rightarrow X A' \\
A \rightarrow B & \rightsquigarrow & A \rightarrow B, \quad B' \rightarrow A' \\
A \rightarrow X Y & \rightsquigarrow & A \rightarrow R A', \quad R \rightarrow X Y \\
A \rightarrow X B & \rightsquigarrow & A \rightarrow X B, \quad B' \rightarrow A' \\
A \rightarrow B X & \rightsquigarrow & A \rightarrow B, \quad B' \rightarrow X A' \\
A \rightarrow B C & \rightsquigarrow & A \rightarrow B, \quad B' \rightarrow C, \quad C' \rightarrow A' \\
A \rightarrow \text{terminal} & \rightsquigarrow & A \rightarrow R A', \quad R \rightarrow \text{terminal} \\
A \rightarrow op(X) & \rightsquigarrow & A \rightarrow R A', \quad R \rightarrow op(X)
\end{array}$$

Fig. (5.4) Production replacement rules for regular approximation

Since all newly created productions are right-linear, after applying this transformation to all components where it is required, all components in the grammar either contain only left- or only right-linear productions. Therefore, the resulting grammar is strongly regular.

For example, consider a nonterminal A that is part of the currently transformed component and has two productions $A \rightarrow XB$ and $A \rightarrow BX$, again with B being in the same component as A and X in a different component. Here, the first production is right-linear and the second production is left-linear, so the grammar is not strongly regular. Now, these productions are replaced according to the fourth and fifth rules in Figure 5.4, which gives us the following productions for A : $A \rightarrow XB$, $B' \rightarrow A'$, $A \rightarrow B$ and $B' \rightarrow XA'$. The new productions of A are all right-linear and therefore, the grammar can be strongly regular.

Implementation

We can again view a grammar as a directed graph as described in Section 5.3.1.

The notion of mutual “reachability”, by which the relation \mathcal{R} defined in Section 3.2 groups the nonterminals, corresponds to SCCs in this graph view of the grammar.

If two nonterminals A and B are mutually reachable in the graph and therefore part of the same SCC, there exists a sequence of productions to produce B from A and vice versa, which, by definition of \mathcal{R} , means they are in the same equivalence class of \mathcal{R} .

Thus, to approximate a grammar, we view it as a directed graph and find its SCCs, determine the components, where not all productions are of the same linearity and apply the transformation described above to those components.

5.4. Strongly Regular Grammar to Automaton

In the previous section, we obtained a strongly regular grammar. This grammar is now converted into a regular expression by first transforming into an automaton like described in this Section. Section 5.5 then describes how this automaton is converted to a regular expression.

5.4.1. Algorithm

Nederhof describes an algorithm to transform an SRG into an equivalent ϵ -NFA [22]. The generated automaton always accepts the same language as the given grammar.

The full algorithm can be seen in Algorithm 1. It creates an NFA $(K, \Sigma, \Delta, s, F)$ with a set of states K , an alphabet Σ , a set of transitions Δ , an initial state s and a set of accepting states F from a given SRG (Σ, N, P, S) with an alphabet Σ , a set of nonterminals N , a set of productions P and a start nonterminal S .

Note that, for the general algorithm an operation production of form $A \rightarrow op(X)$ is treated like a unary production of form $A \rightarrow X$. The operation productions are always handled by one of the loops in lines 21, 29 or 37, because for any operation production initially contained in a cycle, the cycle is broken up by the character set approximation described in Section 5.3.1. Therefore, e.g. an operation production $C \rightarrow op(D)$ with C and D in the same SCC can no longer occur. We provide a detailed description of resolving the effects of the operations in Section 5.4.2.

The MAKE_FA procedure takes two states q_0 and q_1 and a sequence α of symbols - i.e. terminals and nonterminals - and creates an automaton equivalent to the grammar defined by α between those two states.

5. Approach and Implementation

Algorithm 1 Nederhof Algorithm: $\text{SRG } (\Sigma, N, P, S) \rightarrow \text{NFA } (K, \Sigma, \Delta, s, F)$

```

1: let  $\Delta = \emptyset$ ;  $s = \text{create\_state}()$ ;  $f = \text{create\_state}()$ ;  $F = \{f\}$ ;  $K = \{s, f\}$ 
2:  $\text{MAKE\_FA}(s, S, f)$ 
3: procedure  $\text{MAKE\_FA}(q_0, \alpha, q_1)$ 
4:   if  $\alpha = \epsilon$  then
5:     let  $\Delta = \Delta \cup (q_0, \epsilon, q_1)$   $\triangleright$  add  $\epsilon$  transition from state  $q_0$  to state  $q_1$ 
6:   else if  $\alpha = a$ , some  $a \in \Sigma^*$  then
7:     let  $\Delta = \Delta \cup (q_0, \alpha, q_1)$ 
8:   else if  $\alpha = X\beta$ , some  $X \in V, \beta \in V^*$  such that  $|\beta| > 0$  then
9:     let  $q = \text{create\_state}()$ ;
10:     $K = K \cup \{q\}$   $\triangleright$  create some new state  $q$  and add it to the automaton
11:     $\text{MAKE\_FA}(q_0, X, q)$ 
12:     $\text{MAKE\_FA}(q, X, q_1)$ 
13:   else
14:     let  $A = \alpha$   $\triangleright \alpha$  must be a single nonterminal
15:     if  $A \in N_i$  some  $i$  then
16:       for  $B \in N_i$  do
17:         let  $q_B = \text{create\_state}()$ ;  $K = K \cup \{q_B\}$ 
18:       end for
19:       if  $\text{recursive}(N_i) = \text{left}$  then
20:         for  $(C \rightarrow X_1 \dots X_m) \in P$  such that  $C \in N_i \wedge X_1, \dots, X_m \notin N_i$  do
21:            $\text{MAKE\_FA}(q_0, X_1 \dots X_m, q_C)$ 
22:         end for
23:         for  $(C \rightarrow DX_1 \dots X_m) \in P$  such that  $C, D \in N_i \wedge X_1, \dots, X_m \notin N_i$  do
24:            $\text{MAKE\_FA}(q_D, X_1 \dots X_m, q_C)$ 
25:         end for
26:         let  $\Delta = \Delta \cup (q_A, \epsilon, q_1)$ 
27:       else
28:         for  $(C \rightarrow X_1 \dots X_m) \in P$  such that  $C \in N_i \wedge X_1, \dots, X_m \notin N_i$  do
29:            $\text{MAKE\_FA}(q_C, X_1 \dots X_m, q_1)$ 
30:         end for
31:         for  $(C \rightarrow X_1 \dots X_m D) \in P$  such that  $C, D \in N_i \wedge X_1, \dots, X_m \notin N_i$  do
32:            $\text{MAKE\_FA}(q_C, X_1 \dots X_m, q_D)$ 
33:         end for
34:         let  $\Delta = \Delta \cup (q_0, \epsilon, q_A)$ 
35:       end if
36:     else
37:       for  $(A \rightarrow \beta)$  do  $\triangleright A$  is not recursive
38:          $\text{MAKE\_FA}(q_0, \beta, q_1)$ 
39:       end for
40:     end if
41:   end if
42: end procedure

```

5. Approach and Implementation

This recursive process is started in line 2 with the start nonterminal S , a newly created initial state s as q_0 and a newly created accepting state f as q_1 .

For single terminals and ϵ the algorithm adds an according edge between the two nodes q_0 and q_1 in lines 4 to 7. Here, our implementation differs from the original definition because we allow strings and regular expressions as terminals, whereas usually terminals are single characters. We generalize Nederhof's definition by allowing $a \in \Sigma^*$ instead of just $a \in \Sigma$. This changes the type of the generated NFA because it contains edges labeled with multi-character strings instead of just single symbols. This generalization is possible, because we use the resulting automaton as an input to the state elimination algorithm to create a regular expression in Section 5.5. This algorithm uses a generalized NFA definition with regular expressions - and therefore also strings - as edge labels. If we want to use the NFA directly, we can convert it to a usual NFA by replacing an edge labeled with a string of length n with n edges chained together using new intermediate states.

When α contains multiple symbols, a new state q is created and the automaton for the first symbol in α is inserted between q_0 and q and the one for the rest of α between q and q_1 . Note that for our use case the rest of alpha always contains at most 1 nonterminal since our productions have at most 2 nonterminals on their right hand side.

If α consists of just a single terminal A , that is not part of any set of mutually recursive nonterminals, i.e. from A there is no sequence of productions to reach A again, we just continue the recursion with the right hand sides of A 's productions. The created automaton does not need edges or states corresponding to those single non-recursive nonterminals.

$A \rightarrow B$
 $A \rightarrow C$
 $B \rightarrow b$
 $C \rightarrow c$

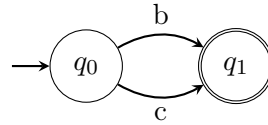


Fig. (5.5) Example grammar with no recursion

Fig. (5.6) Resulting automaton for the grammar in Figure 5.5

To explain why this is the case, consider the grammar in Figure 5.5 creating just the two words “ b ” and “ c ”. Here A , B and C are non-recursive nonterminals, so in the initial procedure call with arguments (q_0, A, q_1) there are just the two recursive calls $\text{MAKE_FA}(q_0, B, q_1)$ and $\text{MAKE_FA}(q_0, C, q_1)$ in line 38. For those calls again the non-recursive case is chosen, such that for the next recursions α equals b or c respectively, which leads to the corresponding edges being created in line 7. As demonstrated, no edges

5. Approach and Implementation

or states are created for any of the 3 nonterminals, only for the two terminals a and b and the resulting automaton in Figure 5.6 accepts the correct language.

For the last remaining case, where α consists of a single nonterminal A that is part of some set of mutually recursive nonterminals N_i , the algorithm first adds a new state for each nonterminal in N_i to the graph.

Then, we differentiate according to the recursion type of N_i , which is obtained by the call to *recursive*(N_i).

Note that sets with neither left nor right recursion can be handled by either case.

Now, for all productions where the left hand side is a nonterminal in N_i , a recursive call depending on the right hand side of the production is performed.

Definition of the Case for Right Recursive Components

Nederhof only defines the case for left recursion in his publication and states that the else part is “the converse of the then part” [22]. This suggests that besides switching the condition for the second loop from $C \rightarrow DX_1 \dots X_m$ to $C \rightarrow X_1 \dots X_mD$, switching the order of the states passed to the recursive calls suffices for handling the right recursive case.

However, only changing e.g. $\text{MAKE_FA}(q_0, X_1 \dots X_m, q_C)$ to $\text{MAKE_FA}(q_C, X_1 \dots X_m, q_0)$ leads to incorrect results. Applying this version of the algorithm to a fully right recursive grammar returns a graph with correct states and correct edges, with the only difference to a correct solution being that the start and the end state are switched. To get correct results, besides switching the argument order, all occurrences of q_0 as an argument to recursive calls need to be replaced with q_1 and vice-versa q_1 with q_0 . Algorithm 1 shows this corrected version of Nederhof’s algorithm.

To explain the differences between the recursive calls in the different cases, consider the grammars in Figures 5.7 and 5.9 and the corresponding automata in Figures 5.8 and 5.10.

The inverting of the states in the recursive calls leads to the edges between states q_2 and q_3 of the automata being inverted, which has no influence on the accepted language.

Switching q_0 and q_1 in the recursive calls is what leads to the needed difference in the resulting automata. In the case of the left recursive grammar, any production sequence of n applications of $B \rightarrow Ab$ has to end with replacing the A on the left hand side of the resulting word with the terminal a to finalize the production rule application. This means that each word has to start with a , which is realized in the automaton by adding an edge labeled with a from q_0 to q_3 due to the recursive call in line 21. For the right recursive grammar conversely, each word has to end with an a due to the bs being generated on the left hand side of the A in $B \rightarrow bA$. Therefore, an edge from q_3 to the final state q_1 is being added by the recursive call in line 29. Accordingly, the corresponding ϵ -edges are added in lines 26 and 34.

$$\begin{aligned} A &\rightarrow a \\ A &\rightarrow B \\ B &\rightarrow Ab \end{aligned}$$

Fig. (5.7) Example grammar with left recursion

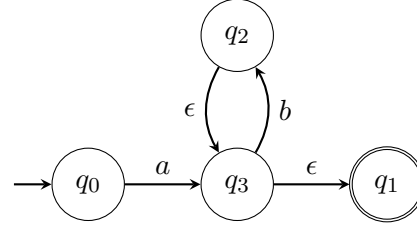


Fig. (5.8) Resulting automaton for the grammar in Figure 5.7

$$\begin{aligned} A &\rightarrow a \\ A &\rightarrow B \\ B &\rightarrow bA \end{aligned}$$

Fig. (5.9) Example grammar with right recursion

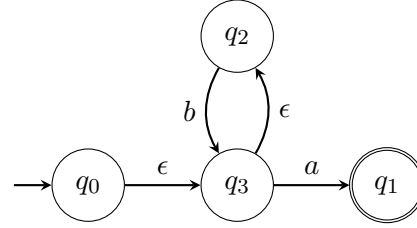


Fig. (5.10) Resulting automaton for the grammar in Figure 5.9

5.4.2. Operation Productions

In the following, we use the two Java operations `reverse` and `replace` as examples for operations with different complexities. Note, however, that we did not implement the complete list of operations on strings the Java standard library contains. To fully support the standard library, one has to define the transformation described in this section for each operation.

As described above, for the Nederhof algorithm, operation productions of form $C \rightarrow op(X)$ are treated like a normal unary production $C \rightarrow X$.

Each operation defines an automaton transformation that changes a given automaton. The new automaton accepts the language obtained by applying the operation to each word in the language of the input automaton. Consider the operation `replace[old, new]`, corresponding to the Java call `s.replace(old, new)`, that returns a copy of the `String s`, where each occurrence of the `char old` is replaced with the `char new`. The automaton transformation for `replace[old, new]` traverses the automaton and replaces each occurrence of `old` on any edge with `new`.

To apply the effect of the different operations onto the created automaton, we first need to find the sub-automata affected by each operation.

To obtain these sub-automata, we taint all nodes and edges if they are created in recursion calls originating from an operation production. If a recursive call of the Nederhof algorithm `MAKE_FA(q_0, X, q_C)` in line 21 is caused by an operation production $C \rightarrow op_1(X)$, we pass op_1 as a taint to the recursive call. All edges and states created further

5. Approach and Implementation

down this recursion path will be tainted with op_1 . In the resulting NFA, for each operation that is part of the given grammar, there's a set of tainted nodes and edges representing the parts of the automaton affected by this operation. These sets each form a sub-automaton of the NFA, onto which the transformation of the corresponding operation can subsequently be applied.

Consider the grammar in Figure 5.11 and the corresponding automaton in Figure 5.12.

$A \rightarrow E$
 $A \rightarrow \text{replace}[f, x](F)$
 $E \rightarrow eE$
 $E \rightarrow e$
 $F \rightarrow fF$
 $F \rightarrow f$

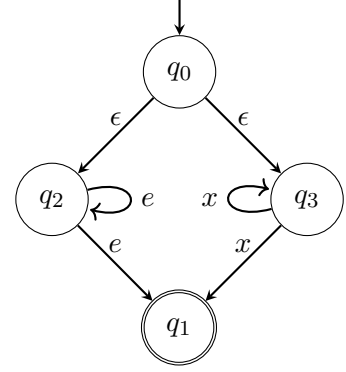
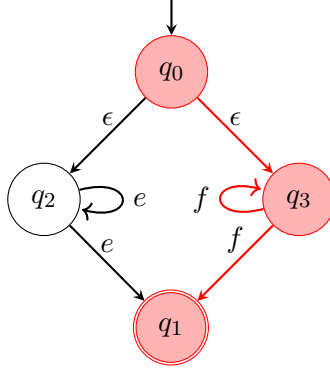


Fig. (5.11) Example grammar with operation production

Fig. (5.12) Resulting automaton for the grammar in Figure 5.11

Fig. (5.13) Automaton in figure 5.12 after applying operation transformation

Here, the production $A \rightarrow E$ leads to the creation of the left path including state q_2 , while for the operation production $A \rightarrow \text{replace}[f, x](F)$ the subsequent algorithm calls create the colored path. All colored edges and states are tainted with the *replace* operation. The created NFA has two similar paths since $A \rightarrow \text{replace}[f, x](F)$ is treated like an $A \rightarrow F$ production analogous to $A \rightarrow E$, just that the resulting edges and adjacent states are tainted.

After completing the NFA creation, we can collect all tainted nodes and apply the automaton transformation defined by $\text{replace}[f, x]$ to this sub-automaton consisting of the states q_0 , q_3 and q_1 .

As mentioned above, for the $\text{replace}[f, x]$ operation this transformation consists of replacing all occurrences of f on tainted edges with x , which gives us the automaton in Figure 5.13 for the given example.

For regular expressions as edge labels the replace operation is more complex.

Take $.*$, for example, which is created when we encounter an unknown value like user input, matches all strings, and therefore also strings containing f . After applying the $\text{replace}[f, x]$ operation to these strings, they can never contain an f , so therefore the edge label should match all strings that contain no f . This can be implemented using a negative character class, so $.*$ is transformed to $[\^f]$ by the $\text{replace}[f, x]$ operation.

5. Approach and Implementation

Similarly, existing character classes need to be transformed, for example $[abf]$ to $[abx]$ and $[\wedge ab]$ to $[\wedge abf]$.

For more complex operation transformations like a *reverse* operation, the operation transformation also includes adding and removing states and edges of the automaton.

Consider the automaton in Figure 5.15, where the colored parts are tainted with the *reverse* operation.

$S \rightarrow A$
 $S \rightarrow reverse(B)$
 $A \rightarrow aA$
 $A \rightarrow c$
 $B \rightarrow bB$
 $B \rightarrow c$

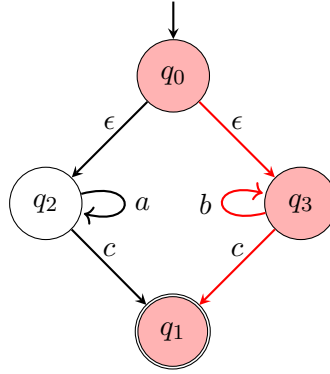


Fig. (5.14) Example grammar with operation production

Fig. (5.15) Resulting automaton for the grammar in figure 5.14

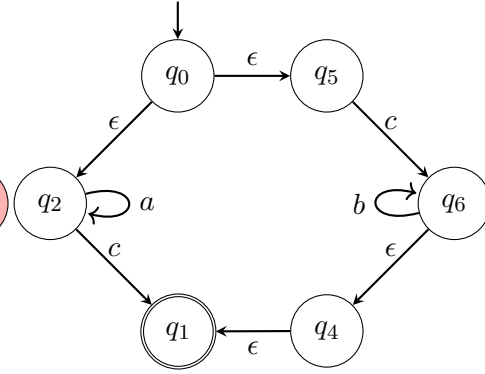


Fig. (5.16) Automaton in figure 5.15 after applying operation transformation

To apply the reverse operation, we first duplicate the tainted sub-automaton and create a new state for each contained state. Here, q_4 is created for q_0 , q_5 is created for q_1 and q_6 is created for q_3 . We also duplicate all tainted edges alongside the states. Now, we reverse the direction of all edges in the duplicated sub-automaton. For example, after duplication there was the edge (q_6, c, q_5) because the original automaton had an edge (q_3, c, q_1) . This edge is now reversed to give us the edge (q_5, c, q_6) that is present in the final automaton.

After reversing all edges, the sub-automaton is connected back to the rest using new ϵ edges. Finally, all tainted edges between the original states are removed together with all states that are not connected to the automaton anymore after this step. In the exemplary automata, after removing the tainted edges, q_3 is disconnected from all other states and therefore removed. The resulting automaton can be seen in Figure 5.16.

5.5. Automaton to Regular Expression

To get a human-readable format for the information we obtained, we transform the automaton we created to a regular expression. Section 5.5.1 describes this conversion while 5.5.2 presents an optimization we used to improve the results.

5.5.1. State Elimination

To transform the automaton we created from an SRG, we use the state elimination strategy, also known as the Brzozowski-McCluskey procedure [1].

In order to apply the procedure, an automaton is first transformed into a generalized nondeterministic finite automaton (GNFA). A GNFA is an NFA where the edges are labeled with regular expressions instead of single symbols. Also, a GNFA must only have a single start state and a single accepting state [14].

To achieve this characteristic, one can add a new start state with a single ϵ transition to the old start state and a new final state with incoming ϵ edges from all previously accepting states.

However, due to the automaton construction using the Nederhof algorithm described above, the automata we obtain already fulfill this property without any need for further modification. We also already use regular expressions as edge labels from the start.

First, we replace each pair of edges $(q_0, r_1, q_1), (q_0, r_2, q_1)$ between two states with a single edge $(q_0, r_1|r_2, q_1)$. After applying this replacement rule exhaustively, there are no two states q_0 and q_1 with more than one direct edge between them.

To eliminate a state q , we “shortcut” the state by replacing all pairs of transitions $(q', r, q), (q, t, q'')$ with a new transition from q' to q'' . The new transition is (q', rt, q'') if q has no loop edge to itself and (q', rs^*t, q'') if it has one with label s [7].

Figures 5.17 and 5.18 contain examples adapted from Esparza [7] that visualize those rules.

After repeatedly applying the two rules and eliminating all other states, the resulting automaton contains only the start and the end state. The single edge between those two states then has the resulting regular expression as a label.

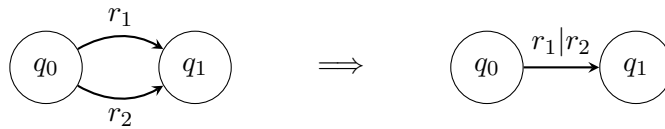
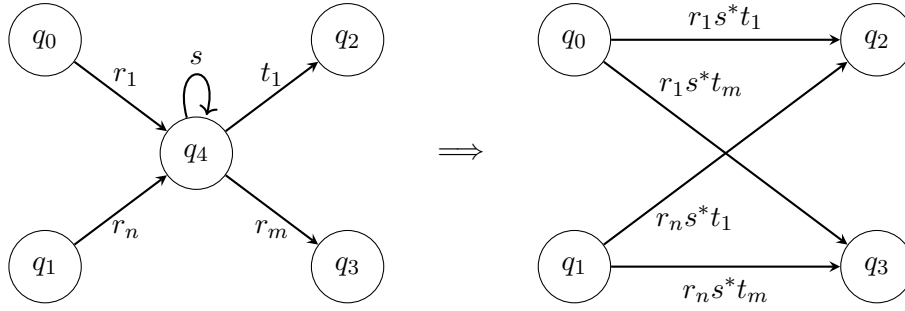


Fig. (5.17) Replacement of edge pairs


 Fig. (5.18) Elimination of state q_4

5.5.2. Delgado Heuristic

The resulting regular expressions often do not have minimal length, but there exists no efficient algorithm that always produces optimal regular expressions.

Minimizing regular expressions is PSPACE-complete [9]. If there were an efficient algorithm to obtain the minimal regular expression for a given NFA, one could first apply Thompson's algorithm [27] for turning a regular expression into an equivalent NFA and then this algorithm. This chaining would then be an efficient algorithm for regular expression minimization algorithm, which contradicts its proven PSPACE-completeness.

However, we can still improve the result we obtain using the state elimination method. The order in which states are eliminated affects the size of the resulting regular expression. There exist different heuristics for choosing an elimination order to reduce the expression size.

We chose a heuristic described by Delgado and Morais [5]. For each state a weight is calculated using Formula 5.3, where In_q is the set of incoming edges of the state q , Out_q the set of outgoing edges, W_e the size of the label on any edge e . Out_q and In_q both do not contain a potential loop on q . W_{loop} is the size of the loop around q if it exists and 0 otherwise.

$$\begin{aligned}
 weight(q) = & \sum_{e \in In_q} (W_e \times (|Out_q| - 1)) + \\
 & \sum_{e \in Out_q} (W_e \times (|In_q| - 1)) + \\
 & W_{loop} \times (|In_q| \times |Out_q| - 1)
 \end{aligned} \tag{5.3}$$

The weight represents the length of the expression added to the result by removing this state. Therefore, in each algorithm run, the state with the smallest weight is chosen for elimination.

5. Approach and Implementation

Delgado and Morais show that using this heuristic produces significantly shorter expressions compared to a naive state elimination [5] with random or undefined ordering. Gruber et al. also show it outperforms almost all other heuristics they considered for their comparison [10]. Improving the algorithm further by implementing a look-ahead additionally to the heuristic also improves the results, but adds more complexity and impairs the algorithm’s performance [5].

Another reduction of the regular expression size can often be obtained by first converting the automaton to a DFA, for example using the powerset construction. For a given NFA with n states, the DFA obtained by using the powerset construction to convert the NFA can have up to 2^n states. However, we observed that for many NFAs generated using Nederhof’s algorithm, the resulting DFAs are significantly smaller than the input NFA or even minimal. This DFA can be minimized using common algorithms like Hopcroft’s or Brzozowski’s algorithm for even better results.

Since it is unclear in which cases conversion to DFAs leads to better results, we can try both approaches for a given query and return the shorter expression to optimize the returned result.

5.6. Hotspot Collection

We also implemented a new pass that traverses the CPG and collects nodes representing string values which might be of interest for further analysis. This hotspot collection provides common starting points for our grammar creation to the user. However, the grammar creation is completely independent of this collection and a grammar can be created for any string node, independent of whether it is part of the collection. We consider all strings that are passed as an SQL query to the Standard Java SQL API as hotspots, as these are the locations where potential SQL injections can occur.

6. Evaluation and Discussion

We first evaluate our approach regarding the obtained results and its performance. Afterward, we discuss its limitations and potential future work to improve and extend our current implementation.

6.1. Evaluation and Benchmarking

In this section, we first analyze the quality of the results of our approach, e.g. concerning the length of the regular expression and whether it describes a correct language that contains all possible values of the analyzed variable in Section 6.1.1. As a metric for the quality of the result, we use the length of the regular expression, because shorter and more concise expressions are usually easier to read and understand for a human user.

Afterward, we measure the execution times of the different steps, including grammar creation, character set and Mohri-Nederhof approximation as well as automaton and regex creation in Section 6.1.2. We analyze the effects of different inputs and variations of our approach.

6.1.1. Quality

We analyze the resulting regular expressions for two synthetic code examples, namely the “Tricky” example from Christensen et al. [3] and one custom example containing simple sanitization logic. We also analyze the results for the SQL injection test cases of the Juliet test suite¹. We choose these examples due to the wide variety of complexity they offer. The Tricky example is specifically crafted to be complex, while the Juliet test cases are comparatively simple. The first therefore shows theoretical limitations of our approach, while the latter two examples are closer to real word applications. Whether these examples are representative of real world applications is discussed in Section 6.2.

Tricky

We adapted the “Tricky” example code Christensen et al. [3] created for their implementation, which can be seen in Listing 6.1. It creates strings of the form

¹<https://samate.nist.gov/SARD/test-suites/111>

6. Evaluation and Discussion

(((((8*7)*6)*5)+4)+3)+2)+1)+0). Since regular languages cannot count the occurrences of characters, a normal regular expression describing those strings can not guarantee, for example, an equal amount of opening and closing brackets. Therefore, a good description using regular languages would, for example, be `\(*<int>(*<int>\\)*(<int>\\)*`, where `<int>` abbreviates the expression `0|(?[1-9][0-9]*)`.

We create the grammar starting at the node representing the variable reference `res` in line 23. After the regular approximation the resulting grammar contains 51 nonterminals and 59 productions. Christensen et al. obtain a different grammar for this example. This difference stems from differences in the definition and implementation of the data flow graph.

The NFA created from the grammar contains 40 states and 51 transitions, of which 39 are ϵ transitions, and can be seen in Figure A.1 of Appendix A. The NFAs created using Nederhof’s algorithm in general often have unnecessary states and transitions, like chains of states only connected by ϵ transitions.

Christensen et al. describe the language they obtained with the expression $\backslash(*\langle\text{int}\rangle([+*]\langle\text{int}\rangle\backslash))*$ [3]. Note that Christensen et al. only create automata and not regular expressions, so this expression is just to describe their created automaton. How their automaton compares to ours is unclear, as they only share this description of the language. We still use their short regular expression as a comparison for our results, both for correctness of the language we accept and size of the expressions.

With a length of 315 characters, the regular expression we obtain is more complex than necessary, but it accepts the same language as the one Christensen et al. gave. The high complexity stems from many optional cases in the expression, which could either be combined into one case or which accept a subset of another option already present in the expression. However, this length is for a sub-optimal scenario, and can be improved as we describe in the following.

Note that our implementation uses the built-in Kotlin functionality to escape strings. This implementation escapes literals by surrounding them with the special characters `\Q` and `\E`, which adds 45 characters compared to escaping using a backslash. To increase readability we use the `<int>` abbreviation and replace the `\Q\E` escape characters with single backslashes in the following regular expressions.

Like mentioned in Section 5.5, converting the NFA into an equivalent DFA significantly improves the result. The corresponding regular expression for this DFA, which can be seen in Figure 6.1, is

$$((\backslash(\backslash(\backslash()*(<\text{int}>|<\text{int}>[*+]))(((<\text{int}>\backslash))[*+])*(((<\text{int}>\backslash))))|((\backslash(\backslash(\backslash()*(<\text{int}>|<\text{int}>)))$$

Minimizing the created DFA gives the automaton in Figure 6.2, which our implementation transforms to the even shorter regular expression $(\backslash()^*\langle\text{int}\rangle([+]\langle\text{int}\rangle\backslash))^*$. In

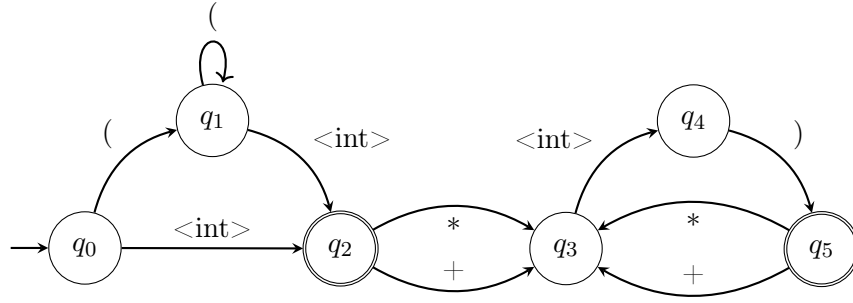


Fig. (6.1) DFA for Tricky example

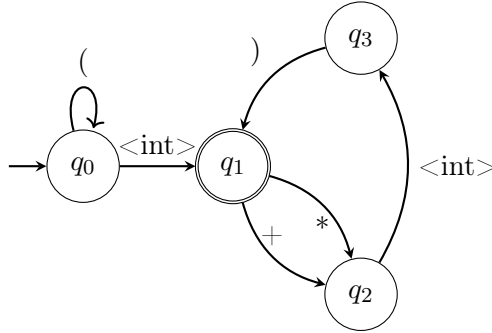


Fig. (6.2) Minimal DFA for Tricky example

essence, this is equivalent to the description by Christensen et al. we mentioned earlier. Therefore, when using some optimizations, our approach can produce very short, human-readable regular expressions even for complex inputs.

The regular expressions mentioned here all accept the same language and just differ in their length and complexity. An overview of the different results and the non-abbreviated regular expressions can be seen in Table A.1.

Note, that neither our nor Christensen et al.'s result does account for the fact, that in the strings generated in the Tricky example all occurrences of $*$ are before the first occurrence of $+$. This is a shortcoming compared to the manually created regular expression mentioned at the beginning of this section. As Christensen et al. mention, this improvement could be achieved by distinguishing the two calls to the `bar` method using a polyvariant analysis [3], which is explained further in Section 6.3.2.

Listing (6.1) Tricky example

```

public class Tricky {
    String bar(int n, int k, String op) {
        if (k==0) {
            return "";
        }
        return op+n+"]"+bar(n-1,k-1,op)+"";
    }
    String foo(int n) {
        String b = "";
        if (n<2) {
            b = b + "(";
        }
        for (int i=0; i<n; i++) {
            b = b + "(";
        }
        String s = bar(n-1,n/2-1,"*");
        String t = bar(n-n/2,n-(n/2-1),"+");
        return b+n+(s+t).replace(']',')');
    }
    public static void main(String args[]) {
        int n = new Random().nextInt();
        String res = new Tricky().foo(n);
        System.out.println(res);
    }
}

```

Listing (6.2) SQL query sanitization example

```

import java.sql.*;
public class DatabaseSanitization {
    public static void main(String[] args) throws SQLException {
        String input = args[1];

        String param = (args[2] == "id") ? "id" : "name";
        String sanitized = sanitize(input);

        Statement stmt = (new Connection()).createStatement();
        stmt.executeQuery(
            "DELETE_*_FROM_users_WHERE_" +
            param + "_=_" + sanitized + "'"
        );
    }

    public static String sanitize(String input) {
        return input.replace('\\', '_').replace('-', '_');
    }
}

```

SQL Query Sanitization

Since the Tricky example is artificially complex, we created the example in Listing 6.2, which resembles a possible real attempt to sanitize an input for an SQL query.

Note that the given code does not compile as `Connection` is an interface and cannot be instantiated, but as this is not relevant for our analysis we ignore it for the sake of brevity.

Our approach returns `(DELETE * FROM users WHERE ((id|name) = '[^'\-]*'))` as a result when analyzing the argument of the `executeQuery` call. This result correctly displays the two possible options `id` and `name` for the parameter.

The transformation of the two `replace` operations also correctly transformed the initial wildcard `.*` corresponding to the unknown `input` to an expression matching any strings that do not contain one of the SQL special characters `'` or `-`. Since an attacker would need to close the opened quote using a single quote character for a successful SQL injection, a further evaluation of this result could show that this sanitization reduces the security risk compared to unfiltered input.

Juliet

The Juliet Test Suite by the National Security Agency's (NSA) Center for Assured Software (CAS) is specifically designed to assess the capabilities of static analysis tools.

6. Evaluation and Discussion

The test cases are grouped by the type of flaw they target, each corresponding to a specific CWE² entry. All 2224 test cases we analyzed target SQL injection vulnerabilities described in CWE-89, as these are flaws, where strings and string operations are the relevant points, while also being relevant risks in real applications.

The Juliet test cases can generally be grouped into ones using bad sinks, where a query string is vulnerable to an SQL injection and ones using good sinks, where prepared statements are used correctly. For the latter, the queries are just literal strings like `"select_*_from_users_where_name=?"`, where the replacement of the question mark with the desired parameter is handled internally by the SQL library.

The vulnerable test cases in the Juliet test suite build an SQL query similar to `"insert_into_users_(status)_values_('updated')_where_name='"+data+"'"`, where `data` is an unsanitized string. The actual semantics of the statements differ, but they all have the same structure where the value is injected as a parameter at the end of the query.

The cases also differ in the data flow from the source of the `data` string and the sink, where it is passed to a database library. Sometimes `data` is unknown and other times a string constant. We are interested in the case, where the value of `data` is unknown, as these cases sometimes pose security vulnerabilities.

For these cases, we obtain regular expressions similar to the example `(insert into users (status) values ('updated') where name='.*')` as a result.

While the interpretation and analysis of the obtained expressions is out of scope for this thesis, it is obvious that this result exposes a severe security risk, because any injection string could be inserted as a value for `data`, signified by the wildcard pattern `.*`.

6.1.2. Performance

Further, we measured the execution time of the different steps in our analysis approach. We measured the steps separately to observe the relative differences between the individual steps and the effects of the optimizations we described. Execution times naturally depend on the machine used to run the analysis. We used a common desktop computer³ for the benchmarks. Like in the previous section, we considered the two synthetic examples and the test cases of the Juliet test suite.

In the following, all mentioned averages are truncated means, which means that before calculating the mean, first the highest $n\%$ and the lowest $n\%$ of values are discarded. This type of truncation is a common method to increase the robustness of the mean, usually with values for n ranging from 10 to 20% [18]. In our benchmarks we use 20% for this cutoff value. This method eliminates the effect of statistical outliers which do

²<https://cwe.mitre.org>

³Intel i7-3770K (8 cores) @ 4.100GHz with 16 GB RAM running Arch Linux 6.1

6. Evaluation and Discussion

not represent an actual measurement but e.g. are influenced by an external factor like CPU scheduling, caching or other unknown influences. For example, when running our benchmarks using JUnit tests, the first test case had the highest durations for all steps including grammar creation, approximation and automaton creation by a factor of 10-20 for some inputs. This effect occurred even when repeatedly testing the same input. Such outliers would skew our measurements and therefore are disregarded for the analysis.

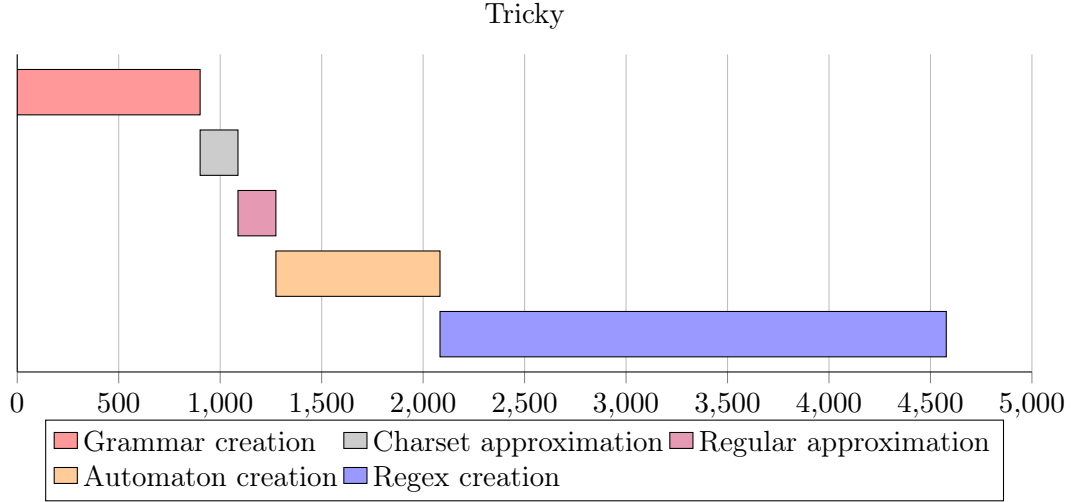


Fig. (6.3) Durations of Tricky example in μs

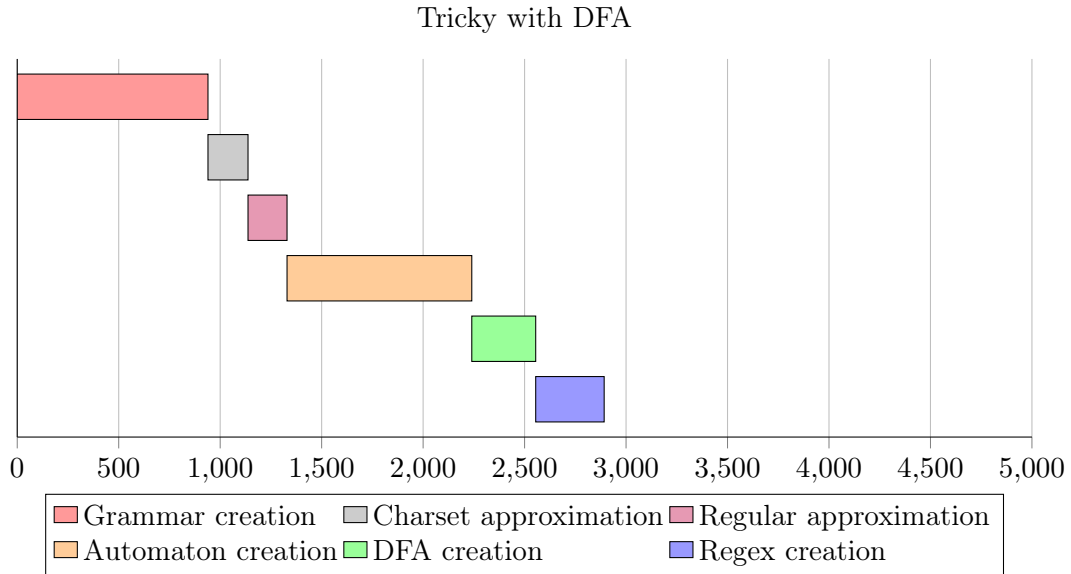


Fig. (6.4) Durations of Tricky example with DFA creation in μs

Tricky

Consider the Figures 6.3 and 6.4. The plots show the durations of each step in the process from CPG to regular expression. For the second plot, we additionally transform the created NFA to a DFA before we create the regular expression, whereas in the first plot, the NFA is converted directly. As visualized by this plot, the DFA creation adds an additional 315 μs to transform the automaton, but drastically reduces the time the state elimination algorithm takes to create a regular expression.

For both plots we averaged the measurements of 100 runs after we trimmed the highest and lowest 20% of values.

SQL Query Sanitization

Figure 6.5 shows the execution time of analyzing the SQL query sanitization example from the previous section. Similar to the Tricky example, the plot shows the average of 100 measurements trimmed by 20%.

Note the different scale of the x axis compared to the previous examples. The difference compared to the other plots is the result of the different complexity of the analyzed code. Since the SQL query sanitization example is considerably less complex than the Tricky example, naturally the analysis execution time is lower.

Figure 6.6 again shows the execution times of the query sanitization example, but with the additional step of creating a DFA before performing the state elimination algorithm.

Due to the characteristics of this example, like a simpler data flow, the resulting NFA already is a DFA. Therefore, the - in this case - unnecessary DFA creation just adds additional time, without significantly changing the execution time of the state elimination algorithm. This shows that whether converting the automaton to a DFA improves the execution time is dependent on properties of the input.

Juliet

We analyzed the execution time of each database related hotspot of all 2224 test cases in the Juliet test suite targeting SQL injection vulnerabilities.

Figure 6.7 shows the averaged execution times of each step, again trimmed by 20%. We can see that the percentage of the total time spent on regex creation is lower compared to the plot for the Tricky example in Figure 6.3, even without the intermediary DFA step. This is due to the fact that the resulting automata for the Juliet test cases are very small, ranging from 2 to 4 states, compared to the 28 nodes of the Tricky NFA.

Again, note the difference in the scale of the x-axis of factor 10, which shows that the Juliet test cases are analyzed significantly faster due to their lower complexity.

6. Evaluation and Discussion

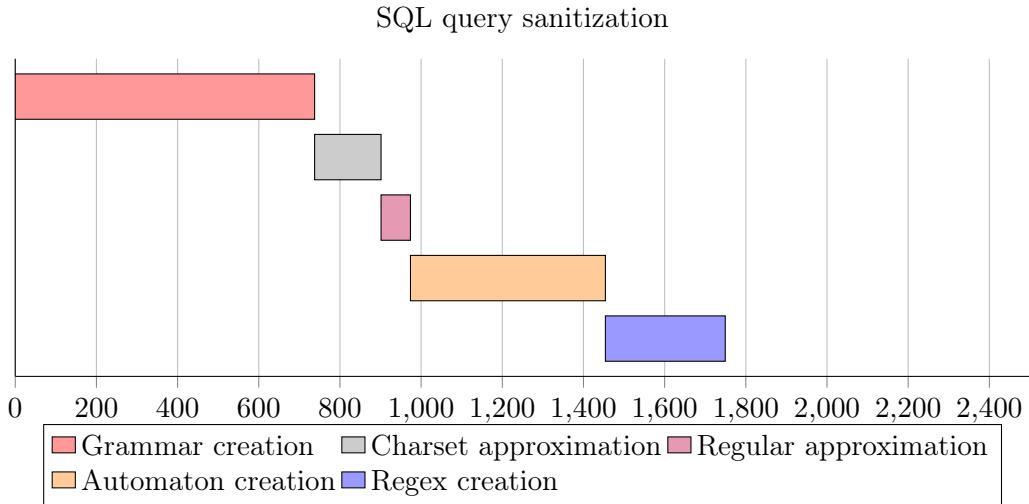


Fig. (6.5) Durations of SQL query sanitization test case in μs

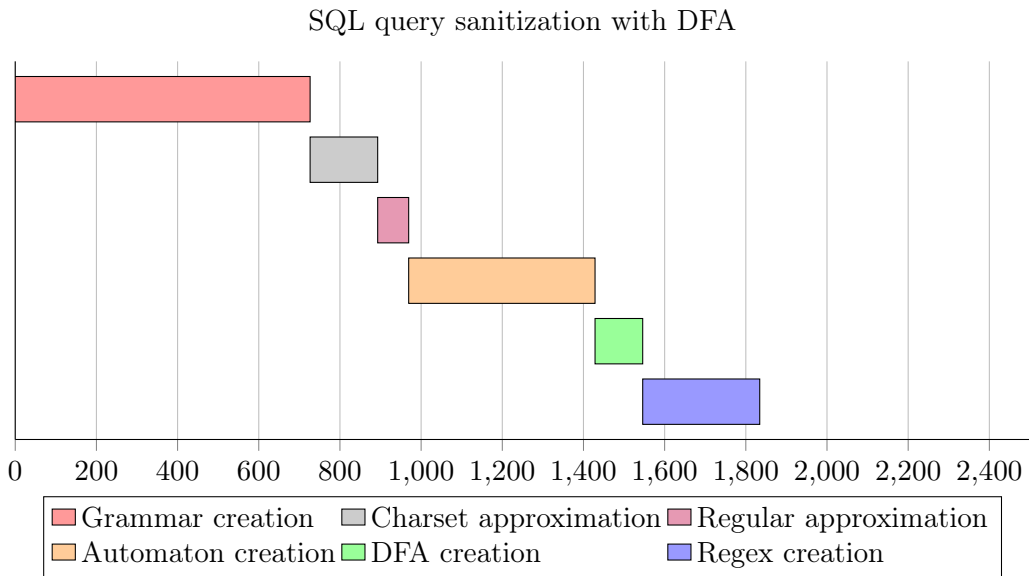
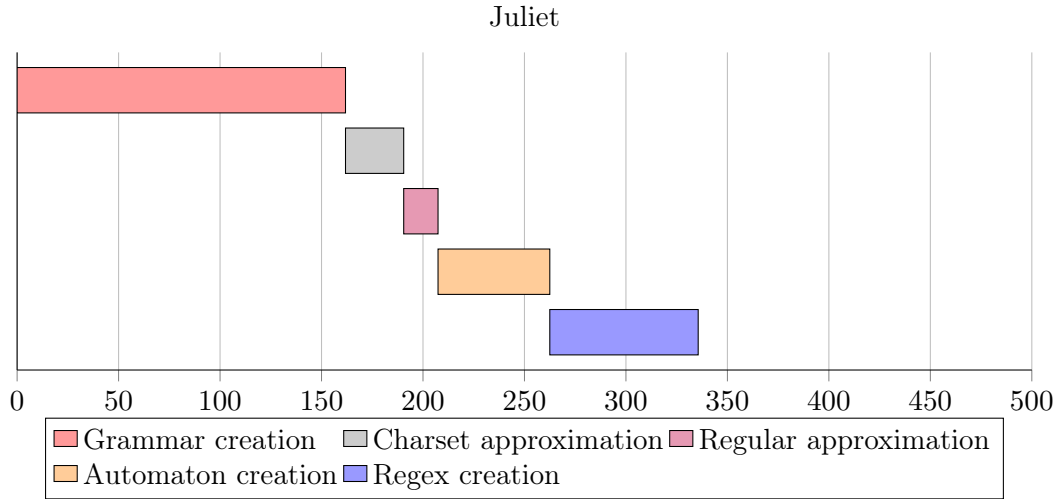


Fig. (6.6) Durations of SQL query sanitization test case with DFA creation in μs

Fig. (6.7) Durations of Juliet test cases in μs

Also note that the relative execution times of the different steps are fairly consistent among all the different examples.

6.2. Limitations of the Evaluation

We have to note that all examples we used to test our implementation are synthetic examples. This is due to the fact that we currently only support a small subset of the Java language, which is not enough to obtain meaningful results from real code. It is unclear how closely the used examples resemble most actual application code and to what extent the observations we made can be transferred to real world use.

However, looking at recent vulnerability reports in the form of CVEs⁴ suggests that comparatively simple SQL injection vulnerabilities still occur in real applications. For example, consider the entry CVE-2022-45932⁵ for a vulnerability in the OpenDaylight⁶ project. Listing 6.3 shows an excerpt of the vulnerable code, where the `deleteRole` function handles a call to an API endpoint and `roleid` is the user input. Here this user input is insufficiently escaped using an HTML escapement utility provided by Apache, that effectively replaces all characters that cannot be displayed in plain HTML. The single quote character however is not changed by this escapement, which allows e.g. the malicious `'_or_1=1--` as a value of `escaped`.

This code has striking similarities to our synthetic example from Listing 6.2, as malicious user input is (insufficiently) sanitized and naively inserted into the query. As we

⁴<https://www.cve.org/>

⁵<https://cve.mitre.org/cgi-bin/cvename.cgi?name=CVE-2022-45932>

⁶<https://www.opendaylight.org/>

demonstrated, our implementation can provide useful information for vulnerabilities of this complexity. This suggests that the approach may be viable for real world applications when extended to support the complete Java standard library.

Listing (6.3) Adapted excerpt from vulnerable code

```
deleteRole(String roleid){                                1
    String escaped = StringEscapeUtils.escapeHtml4(roleid); 2
    String query =                                         3
        String.format(                                     4
            "DELETE FROM " + TABLE + " WHERE " + COL_ID + " = '%s'", 5
            escaped                                           6
        );                                                  7
    stmt.executeUpdate(query);                             8
}                                                            9
```

6.3. Future Work

We discuss several limitations of our implementation and present potential approaches, how future work can mitigate these limitations, and extend our implementation.

6.3.1. Assertions

Our implementation currently does not try to evaluate assertions like `s.isEmpty()` due to limitations in the creation of the CPG we use.

Consider the example in Listing 6.4 where `getSomeKnownValue()` returns some value we can analyze, which is henceforth abbreviated with the generic `<val>`.

Listing (6.4) Assertion Example

```
String s1 = getSomeKnownValue();                        1
if(s2.isEmpty()){                                        2
    s4 = s3 + "empty";                                   3
}                                                        4
System.out.println(s5);                                  5
```

In our CPG the only incoming DFG edge of s^3 in line 3 is an edge from s^1 in line 1. However, there is an implicit information flow from s^2 in line 2 to s^3 , as the result of applying the `isEmpty` operation on s^2 influences the information we can get about s^3 . If there was a DFG edge to the `s2.isEmpty()` call from s^3 instead of the edge from s^1 , we could include the operation in our analysis.

6. Evaluation and Discussion

For example, for such an edge, we could add a new type of production comparable to the existing operation productions, from the nonterminal representing s^3 to the one representing s^2 .

To resolve such an assertion production $A \rightarrow \text{assertion}(B)$ we could implement transformations similar to the existing operation productions. For this example, the transformation of the *isEmpty* assertion would always return just the empty string. In Listing 6.4, we could always infer that s^3 is empty, which is clear from the code.

For this example we currently get the regular expression $(\langle \text{val} \rangle \text{empty}) | (\langle \text{val} \rangle)$ as a result. Consider $\langle \text{val} \rangle$ to be $\text{abc}|\epsilon$, which gives us $((\text{abc}|\epsilon)\text{empty}) | (\text{abc}|\epsilon)$. Here the first part $((\text{abc}|\epsilon)\text{empty})$ corresponds to the value of s^4 , which is a possible value of the analyzed s^5 and the second part $(\text{abc}|\epsilon)$ to the value of s^1 , which is the result if the condition evaluates to false.

With the mentioned additional DFG edges and the described logic, we could sharpen this result. As mentioned above, the value of s^3 would be ϵ due to the *isEmpty* assertion transformation and therefore s^4 would be a concatenation of ϵ and the string "empty", so just *empty*.

This gives us *empty* | $(\text{abc}|\epsilon)$ as a result, which is more precise, as, for example, the unobtainable *abcempty* is not part of the language of this regular expression.

Similar transformations could also be defined for more complex assertions like `s.length() == 1`.

However, as mentioned above, this is currently not possible because the CPG is missing the required DFG edges representing this implicit information flow from an assertion to subsequent variable usages.

Christensen et al. implement a similar feature in their implementation.

6.3.2. Polyvariance

Polyvariance is an analysis strategy, where functions are analyzed more than once, usually once for every call site [24]. The best result we currently obtain for the Tricky example is the regular expression $(\langle \rangle * \langle \text{int} \rangle ([*+] \langle \text{int} \rangle \langle \rangle) *)$, that does not differentiate between $*$ and $+$. Using a polyvariant analysis could improve this, because the two calls to the *bar* function would be analyzed separately with respect to their corresponding arguments.

Currently, we only follow the DFG edges, of which the parameter of the function has one to each variable passed as a parameter, here $+$ and $*$. These two edges are not differentiated and the same result is used for both calls to the function. By differentiating between the corresponding arguments for the analysis of each call we could sharpen the result to $(\langle \rangle * \langle \text{int} \rangle (\langle \rangle * \langle \text{int} \rangle \langle \rangle)) * (\langle \rangle + \langle \text{int} \rangle \langle \rangle) *$.

6.3.3. More Extensive Implementation

We currently do not further analyze values stored in arrays but rather just insert a regular expression generally describing the type stored in the array.

This is due to the fact that the CPG does not contain DFG edges to differentiate which field of an array is accessed in an array subscription expression like `myArray[5]`. A more advanced analysis of arrays, that for example inserts data flow edges from an access to an array at index i to a value that was stored at index i earlier would enable better analysis of these values.

We mostly focused on simple strings, but in general everything we described can be applied to string builders, string buffers, string writers and similar classes. For example, during grammar creation, a `ConcatProduction` could not only be created for `s1 + s2` with two strings, but also for `sb.append(s)` with a `StringBuilder sb` and some other string.

Currently, however, the DFG does not account for method calls on mutable objects. For example, for the mentioned `sb.append` call, there is no DFG edge to the original reference of `sb`, even though the append call changes the value referenced by it. Therefore, the value change is not reflected by the DFG, which makes extracting information hard. A naive approach to obtain this information could include collecting all method calls performed on such an object, ordering them using the evaluation order edges of the CPG, and then evaluating their effects in that order. While extending the DFG to include such edges for operations on mutable objects is possible, one would have to further evaluate which changes to our approach might be needed for it to be able to correctly handle them.

As this is just a proof of concept, we also implemented only a few operations on strings to showcase the approach. For a fully functioning analysis, the other operations included in the Java standard library need to be supported.

6.3.4. State Elimination Heuristics

Moreira et al. analyzed and evaluated different existing heuristics and proposed new ones for choosing a good order to eliminate states during the state elimination algorithm [21].

While they conclude that the method by Delgado and Morais [5] we use is similar to their newly proposed heuristics, they also note that the new strategies outperform Delgado's for small automata. Since the automata we handle are comparatively small, adapting our implementation to use one of the proposed new heuristics could improve the quality of the results.

Moreira et al. also observe that the strategy of taking the best result of using all three heuristics as a final result leads to a gain of 25%. This approach of trying multiple strategies is also used by the Vcsn platform for computations on finite state machines [6].

6.3.5. Automata Centric Approach

We chose to provide the information we extracted only as regular expressions due to regular expressions being widely used and supported. However, representing the information as DFAs instead of converting the automata to regular expressions has some advantages due to theoretical properties of DFAs.

In most programming languages, regular expression objects, for example, Kotlin’s `Regex` object, can determine whether a given string matches the expression. Using a sufficiently advanced automaton implementation more advanced checks can be performed. After analyzing a given hotspot and obtaining an automaton M , instead of just matching a given string as a query, the input can be a regular expression. This regular expression can then be turned into another DFA N . Since DFAs are closed under intersection and complement, we can build the DFA $R = M \cap \overline{N}$, which accepts words that are in the language of the analysis result, but not in the language of the query expression. Now we can check whether R ’s language is empty to determine whether all strings of the query are possible values of the analyzed node. Furthermore, if R ’s language is not empty, we can generate a word from this language as an example for a string that is a possible value of the analyzed node, but not part of the query language. Additionally, we can check whether M and N are equivalent to determine, whether the query expression matches the computed result.

We currently focus on producing human readable results in the form of regular expressions. For a fully automated analysis however, these advanced query possibilities could be leveraged to produce better feedback for users. Since we use NFAs, or optionally already DFAs, as an intermediate representation, future work could increase the capabilities of our automata implementation and implement the mentioned features.

7. Conclusion

In this thesis, we presented a method to obtain information about the values of strings from the data flow graph of an analyzed program. We provide a proof of concept implementation as an extension of an existing Code Property Graph implementation used for static analysis.

We adapted part of an existing approach to obtain a strongly regular grammar from the graph. We then converted this grammar into an automaton using an algorithm we adapted and extended for our use case. Further, this automaton is transformed into a regular expression, which describes the analyzed string. We used approximations of different precision to model the effects of concatenation and other operations on strings.

Additionally, we described different methods, like intermediate conversion to a DFA and a heuristic for state elimination, to potentially increase the performance of our implementation.

We also showed that, even for a complex example, our implementation provides useful results which could be used to detect security vulnerabilities like SQL injections. Furthermore, we tested and benchmarked our implementation using different examples and the well known Juliet test suite, which showed the viability of our performance optimizations and general approach. Moreover, we summarized limitations of our implementation and provided starting points for potential further research.

The information our approach, especially with further enhancements to our implementation, provides can be beneficial for static analysis, especially in the context of preventing common security vulnerabilities.

A. Evaluation Results

Tricky Results

Method	Automaton Size (states, transitions)	Expression length	Regular expression
NFA	40, 51	315	$(((\backslash Q \backslash E)? \backslash Q \backslash E) * \backslash Q \backslash E \backslash Q \backslash E)? (0 (-? [1-9] [0-9] *))) (([++]) (0 (-? [1-9] [0-9] *)))$ $) ((\backslash Q) \backslash E [++]) (0 (-? [1-9] [0-9] *))) * (((\backslash Q) \backslash E [++]) (0 (-? [1-9] [0-9] *))) ((\backslash Q) \backslash E [++]) (0 (-? [1-9] [0-9] *))) * \backslash Q \backslash E \backslash Q \backslash E) $ $(((\backslash Q \backslash E)? \backslash Q \backslash E) * \backslash Q \backslash E \backslash Q \backslash E)? (0 (-? [1-9] [0-9] *))) (([++]) (0 (-? [1-9] [0-9] *))) * (((\backslash Q) \backslash E [++]) (0 (-? [1-9] [0-9] *))) * \backslash Q \backslash E)?$
DFA	6, 10	176	$(((\backslash Q \backslash E \backslash Q \backslash E) * (0 (-? [1-9] [0-9] *))) 0 (-? [1-9] [0-9] *))) [++] ((0 (-? [1-9] [0-9] *))) \backslash Q \backslash E [++]) * (((0 (-? [1-9] [0-9] *))) \backslash Q \backslash E) $ $(((\backslash Q \backslash E \backslash Q \backslash E) * (0 (-? [1-9] [0-9] *))) 0 (-? [1-9] [0-9] *)))$
minimized DFA	4, 6	62	$((\backslash Q \backslash E) * (0 (-? [1-9] [0-9] *))) ([++] ((0 (-? [1-9] [0-9] *))) \backslash Q \backslash E)) *$

Table (A.1) Comparison of different approaches to obtain regular expression for Tricky example

A. Evaluation Results

Tricky NFA

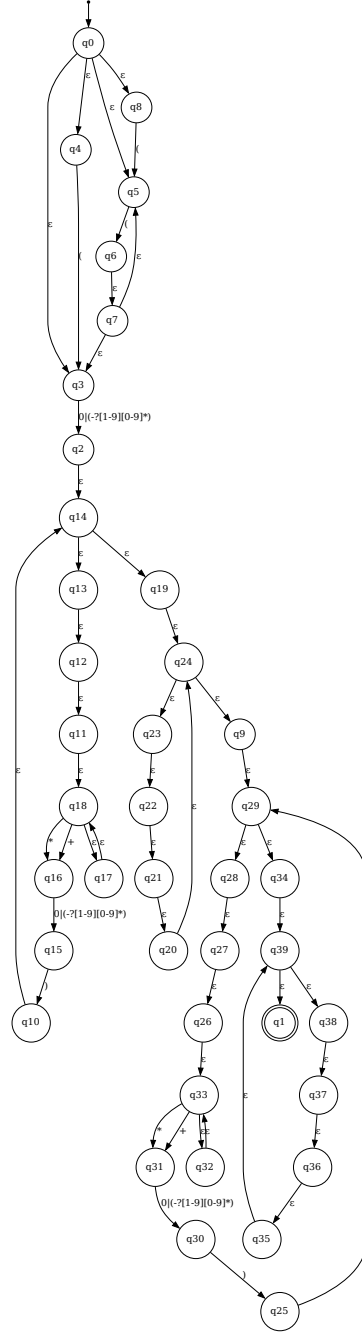


Fig. (A.1) NFA created from the Tricky example from Section 6.1.1

Bibliography

- [1] J. A. Brzozowski and E. J. McCluskey. “Signal Flow Graph Techniques for Sequential Circuit State Diagrams.” In: *IEEE Transactions on Electronic Computers* EC-12.2 (1963), pp. 67–76. DOI: 10.1109/PGEC.1963.263416.
- [2] C. Câmpeanu, K. Salomaa, and S. Yu. “A formal study of practical regular expressions.” In: *International Journal of Foundations of Computer Science* 14.06 (2003), pp. 1007–1018.
- [3] A. S. Christensen, A. Møller, and M. I. Schwartzbach. “Precise Analysis of String Expressions.” In: *Proc. 10th International Static Analysis Symposium (SAS)*. Vol. 2694. LNCS. Available from <http://www.brics.dk/JSA/>. Springer-Verlag, June 2003, pp. 1–18.
- [4] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to algorithms*. MIT press, 2022.
- [5] M. Delgado and J. Morais. “Approximation to the smallest regular expression for a given regular language.” In: *Implementation and Application of Automata: 9th International Conference, CIAA 2004, Kingston, Canada, July 22-24, 2004, Revised Selected Papers 9*. Springer. 2005, pp. 312–314.
- [6] A. Demaille, A. Duret-Lutz, S. Lombardy, and J. Sakarovitch. “Implementation Concepts in Vaucanson 2.” In: *Proceedings of Implementation and Application of Automata, 18th International Conference (CIAA’13)*. Ed. by S. Konstantinidis. Vol. 7982. Lecture Notes in Computer Science. Halifax, NS, Canada: Springer, July 2013, pp. 122–133. ISBN: 978-3-642-39274-0. DOI: 10.1007/978-3-642-39274-0_12.
- [7] J. Esparza. “Automata theory – An algorithmic approach.” Lecture Notes, <https://www7.in.tum.de/~esparza/autoskript.pdf>. Aug. 2017.
- [8] C. Gould, Z. Su, and P. Devanbu. “Static checking of dynamically generated queries in database applications.” In: *Proceedings. 26th International Conference on Software Engineering*. IEEE. 2004, pp. 645–654.
- [9] G. Gramlich and G. Schnitger. “Minimizing nfa’s and regular expressions.” In: *Journal of Computer and System Sciences* 73.6 (2007), pp. 908–923.
- [10] H. Gruber, M. Holzer, and M. Tautschnig. “Short regular expressions from finite automata: Empirical results.” In: *Implementation and Application of Automata: 14th International Conference, CIAA 2009, Sydney, Australia, July 14-17, 2009. Proceedings 14*. Springer. 2009, pp. 188–197.

Bibliography

- [11] W. G. J. Halfond and A. Orso. “AMNESIA: Analysis and Monitoring for NEutralizing SQL-Injection Attacks.” In: *Proceedings of the 20th IEEE/ACM International Conference on Automated Software Engineering*. ASE ’05. Long Beach, CA, USA: Association for Computing Machinery, 2005, pp. 174–183. ISBN: 1581139934. DOI: 10.1145/1101908.1101935.
- [12] W. G. J. Halfond and A. Orso. “Preventing SQL Injection Attacks Using AMNESIA.” In: *Proceedings of the 28th International Conference on Software Engineering*. ICSE ’06. Shanghai, China: Association for Computing Machinery, 2006, pp. 795–798. ISBN: 1595933751. DOI: 10.1145/1134285.1134416.
- [13] W. G. J. Halfond, A. Orso, and P. Manolios. “Using Positive Tainting and Syntax-Aware Evaluation to Counter SQL Injection Attacks.” In: *Proceedings of the 14th ACM SIGSOFT International Symposium on Foundations of Software Engineering*. SIGSOFT ’06/FSE-14. Portland, Oregon, USA: Association for Computing Machinery, 2006, pp. 175–185. ISBN: 1595934685. DOI: 10.1145/1181775.1181797.
- [14] Y.-S. Han and D. Wood. “The generalization of generalized automata: Expression automata.” In: *International Journal of Foundations of Computer Science* 16.03 (2005), pp. 499–510.
- [15] H. Hosoya and B. C. Pierce. “XDuce: A statically typed XML processing language.” In: *ACM Transactions on Internet Technology (TOIT)* 3.2 (2003), pp. 117–148.
- [16] C. Kirkegaard, A. Møller, and M. I. Schwartzbach. “Static Analysis of XML Transformations in Java.” In: *IEEE Transactions on Software Engineering* 30.3 (Mar. 2004), pp. 181–192.
- [17] J. Kodumal and A. Aiken. “Banshee: A Scalable Constraint-Based Analysis Toolkit.” In: *Static Analysis*. Ed. by C. Hankin and I. Siveroni. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 218–234. ISBN: 978-3-540-31971-9.
- [18] U. Krenel. *Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Vol. 8. Wiesbaden: Vieweg, 2005, p. 171. ISBN: 3-8348-0063-5. DOI: 10.1007/978-3-663-09885-0.
- [19] V. B. Livshits and M. S. Lam. “Finding Security Vulnerabilities in Java Applications with Static Analysis.” In: *USENIX security symposium*. Vol. 14. 2005, pp. 18–18.
- [20] M. Mohri and M.-J. Nederhof. “Regular approximation of context-free grammars through transformation.” In: *Robustness in language and speech technology*. Springer, 2001, pp. 153–163.
- [21] N. Moreira, D. Nabais, and R. Reis. “State elimination ordering strategies: Some experimental results.” In: *arXiv preprint arXiv:1008.1656* (2010).
- [22] M.-J. Nederhof. “Regular approximation of CFLs: a grammatical view.” In: *Advances in Probabilistic and other Parsing Technologies*. Springer, 2000, pp. 221–241.
- [23] Open Worldwide Application Security Project (OWASP). *OWASP Top 10*. 2021. URL: <https://owasp.org/Top10/> (visited on 02/25/2023).

Bibliography

- [24] J. PALSBERG and C. PAVLOPOULOU. “From Polyvariant flow information to intersection and union types.” In: *Journal of Functional Programming* 11.3 (2001), pp. 263–317. DOI: 10.1017/S095679680100394X.
- [25] N. Tabuchi, E. Sumii, and A. Yonezawa. “Regular Expression Types for Strings in a Text Processing Language.” In: *Electronic Notes in Theoretical Computer Science* 75 (2003). TIP’02, International Workshop in Types in Programming, pp. 95–113. ISSN: 1571-0661. DOI: [https://doi.org/10.1016/S1571-0661\(04\)80781-3](https://doi.org/10.1016/S1571-0661(04)80781-3).
- [26] R. Tarjan. “Depth-First Search and Linear Graph Algorithms.” In: *SIAM Journal on Computing* 1.2 (1972), pp. 146–160. DOI: 10.1137/0201010. eprint: <https://doi.org/10.1137/0201010>.
- [27] K. Thompson. “Programming Techniques: Regular Expression Search Algorithm.” In: *Commun. ACM* 11.6 (June 1968), pp. 419–422. ISSN: 0001-0782. DOI: 10.1145/363347.363387.
- [28] G. Wassermann and Z. Su. “Sound and precise analysis of web applications for injection vulnerabilities.” In: *ACM-SIGPLAN Symposium on Programming Language Design and Implementation*. 2007.
- [29] K. Weiss and C. Banse. *A Language-Independent Analysis Platform for Source Code*. 2022. DOI: 10.48550/ARXIV.2203.08424.

List of Figures

3.1. Part of the CPG for the code in Listing 3.1	8
5.1. The general approach for obtaining regular expressions	14
5.2. Example grammar	17
5.3. Definitions of union and intersect using standard set operations.	20
5.4. Production replacement rules for regular approximation	21
5.5. Example grammar with no recursion	24
5.6. Resulting automaton for the grammar in Figure 5.5	24
5.7. Example grammar with left recursion	26
5.8. Resulting automaton for the grammar in Figure 5.7	26
5.9. Example grammar with right recursion	26
5.10. Resulting automaton for the grammar in Figure 5.9	26
5.11. Example grammar with operation production	27
5.12. Resulting automaton for the grammar in Figure 5.11	27
5.13. Automaton in figure 5.12 after applying operation transformation	27
5.14. Example grammar with operation production	28
5.15. Resulting automaton for the grammar in figure 5.14	28
5.16. Automaton in figure 5.15 after applying operation transformation	28
5.17. Replacement of edge pairs	29
5.18. Elimination of state q_4	30
6.1. DFA for Tricky example	34
6.2. Minimal DFA for Tricky example	34
6.3. Durations of Tricky example in μs	38
6.4. Durations of Tricky example with DFA creation in μs	38
6.5. Durations of SQL query sanitization test case in μs	40
6.6. Durations of SQL query sanitization test case with DFA creation in μs	40
6.7. Durations of Juliet test cases in μs	41
A.1. NFA created from the Tricky example from Section 6.1.1	48

List of Tables

3.1. The regular expression metacharacters we use	7
4.1. Comparison of the presented related work	13
5.1. Utilized production types	15
A.1. Comparison of different approaches to obtain regular expression for Tricky example	47

List of Listings

3.1. CPG example code	8
5.1. Grammar creation example code	15
6.1. Tricky example	35
6.2. SQL query sanitization example	36
6.3. Adapted excerpt from vulnerable code	42
6.4. Assertion Example	42

List of Acronyms

AST abstract syntax tree	8
CFG context-free grammar	5
CFL context-free language	10
CPG Code Property Graph	2
DAG directed acyclic graph	18
DFA deterministic finite automaton	2
DFG data flow graph	2
GNFA generalized nondeterministic finite automaton	29
MLFA multi-level automaton	9
NFA nondeterministic finite automaton	3
PQL Process Query Language	11
SCC strongly connected component	17
SRG strongly regular grammar	14