



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

Thanks everyone!

Abstract

In the last couple of years, I have supervized numerous bachelor's and master's thesis and various seminars. This led to a broad observation of typical questions and issues the students faced when writing their thesis or papers. Surprisingly, they are always quite similar. This template aims to give advise to future sstudents in order to answer the most frequent questions and avoid the most common mistakes. It provides the TUM template which has already been accepted many times, shows the most basic outline and some tips on the contents of each chapter. It further contains some tips on the style of scientific works. An evaluation on a small set of students showed that this guideline can assist in making progress faster. However, we found that we have to keep improving the tips to achieve better results.

Your abstract goes here. The typical structure is:

- Broad description of the current state
- Gap in the current state
- Your contribution

Contents

Acknowledgments	i
Abstract	ii
1 Introduction	1
2 Problem Description	2
3 Background	3
4 Approach and Implementation	5
4.1 Hotspot Collection	5
4.2 Grammar Creation	5
4.3 Regular Approximation	7
4.3.1 Character Set Approximation	7
4.3.2 Mohri-Nederhof Approximation	9
4.4 Transformation to Regular Expression	11
4.4.1 Strongly Regular Grammar to Automaton	11
4.4.2 Automaton to Regular Expression	14
5 Evaluation and Discussion	18
6 Related Work	19
7 Conclusion	21
Bibliography	22

1 Introduction

Your introduction goes here

- Generic description of the broad field of research
- Current state of research
- What's the gap that you're trying to fill?
- Short motivation
- Summary of the most important results
- Your contribution
- Structure of the thesis

1-2.5 pages

This text is not too detailed. Start quite high-level, then narrow down until you reach your topic. After the introduction, the reader must want to read the rest of your thesis and understand the relevance. However, it doesn't have to be super technical.

2 Problem Description

The introduction is a bit like a teaser. Here, you dig more into details, also technical ones. After this chapter, the reader must understand why you do this work, why it's important, what makes it difficult and what you want to achieve.

- What's the problem that you're trying to solve?
- What is your goal?
- What is/are the research question(s)?
- What are special problems?

Probably 1-3 pages

3 Background

The library¹ we extend in this thesis extracts a Code Property Graph (CPG) out of source code of a set of different programming languages.

The 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. The 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 by multiple passes [12].

Users of the library can extend this functionality by adding additional passes, which is how we implement the hotspot collection in this thesis.

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

```
String s = "xyz";  
System.out.println(s);
```

Listing 3.1: Example code

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

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

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

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

3 Background

The nodes connected by those edges effectively form a subgraph of the CPG, the data flow graph (DFG), from which we then extract the information on string values.

4 Approach and Implementation

The general approach for our implementation is adapted from the one described by Christensen et al. [2]. Conceptually, we first extract a context free grammar (CFG) from the DFG. In multiple steps using different methods we then approximate this grammar into a strongly regular grammar. From this grammar we can create a regular expression object to provide to the user for further analysis.

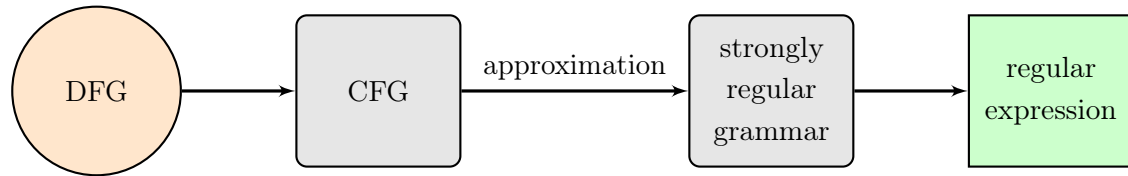


Fig. 4.1: The general approach for obtaining regular expressions

4.1 Hotspot Collection

We 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 includes all strings that are passed as a query to the Java SQL library and all strings in return statements.

4.2 Grammar Creation

To create the grammar for a given CPG node, we traverse the DFG backwards, starting at the given node. The starting node can be one of the hotspot nodes collected by the aforementioned `Pass`, but in general the grammar creation is independent of the hotspot collection. For each visited node, we add a `Nonterminal` and the fitting productions to our grammar.

Our Grammar contains the following five types of productions:

- **UnitProduction:** $X \rightarrow Y$ for references between nodes
- **ConcatProduction:** $X \rightarrow Y Z$ for concatenation of two nodes

4 Approach and Implementation

- **TerminalProduction**: $X \rightarrow \langle \text{terminal} \rangle$ for literal string values and other terminal symbols
- **UnaryOpProduction**: $X \rightarrow op(Y)$ for unary operations on strings
- **BinaryOpProduction**: $X \rightarrow op(Y, Z)$ for binary operations on strings

Here $\langle \text{terminal} \rangle$ 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.

```

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

Listing 4.1: Example code

Consider the code example in Listing 4.1 for the following explanations of the different productions.

UnitProductions mostly represent references between nodes where the underlying string is not changed. In Listing 4.1 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 string concatenation using the + operator. For the example in Listing 4.1 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 "abc". **TerminalProductions** also occur at CPG 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 lanuage `.*`, matching all strings.

UnaryOpProductions and **BinaryOpProductions** represent function calls or other operators. The CPG for 4.1 contains a **CallExpression** representing the function call of the library function `trim`. We then create an **Operation** object representing this operation and the **UnaryOpProduction** $X \rightarrow trim(Y)$, where X is the nonterminal corresponding to the node representing s^4 and Y to the one representing s^5 . The **Operation** objects also contain information about possible arguments and implement the character set transformation and regular approximation needed for the approximation of the grammar described

in Section 4.3. This language agnostic representation of string operation 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**.

Improvements

Unlike Christensen et al. [2], we do not consider the total DFG when extracting the grammar. While they parse the whole graph into a data structure, to later extract automata for specific nodes, 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 often the majority of a large program is not relevant for a specific node, this reduces the amount of nodes we need to handle and 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, it uses a **ValueEvaluator** to try, whether the value the node represents is known. 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 integrals.

4.3 Regular Approximation

4.3.1 Character Set Approximation

To use the Mohri-Nederhof approximation algorithm [7], we need to eliminate all cycles in our grammar that contain operation productions. All nonterminals are assigned a character set, containing all characters that make up the words in the language of the corresponding nonterminal. 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 char `o` is replaced by a known char `n` has the following character set transformation, whereas for a **replace** operation, where the newly inserted char is not known, S is transformed to Σ if the replaced char is contained in S .

4 Approach and Implementation

$$T_{replace[o,n]}(S) = \begin{cases} (S \setminus \{o\}) \cup \{n\}, & \text{if } o \in S \\ S, & \text{if } o \notin S \end{cases} \quad (4.1)$$

These approximations, together with the terminals where the character set is known, for example a string literal, can be used in a fixed point computation to assign a character set $C(X)$ to each nonterminal X .

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 $C(X)^*$.

We find those operation cycles by viewing the grammar as a graph and determining the strongly connected components (SCCs) of this graph. Now for each nonterminal N in a given component C , we check, whether it has an operation production, and if yes, whether one of the nonterminals on its right-hand side is also part of C . If this is the case, by definition of SCCs, N is reachable from this nonterminal and therefore the operation production is part of a cycle.

To determine the SCCs, we use Tarjan's algorithm [10]. This algorithm topologically sorts the returned components in reverse order, which is necessary for the fixpoint computation used to find the charsets. During the computation, for a given nonterminal N , its charset is updated using the charsets 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 been handled earlier.

To represent character sets easily, we have two different implementations, both conforming to a common `CharSet` interface that requires functions like `union : CharSet -> CharSet` and `intersect : CharSet -> CharSet`.

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 following set operations:

4 Approach and Implementation

$$\begin{array}{lll}
\text{SigmaCharSet union SigmaCharSet} & \hat{=} (\Sigma \setminus A) \cup (\Sigma \setminus B) = \Sigma \setminus (A \cap B) \\
\text{SigmaCharSet union SetCharSet} & \hat{=} (\Sigma \setminus A) \cup S & = \Sigma \setminus (A \setminus S) \\
\text{SetCharSet union SetCharSet} & \hat{=} & S_1 \cup S_2 \\
\text{SigmaCharSet intersect SigmaCharSet} & \hat{=} (\Sigma \setminus A) \cap (\Sigma \setminus B) = \Sigma \setminus (A \cup B) \\
\text{SigmaCharSet intersect SetCharSet} & \hat{=} (\Sigma \setminus A) \cap S & = S \setminus A \\
\text{SetCharSet intersect SetCharSet} & \hat{=} & S_1 \cap S_2
\end{array}$$

This approach reduces the storage needed to represent the commonly occurring type of character sets, where only a few characters are removed from Σ . It also simplifies the creation of a regular expression from the character set, since the approach of using a character class containing all characters in the set produces very large character classes for sets with cardinality close to $|\Sigma|$. Using our approach, we can represent a **SigmaCharSet** using negated character classes. Since most character sets either contain a comparatively small amount of given chars, or all chars except a few this reduces the average length of the resulting regular expressions. For example the **SetCharSet** that represents the set $\{'a', 'b', 'c'\}$ gives us the regular expression $[abc]^*$, while the **SigmaCharSet** representing $\Sigma \setminus \{'0', '1', '2'\}$ corresponds to $[^012]^*$.

4.3.2 Mohri-Nederhof Approximation

Strongly Regular Grammars

Mohri and Nederhof [7] describe an algorithm to transform a CFG into a strongly regular grammar that approximates the given CFG.

They define strongly regular grammars as follows:

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

$$ARB \Leftrightarrow (\exists \alpha, \beta \in V^* : A \xrightarrow{*} \alpha B \beta) \wedge (\exists \alpha, \beta \in V^* : B \xrightarrow{*} \alpha A \beta) \quad (4.2)$$

Here V is $\Sigma \cup N$, so 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. $A \xrightarrow{*} \alpha B \beta$ means, that there exists a sequence of productions starting at the symbol A to produce a set of symbols that contain B . Therefore \mathcal{R} groups all nonterminals into disjoint equivalence classes, where each nonterminal in a class can be produced by each other nonterminal in the class. Those nonterminals are called mutually recursive.

A grammar is strongly regular if the production rules in each such equivalence class are either all right-linear or left-linear.

4 Approach and Implementation

A production rule is right-linear if 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 nonterminal is on the left side of the production result.

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.

Therefore, to transform a CFG into a strongly regular grammar, 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 transformation approach for productions with an arbitrary number of nonterminals on the left hand side [7]. 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 set of rules described by Christensen et al.[2].

For each nonterminal A in a given equivalence class M add a new nonterminal A' .

Replace all productions of A with the following new productions, where 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.

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

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.

Implementation

We can view a grammar as a directed graph, with the nonterminals as nodes and an edge from a node A to a node B iff there is a production with A on its left-hand side and B contained in its right-hand side, so a production of form $A \rightarrow \alpha B \beta$.

The aforementioned notion of mutual "reachability", by which \mathcal{R} 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 is 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 mentioned above to those components.

4.4 Transformation to Regular Expression

4.4.1 Strongly Regular Grammar to Automaton

Algorithm

Nederhof describes an algorithm to transform a strongly regular grammar into an equivalent nondeterministic finite automaton (NFA) in [8]. More specifically, the algorithm creates an ϵ -NFA. The generated automaton always accepts the same language as the given grammar.

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

The `create_state` function used in the pseudo code just creates a new state object which can then be added to the automaton.

Note that for the general algorithm an operation production of form $A \rightarrow op(X)$ is treated like an 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 4.3.1. Therefore an operation production $C \rightarrow op(DX_1)$ with C and D in the same SCC can no longer occur. How the effects of the operation productions are resolved is described in the following section.

The `MAKE_FA` procedure takes two states q_0 and q_1 and a sequence α of symbols - terminals and nonterminals - and creates an automaton equivalent to the grammar

4 Approach and Implementation

starting at α between those two states.

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

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.

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 .

If α consists of just a single terminal A , that is not part of any set of mutually recursive nonterminals, so 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. 4.2: Example grammar with no recursion

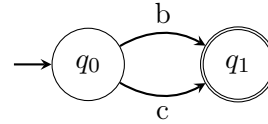


Fig. 4.3: Resulting automaton for the grammar in Figure 4.2

To show why this is the case, consider the grammar in Figure 4.2 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 the corresponding edges being created in line 7. As demonstrated no edges or states are created for any of the 3 nonterminals, only for the two terminals a and b and the resulting automaton in Figure 4.3 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 $\text{recursive}(N_i)$.

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. To explain the differences between the recursive calls in the different cases consider the grammars in

4 Approach and Implementation

Figures 4.5 and 4.7.

Nederhof only defines the case for left recursion in his publication and states that the else part is "the converse of the then part" [8]. 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. $\text{MAKE_FA}(q_0, X_1 \dots X_m, Q_C)$ 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 .

The inverting of the states in the recursive calls leads to the edges between states q_2 and q_3 of the automata in Figures 4.6 and 4.8 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 finale 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.

Operation Productions

As described above, operation productions of form $C \rightarrow op(X)$ are treated like a normal production $C \rightarrow X$. To apply the effect of the different operations onto the created automaton, we taint all nodes and edges if they are created in recursion calls after an operation production.

If a recursive call $\text{MAKE_FA}(q_0, X, q_C)$ 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 down this recursion path will be tainted with op_1 . This allows us to apply the automaton-transformation required to apply the effect of an operation op_1 to the sub-automaton tainted with op_1 .

Consider the grammar in Figure 4.8 and the corresponding automaton in Figure 4.9. Here the production $A \rightarrow F$ leads to the creation of the left path including state q_2 , while

$A \rightarrow a$
 $A \rightarrow B$
 $B \rightarrow Ab$

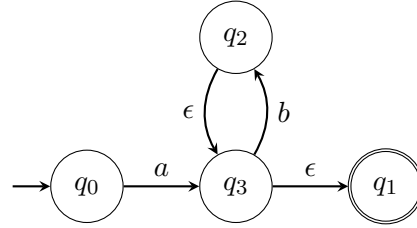


Fig. 4.5: Example grammar with left recursion Fig. 4.6: Resulting automaton for the grammar in Figure 4.5

$A \rightarrow a$
 $A \rightarrow B$
 $B \rightarrow bA$

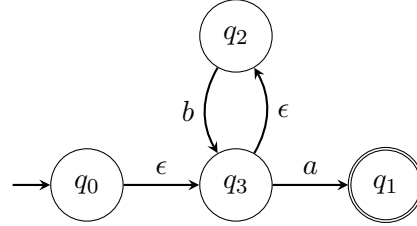


Fig. 4.7: Example grammar with right recursion Fig. 4.8: Resulting automaton for the grammar in Figure 4.7

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 identical paths, since $A \rightarrow \text{replace}[f, x](F)$ is treated like a second $A \rightarrow F$ production, 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 .

For the $\text{replace}[\text{old}, \text{new}]$ operation this transformation consists of replacing all occurrences of *old* on tainted edges with *new*, which gives us the automaton in Figure 4.10 for the given example.

4.4.2 Automaton to Regular Expression

State elimination

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

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 end state [6].

To achieve this characteristic, one can add a new start state with a single ϵ transition

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

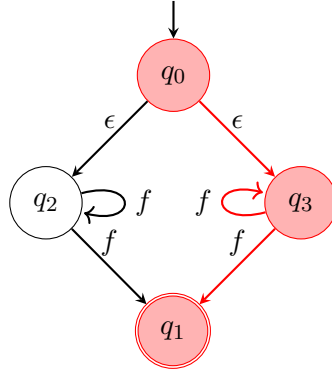


Fig. 4.8: Example grammar with operation production

Fig. 4.9: Resulting automaton for the grammar in Figure 4.8

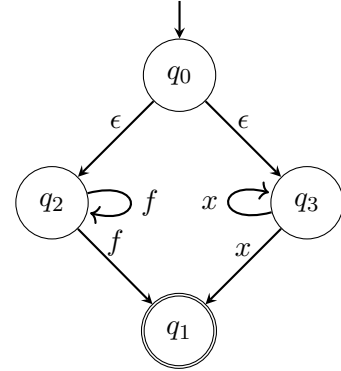


Fig. 4.10: Automaton in figure 4.9 after applying operation transformation

to the old start state and a new finale 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 use regular expressions as edge labels from the start.

We first 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 [4].

Figures 4.11 and 4.12 contain examples adapted from Esparza [4] 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.

Delgado heuristic

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 minimize the expression size.

We chose a heuristic described by Delgado and Morais [3]. For each state a weight is calculated using the following formula, where In_q is the set of incoming edges of 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

4 Approach and Implementation

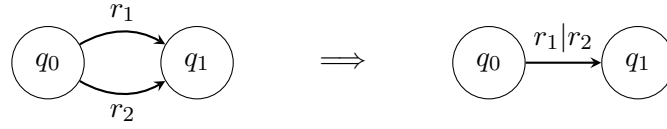


Fig. 4.11: Replacement of edge pairs

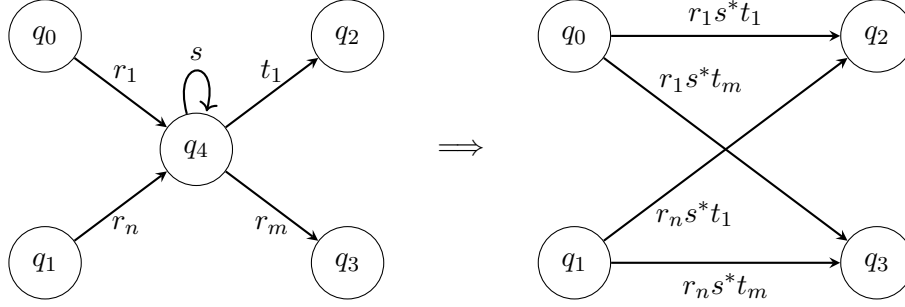


Fig. 4.12: Elimination of state q_4

not contain a potential loop on q . W_{loop} is the size of the loop around q if it exists and 0 otherwise.

$$weight(q) = \sum_{e \in In_q} (W_e \times (|Out_q| - 1)) + \sum_{e \in Out_q} (W_e \times (|In_q| - 1)) + W_{loop} \times (|In_q| \times |Out_q| - 1)$$

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.

Delgado and Morais show that using this heuristic produces significantly shorter expressions compared to the naive state elimination. Gruber et al. also show it outperforms almost all other heuristics they considered [5]. Improving the algorithm further by implementing a look-ahead additionally to the heuristic also improves the results, but adds more complexity and impairs the algorithms performance [3].

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 Evaluation and Discussion

- How did you test/evaluate your PoC?
 - E.g. case studies, large-scale studies, test bench, etc.
 - What did you do to verify results (if applicable)
- What did you learn from these tests? Depends on your work. E.g.
 - TP/TN/FP/FN rates
 - Performance
 - Results of your studies
 - Interpretation of the results, lessons learned
- Limitations of the approach and your implementation. Any ideas on how to fix them?

Probably 5-15 pages

6 Related Work

The challenge of statically obtaining information about the values of strings is not new and over the years there have been different approaches to it.

We follow the approach by Christensen et al. [2]. 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. They use the same approximation methods for obtaining regular languages from the generated context free grammars, but instead of making the regular languages available as a regular expression they generate automata. Furthermore they introduce a novel formalism, the multi-level automaton (MLFA) which allows easy extraction of these automata for different hotspots. Due to the aforementioned on-demand generation of the grammar, we don't need this extraction for single hotspots the MLFA provides in our implementation. The authors provide a feature rich implementation¹ of their approach and show that it efficiently produces useful results.

Tabuchi et al. [9] describe a type system for a minimal functional calculus, where strings have a regular expression as their type. They show that their proposed type system can produce good results when applied to their minimal calculus. While we considered implementing this approach for the analysis, there are some problems, especially due to our different requirements and 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 those, 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 even 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 adaption easier.

Wassermann and Su [11] present an approach comparable to ours, where they also

¹<https://www.brics.dk/JSA/>

6 *Related Work*

characterize values of string variables using context free grammars. They specifically target SQL injection vulnerabilities by using the generated CFGs to check whether user input 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.

7 Conclusion

Summarize your main contributions and observations. Further research directions?
 ≤ 1 page

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] 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.
- [3] 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.
- [4] J. Esparza. “Automata theory – An algorithmic approach.” Lecture Notes, <https://www7.in.tum.de/~esparza/autoskript.pdf>. Aug. 2017.
- [5] 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.
- [6] 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.
- [7] 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.
- [8] M.-J. Nederhof. “Regular approximation of CFLs: a grammatical view.” In: *Advances in Probabilistic and other Parsing Technologies*. Springer, 2000, pp. 221–241.
- [9] 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).
- [10] 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>.

Bibliography

- [11] 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.
- [12] K. Weiss and C. Banse. *A Language-Independent Analysis Platform for Source Code*. 2022. DOI: 10.48550/ARXIV.2203.08424.