# SIGMA - Normalization

Isaac Griffith and Rosetta Roberts
Empirical Software Engineering Laboratory
Informatics and Computer Science
Idaho State University
Pocatello, Idaho, 83208
Email: {grifisaa@isu.edu, roberose@isu.edu}

*Abstract*—**Introduction:**
**Objective:**
**Methods:**
**Results: What are the main findings? Practical implications?**
**Limitations: What are the weaknesses of this research?**
**Conclusions: What is the conclusion?**
*Index Terms*—**Island Grammars, Automated Grammar Formation, Software Language Engineering**

## I. INTRODUCTION

Multilingual parsing is an open problem that is currently being worked on. One method of multilingual parsing that has been developed is by creating island grammars for the combined grammars [1]. An automated method for doing this is currently being developed [SIGMA REFERENCE HERE]. One challenge to automating the creation of island grammar based multilingual parsers is detecting similar parts of grammars and combining them. This process becomes easier with an appropriate normalization process and normal form.

In this project, we propose to design and verify a procedure for appropriate normalization.

**RG** Design a normal form and normalization process for grammars that is conducive to identifying and merging similar parts across grammars to be used for merging grammars for the automated creation of multilingual parserss.

*Organization*

The remainder of this paper is organized as follows. Section II discusses the theoretical foundations of this work while also discussing other related studies. Section V details the design of experiments which evaluate the normalization steps of SIGMA. Section VI details the threats to the validity of this study. Finally, this paper is concluded in Section VII.

## II. BACKGROUND AND RELATED WORK

### A. Theoretical Foundations

<<<<< Updated upstream Context-free grammars are defined as $G = (V, \Sigma, P, S)$, where $V$ is the set of non-terminal symbols, $\Sigma$ is the set of terminal symbols, $P$ is the set of productions, and $S$ is the start production [2]. In this paper, we use a modified syntax for productions $\Phi \rightarrow R$. $\Phi$ is any non-terminal symbol and $R$ is a rule. Rules can be either a symbol, $\epsilon$ (the empty string), rules concatenated together, or rules unioned together with the | operator. Rather than using multiple productions when a symbol can produce multiple rules, the rules are combined with the | operator. For example productions normally written as

$$\langle A \rangle \rightarrow \texttt{a}$$
$$\langle A \rangle \rightarrow \texttt{b}$$

are written as the following instead

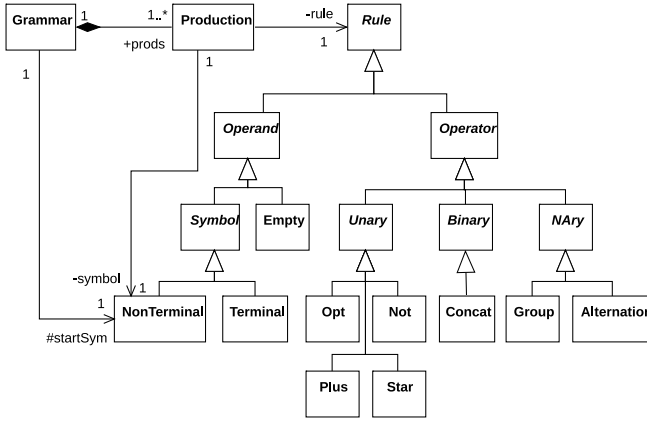$$\langle A \rangle \rightarrow \texttt{a | b}.$$

======= Context-free grammars are defined as $G = (V, \Sigma, P, S)$, where $V$ is the set of non-terminal symbols, $\Sigma$ is the set of terminal symbols, $P$ is the set of productions, and $S$ is the start production [2]. In this paper, a subset of Extended Backus-Naur Form (EBNF) to represent productions [?]. In each production $\Phi \rightarrow R$, $\Phi$ is a non-terminal symbol and $R$ is an expression representing what the non-terminal symbol can be expanded into. An expression can be either a non-terminal symbol, a terminal symbol or string, the empty string ($\epsilon$), expressions concatenated together, expressions separated with the | operator, or an expression surrounded by parenthesis. The | operator represents the union of sentences that can be expressed by each expression. The concatenation of expressions represents the sentences that can be formed by concatenating the sentences that each expression can form. »»»> Stashed changes

Various normal forms of grammars have been proposed by people. These normal forms have mainly been introduced to make parsing and manipulating grammars easier. The most used normal form is the Chomsky Normal Form (CNF) [3]. Grammars are transformed into CNF and other normal forms through simple transformations which preserve the language of the grammar. Normal forms for finding and merging similar parts across grammars have not been designed.

## III. APPROACH

### A. Domain Object Model

- Use sets for children of union operators.
- Union and concatenation operators are n-ary operators.
- Epsilon/empty string

```
 1: procedure NORMALIZE(𝒢)
 2:     repeat
 3:         𝒢 ← ELIMINATEUNUSEDPRODUCTIONS(𝒢)
 4:         𝒢 ← SIMPLIFYPRODUCTIONS(𝒢)
 5:         𝒢 ← MERGEEQUIVPRODUCTIONS(𝒢)
 6:         𝒢 ← ELIMINATEUNITPRODUCTIONS(𝒢)
 7:         𝒢 ← EXPANDPRODUCTIONS(𝒢)
 8:         𝒢 ← COLLAPSEPRODUCTIONS(𝒢)
 9:     until UNCHANGED(𝒢)
10: end procedure
```

## B. Design

To design our normal form, we decided that it should have the following properties.

1. Our domain object model represents each rule as a tree of operators and operands. Searching for similar rules and productions is simplified if each rule is composed of a single operator and its operands because this allows set and list based comparisons, which are simpler, to be performed rather than tree based comparisons. Because of this, we require that the normal form has at most single operator for each rule of the normalized grammar.
2. The size increase induced by normalization is minimal. Larger grammars are more difficult to process.
3. The normalized grammar is unambiguous given a grammar. I.e. there is exactly one normalized grammar for each grammar.
4. One anticipated difficult for searching similar portions of grammars is that refactoring that individual developer introduce that don't change the meaning of the grammar will result in portions of grammars being less similar. To mitigate this, we desire that certain transformations on the input grammar before normalization do not change the normalization result. The transformations we chose are the following

   1. Refactoring common terms of rules into a separate rule.
   2. Duplicating a rule.
   3. Introducing an unused non-terminal symbol and its production.
   4. Replacing a non-terminal symbol with a non-terminal symbol that produces that non-terminal symbol.
   5. Replacing all usages of a non-terminal symbol with its rule.

To meet these requirements, we decided that the normal form would have each production as one of the following forms.

1. $\langle Form_1 \rangle \rightarrow \langle A \rangle\ a\ \ldots$, where each term is a terminal symbol or a non-terminal symbol with an $F_2$ production and there are at least two terms in the rule.
2. $\langle Form_2 \rangle \rightarrow \langle A \rangle\ \ |\ \ a\ \ |\ \ \ldots$, where each term is a terminal symbol, the empty string, or a non-terminal symbol with an $F_1$ production and there are at least two

terms in the rule except for the special case when there is only production.

The reason we chose these two forms is that because rules of these forms are relatively easy to compare. This allows rules to easily be compared and merged. The restriction that non-terminal symbols referenced in each rule must have productions of the opposite form is so that examples

To meet requirement 4, our normalization process is performed using only transformation that are the inverse of transformations we do not want to affect our normalization process. To ensure that the size increase is minimal, we do not attempt to reverse transformations that rely on the distributive property between the concatenation and union | operators. To reverse transformations of this type, it is required to distribute productions. For example,

$$
\begin{aligned}
\langle A \rangle &\rightarrow a\ \langle B \rangle \\
\langle B \rangle &\rightarrow b\ |\ c
\end{aligned}
$$

would have to be transformed to

$$
\langle A \rangle \rightarrow a\,b\ |\ a\,c.
$$

Performing transformations of this kind repeatedly would result in an unreasonable increase in the number of productions. In addition, it would be unable to handle the case when a production indirectly references itself.

## IV. NORMALIZATION ALGORITHM

The following algorithm defines the approach for normalizing a given grammar. The normalization process defined here facilitates the ability to merge productions, in pursuit of the overarching goal of automated generation of Island [X], Tolerant [X], Bridge [X], and Bounded Seas [X] grammars.

This algorithm assumes that the source grammar, $G$, was initially in some defined formalism such as ANTLR [X], EBNF [X], BNF [X], SDF [X], TXL [X], etc. The grammar was then read in and processed to conform to the metamodel depicted in Figure **??**. Assuming that the grammar meets this condition, the goal of this algorithm is then to reformat the grammar such that each production is of one of Form$_1$ or Form$_2$

**Algorithm 2** Eliminate Unused Productions

1: **function** ELIMINATEUNUSEDPRODUCTIONS($\mathcal{G}$)
2:     $H \leftarrow (V, E)$
    ▷ Create empty graph
3:     **for all** $v \in \mathcal{G}.V$ **do**
4:         $\mathcal{H}.V \leftarrow \mathcal{H}.V \cup \{v\}$
5:         ADDRULETOGRAPH($\mathcal{G}, \mathcal{H}, \mathcal{G}.P(v)$)
6:         $\mathcal{H}.E \leftarrow \mathcal{H}.E \cup \{(v, \mathcal{G}.P(v))\}$
7:     **end for**
8:     DFSMARK($\mathcal{G}.S$)
9:     $\mathcal{G}.V \leftarrow \{ v \in \mathcal{G}.V \mid \text{MARKED}(v) \}$
10:     $\mathcal{G}.P \leftarrow \{ (v, \mathcal{G}.P(v)) \mid v \in \mathcal{G}.V \}$
11: **end function**
12: **function** ADDRULETOGRAPH($\mathcal{G}, \mathcal{H}, r$)
13:     $\mathcal{H}.V \leftarrow \mathcal{H}.V \cup \{r\}$
14:     **if** ISOPERATOR($r$) **then**
15:         **for all** $c \in \text{OPERANDS}(r)$ **do**
16:             ADDRULETOGRAPH($\mathcal{G}, \mathcal{H}, c$)
17:             $\mathcal{H}.E \leftarrow \mathcal{H}.E \cup \{(r, c)\}$
18:         **end for**
19:     **end if**
20: **end function**

---

**Algorithm 3** Depth First Marking

1: **function** DFSMARK($start$)
2:     $\mathcal{S} \leftarrow [start]$
3:     **while** $\mathcal{S} \neq \varnothing$ **do**
4:         $p \leftarrow \text{POP}(\mathcal{S})$
5:         MARK($p$)
6:         **for all** $s \in \text{SUCC}(p)$ **do**
7:             **if** !ISMARKED($s$) **then**
8:                 PUSH($\mathcal{S}, s$)
9:             **end if**
10:         **end for**
11:     **end while**
12: **end function**

---

**Algorithm 4** Simplify Productions

1: **function** SIMPLIFYPRODUCTIONS($\mathcal{G}$)
2:     **for all** $v \in \mathcal{G}.V$ **do**
3:         $\mathcal{G}.P(v) \leftarrow$ SIMPLIFYRULE($\mathcal{G}.P(v)$)
4:     **end for**
5: **end function**
6: **function** SIMPLIFYRULE($r$)
    ▷ Replace empty terminal string with $\epsilon$
7:     **if** ISTERMINAL($r$) $\wedge$ ISEMPTY($r$) **then**
8:         **return** $\epsilon$
9:     **end if**
10:     **if** ISOPERATOR($r$) **then**
11:         **let** $C$ be CHILDREN($r$)
12:         $C \leftarrow \{ \text{SIMPLIFYRULE}(c) \mid c \in C \}$
13:         **if** ISCONCATENATE($r$) **then**
14:             $C \leftarrow \{ c \in C \mid c \neq \epsilon \}$
15:             **if** $|C| = 0$ **then**
16:                 **return** $\epsilon$
17:             **end if**
18:         **end if**
19:         **if** ISCONCATENATE($r$) $\vee$ ISUNION($r$) **then**
        ▷ Replace operators with single operand with operand
20:             **if** $|C| = 1$ **then**
21:                 **let** $\{c\}$ be $C$
22:                 **return** $c$
23:             **else**
24:                 **return** $r$
25:             **end if**
26:         **end if**
27:     **end if**
28: **end function**

replacing operators with only one operand with their operator. This process is embodied in Algorithm 4.

### C. Merging Equivalent Productions

Productions that have identical rules are replaced by a single production. This new production is given a name derived from the productions that were merged to create it. The algorithm for this is shown in Alg. 5.

### D. Eliminating Unit Productions

All non-terminals with productions of one of the following two forms will have their non-terminal symbols replaced by their rules, and their productions eliminated.

$$\langle \text{a} \rangle \quad \rightarrow \quad \langle \text{b} \rangle$$
$$\langle \text{a} \rangle \quad \rightarrow \quad \text{a}$$

Elimination of productions of the first form, is derived from Chomsky Normal Form (CNF) [X]. Eliminations of productions of the second form, a derivation from CNF, allows the simplification process to simplify rules of the following form:

The normalization process, as defined in Algorithm 1, repeatedly executes six processes until the grammar stabilizes. These six processes are: i) eliminating unused rules, ii) simplifying productions, iii) merging equivalent rules, iv) eliminating unit rules, v) expanding productions, and vi) collapsing compatible productions.

### A. Eliminating Unused Rules

This process removes all productions that are not produced, directly or indirectly, from the start production. This is accomplished by enumerating all symbols producuable from the start symbol via a depth first search (see Algorithm 3) and then creating a new grammar using only the enumerated symbols, as shown in Algorithm 2.

### B. Simplifying Productions

This process aims to simplify productions. This is achieved by removing unnecessary $\varepsilon$'s concatenated with other rules and

**Algorithm 5** Merge Equivalent Productions

1: **function** MERGEEQUIVPRODUCTIONS($\mathcal{G}$)
2:    $pairs \leftarrow \varnothing$
3:    **for** $i \in [0, |\mathcal{G}.\Sigma|)$ **do**
4:       **for** $j \in (i, |\mathcal{G}.\Sigma|)$ **do**
5:          **if** $i \neq j$ **then**
6:             $pairs \leftarrow pairs \cup (\mathcal{G}.\Sigma[i], \mathcal{G}.\Sigma[j])$
7:          **end if**
8:       **end for**
9:    **end for**
10:    **for all** $p \in pairs$ **do**
11:       **if** $p.left.rule = p.right.rule$ **then**
12:          COMBINEANDREPLACE($p.left, p.right$)
13:       **end if**
14:    **end for**
15: **end function**

---

**Algorithm 6** Eliminate Unit Productions

1: **function** ELIMINATEUNITPRODUCTIONS($\mathcal{G}$)
2:    **for all** $p \in \mathcal{G}.\Sigma$ **do**
3:       **if** $|p.rule| = 1$ **then**
4:          REPLACE($uses(p), p.rule$)
5:       **end if**
6:    **end for**
7: **end function**

$$\begin{aligned}\langle a \rangle &\rightarrow \langle b \rangle \; \texttt{a} \; \texttt{b} \\ \langle b \rangle &\rightarrow \epsilon\end{aligned}$$

### E. Expanding Productions

Productions that have nested rules have all nested content replaced by with a non-terminal. The new non-terminal defines a production pointing to their content.

### F. Collapsing Compatible Productions

The final step of the normalization process combines productions that are associative with each other. This ensures that any non-terminal symbols referenced by a rule will not define a duplicate production. The following provides an example:

$$\begin{aligned}\langle A \rangle &\rightarrow \texttt{a} \; \langle B \rangle \\ \langle B \rangle &\rightarrow \texttt{b} \; \texttt{c} \\ \langle C \rangle &\rightarrow \texttt{c} \; | \; \langle D \rangle \\ \langle D \rangle &\rightarrow \texttt{d} \; | \; \texttt{e}\end{aligned}$$

would then collapse to form:

$$\begin{aligned}\langle A \rangle &\rightarrow \texttt{a} \; \texttt{b} \; \texttt{c} \\ \langle C \rangle &\rightarrow \texttt{c} \; | \; \texttt{d} \; | \; \texttt{e}\end{aligned}$$

**Algorithm 7** Expand Productions

1: **function** EXPANDPRODUCTIONS($\mathcal{G}$)
2:    **repeat**
3:       $changed \leftarrow \bot$
4:       **for all** $p \in \mathcal{G}.\Sigma$ **do**
5:          **if** ISCONCAT($p.rule$) **then**
6:             **for all** $g \in p.rule$ **do**
7:                **if** ISGROUP($g$) **then**
8:                   CREATEANDREPLACEWITH-PROD($g$)
9:                   $changed \leftarrow \top$
10:                **end if**
11:             **end for**
12:          **else if** ISALT($p.rule$) **then**
13:             **for all** $a \in p.rule$ **do**
14:                CREATEANDREPLACEWITHPROD($a$)
15:                $changed \leftarrow \top$
16:             **end for**
17:          **end if**
18:       **end for**
19:    **until** $changed = \bot$
20: **end function**

---

**Algorithm 8** Collapse Productions

1: **function** COLLAPSEPRODUCTIONS($\mathcal{G}$)
   ▷ Split productions into form1 and form2
2:    $f_1 \leftarrow$ COLLECT("form1")
3:    $f_2 \leftarrow$ COLLECT("form2")
4:    **for all** $p \in f_1$ **do**
5:       **if** ONLYTERMINALS($p.rule$) **then**
6:          REPLACEF1USESWITHRULE($p$)
7:       **end if**
8:    **end for**
9:    **for all** $p \in f_2$ **do**
10:       **if** ONLYTERMINALS($p.rule$) **then**
11:          REPLACEF2USEWITHRULE($p$)
12:       **end if**
13:    **end for**
14: **end function**

## V. EXPERIMENTAL DESIGN

### A. Pilot Study

««««< Updated upstream To evaluate the above approach, we performed a small pilot study on three grammars. We selected these grammars from the ANTLR grammar repository. To select these three grammars, we looked for three grammars of varying sizes and applications. The three grammars that we chose were the brainfuck, Java™, and XML grammars (Table **??**). We chose the brainfuck grammar because the small size of it allows it to be easily manually inspected. Java™'s grammar was chosen because it has a relatively complex grammar and the Java™programming language is used in corporate environments. XML's grammar was chosen because

| Language |
|---|
| Java™7 |
| Brainfuck |
| XML |

TABLE I
LANGUAGES USED IN PILOT STUDY.

it is relatively simple, but not as simple as brainfuck's, while also having significant uses.

======= To evaluate the above approach, we performed a small pilot study on three grammars. We selected these grammars from the ANTLR grammar repository. To select these three grammars, we looked for three grammars of varying sizes and applications. The three grammars chosen were the Brainfuck, XML, Java™ grammars. »»»> Stashed changes

Brainfuck is an esoteric language notable for its extreme simplicity [CITE]. It is Turing-complete despite only having 8 commands. We chose this language because its grammar is extremely small which allows it to be easily inspected. XML was also chosen because of its grammar's small size. However, it is still significantly more complex than Brainfuck. XML is commonly used for sending information between applications [CITE something about SOAP] and configuration files [CITE]. Java™ is a general purpose programming language used all over the world [CITE]. Its grammar is significantly more complicated than either of the two previously mentioned grammars. We chose to include this grammar because applications with a need for multilingual parsing would likely include Java as one of their languages [CITE something about JSP, something about embedded languages in strings, something about mixed language systems (e.g. migration to SCALA or Kotlin)].

To evaluate each grammar, we normalize each grammar. Before and after normalization, we measure the number of productions. After normalization, each grammar is checked manually. In this checking process, each rule is verified to be of either Form 1 or Form 2. In addition, we examine each normalized grammar for unexpected rules.

## VI. THREATS TO VALIDITY

## VII. CONCLUSIONS AND FUTURE WORK

The timeline to complete this study is as follows:

We intend to publish these results at one of the following conferences:

## ACKNOWLEDGEMENTS

## REFERENCES

[1] N. Synytskyy, J. R. Cordy, and T. R. Dean, "Robust multilingual parsing using island grammars," in *Proceedings of the 2003 Conference of the Centre for Advanced Studies on Collaborative Research*. IBM Press, 2003, pp. 266–278.
[2] M. Haoxiang, *Languages and Machines: An Introduction to the Theory of Computer Science*, 3rd ed. Boston, MA, USA: Addison-Wesley Longman Publishing Co. Inc., 1988.
[3] N. Chomsky, "On certain formal properties of grammars," *Information and Control*, vol. 2, no. 2, pp. 137–167, Jun. 1959.