

P4 Projekt

SPLAT
P4 PROJEKT
GROUP SW407F13
SOFTWARE
DEPARTMENT OF COMPUTER SCIENCE
AALBORG UNIVERSITY
MAY 2013



AALBORG UNIVERSITY
STUDENT REPORT

Titel:

SPLAT - Special Programming Language for Arduino Tipple-mixer

Project period:

P4, spring 2012

Project group:

SW407F13

Department of Computer Science

Selma Lagerlöfs Vej 300

DK-9220 Aalborg East

<http://www.cs.aau.dk/en>

Group members:

Aleksander Sørensen Nilsson

Christian Jødal O'Keeffe

Kasper Plejdrup

Mette Thomsen Pedersen

Niels Brøndum Pedersen

Rasmus Fischer Gadensgaard

Synopsis:

FiXme Fatal: synopsis mangler

Supervisor:

Ricardo Gomes Lage

Total number of pages:

35

Project end:

29th of May, 2013

The content of the report is freely available, but can only be published (with source reference) with an agreement with the authors.

Prolog

Aalborg March 22, 2013

FiXme Fatal: pr
mangler

Aleksander Sørensen Nilsson

Christian Jødal O’Keeffe

Kasper Plejdrup

Mette Thomsen Pedersen

Niels Brøndum Pedersen

Rasmus Fischer Gadensgaard

Contents

Prolog	v
1 Introduction	1
2 Problem statement	3
3 Analysis	5
3.1 Current language	5
3.2 Embedded systems	5
3.3 Arduino platform	5
4 Theory	7
4.1 Language	7
4.1.1 Paradigms of Programming Language	7
4.2 Compilers	9
4.3 Syntax analysis	9
4.3.1 Grammar	9
4.3.2 Lexical analyzer	11
4.4 Semantics	11
4.5 Contextual analysis	11
4.6 Code generation	11
5 Design	13
5.1 Syntax design	13
5.2 Choice of grammar	13
5.3 Semantics of SPLAT	17
5.3.1 Scoping	17
5.3.2 Type Checking	18
5.4 Code examples	18
6 Implementation	19
6.1 Known lexers and parsers	19
6.1.1 Lexer	19
6.1.2 Parser	20
6.1.3 Lexer and parser	20
6.1.4 Comparison table	21
6.2 Scanner class creation	21
6.3 Parser generation	21
6.4 Class generator classgen	21
6.5 Scope and type checking	21
6.6 Code generation	21

6.7 Test/evaluation	21
7 Conclusion	23
Bibliography	25
8 Appendix	27

Introduction 1

FiXme Fatal: inc
mangler

Problem statement 2

FiXme Fatal: pr
statement mang

Analysis 3

- 3.1 Current language
- 3.2 Embedded systems
- 3.3 Arduino platform

Theory 4

4.1 Language

4.1.1 Paradigms of Programming Language

In computer science, four main paradigms of programming languages exists [Nørmark, 2010]. In this section these paradigms will be shortly described followed by a subsection, explaining the choice of programming paradigm of the language in this project.

Imperative Programming

Imperative programming is a very sequential or procedural way to program, in the sense that a step is performed, then another step and so on. These steps are controlled by control-structures for example the if-statement. An example of a imperative programming language is C. Imperative programming language describes programs in terms of statements which alter the program state. This makes imperative languages very simple, and are also a good starting point for new programmers.

Functional Programming

Functional programming originates from the theory of functions in mathematics. In functional programming all computations are done by calling functions. In functional programming languages calls to a function will always yield the same result, if the function are called with the same parameters as input. This is in contrast to imperative programming where function calls can result in different values depending on the state of the program at that given time. Some examples of functional programming languages are Haskell and OCaml.

Logic Programming

Logic programming is fundamentally different from the imperative-, functional-, and object-oriented programming languages. In logic programming, it cannot be stated how a result should be computed, but rather the form and characteristics of the result. An example of a logic programming language is Prolog.

Object-Oriented Programming

Object-Oriented programming is based on the idea of data encapsulation, and grouping of logical program aspects. The concept of parsing messages between objects are also a very desirable feature when programs become of certain size. In object-oriented programming,

Readability	How easy it is to understand and comprehend a computation
Write-ability	How easy it is for the programmer to write a computation clearly, correctly, concisely and quickly
Reliability	Assures a program behaves the way it is suppose to
Orthogonality	A relatively small set of primitive constructs can be combined legally in a relatively small number of ways
Uniformity	If some features are similar they should look and behave similar
Maintainability	Errors can be found and corrected and new features can be added easily
Generality	Avoid special cases in the availability or use of constructs and by combining closely related constructs into a single more general one
Extensibility	Provide some general mechanism for the programmer to add new constructs to a language
Standardability	Allow programs to be transported from one computer to another without significant change in language structure
Implementability	Ensure a translator or interpreter can be written

Table 4.1: Brief explanation of language characteristics [Sebesta, 2009]

each class of object can be given methods, which is a kind of functions which can be called on that object. For example the expression `foo.Equals(bar)`, would call the `Equals`-method in the class of `'foo'`, and evaluate if `'bar'` equals `'foo'`. It is also easy in object-oriented languages to specify access-levels of classes, and thereby protect certain classes from external exposure. Classes can inherit from other classes. For example one could have a `'Car'`-class, which inherits all properties and methods of a `'Vehicle'`-class. This allows for a high degree of code-reuse.

Choice of Paradigm in This Project

For this project, an imperative approach has been chosen. The reason for this is that the programming language of this project should be very easy to understand for newcomers to programming. Also the programs in this programming language will likely remain of a relatively small length, which does not make object-orienting desired.

Design Criteria in this Project

To determine how a programming language should be syntactically described, the trade-offs of designing a programming language must be taken into care. The different characteristics of a programming language, which will be used to evaluate trade-offs can be seen on table ??.

These characteristics are used to evaluate the the trade-offs of programming language. An overview of these can be seen on table 4.2.

Based on these trade-offs, it is clear that having a simple programming language affects both readability, writability and reliability. This is because having a very simple-to-understand language, might not make it very writable. On the other hand, having a simple-to-write programming language, might not make it very readable. An example of this is the if-statement in C, which can be written both with the `'if'`-keyword, or more compact. This can be seen by comparing listing 4.1 with listing 4.2, which both yield the same result. It is then clear, that the compact if-statement might be faster to write, but slower to read and understand, and opposite with the if-statement.

Characteristic	Readability	Writability	Reliability
Simplicity	•	•	•
Orthogonality	•	•	•
Data types	•	•	•
Syntax design	•	•	•
Support for abstraction		•	•
Expressivity		•	•
Type checking			•
Exception handling			•
Restricted aliasing			•

Table 4.2: Overview of trade-offs [Sebesta, 2009]

```

1  if (x > y)
2  {
3      res = 1;
4  }
5  else
6  {
7      res = 0;
8  }

```

Listing 4.1: Simple example of if-statement in C using the 'if'-keyword

```

1  res = x > y ? 1 : 0;

```

Listing 4.2: Simple example of if-statement in C without using the 'if'-keyword

When defining the syntax of a programming language, it should balance these characteristics to achieve the right amount of trade-offs for that particular language. For the language of this project, it is important that the language is simple to read and understand, because the target group is the hobby-programmer, who might not have much experience in programming.

4.2 Compilers

4.3 Syntax analysis

4.3.1 Grammar

A grammar is used to define the syntax of a language. A context-free grammar (CFG) is a 4-tuple (V, Σ, R, S) finite language defined by [Sipser, 2013]:

1. V is a finite set called the variables
2. Σ is a finite set, disjoint from V called the terminals
3. R is a finite set of rules, with each rule being a variable and a string of variables and terminals
4. $S : S \in V$ is a start variable

The most common way of writing a CFG is by using Backus Naur Form (BNF) or Extended Backus Naur Form (EBNF). BNF is named after John Backus who presented the notation, and Peter Naur who modified Backus' method of notation slightly [Sebesta, 2009]. By using the BNF-notation it is possible to describe a CFG. It is preferred to have an unambiguously grammar. A CFG is ambiguous if a string derived in the grammar has two or more different leftmost derivations [Sipser, 2013]. An unambiguously grammar will ensure that a program running through a string using CFG can only read the string in one way.

A CFG is a part of the $LL(k)$ grammar classes if it is possible to produce the leftmost derivation of a string by looking at most k tokens ahead in the string. LL algorithms work on the same subset of free grammars which means that LL parsers work on $LL(k)$ grammars. $LL(k)$ means that the grammar needs to be left-recursive free which makes it possible to create a top-down leftmost derivation parser. The $LL(1)$ have properties that makes the grammar attractive for simple compiler construction. A property is that $LL(1)$ grammars are fairly easy compared to $LL(k)$ where $k > 1$ to implement because the parser analyser only has to look one element ahead in order to determine what parser action there should be taken. $LL(1)$ is also relatively faster than $LL(k)$ where $k > 1$ based on the same reason, that the parser only has to look one element ahead. A disadvantage of the LL grammars is that the parser finds syntax errors towards the end of parsing process where a LR parser is faster at detecting the syntax errors. LL is also inferior compared to LR in terms of describing a language based on the idea that LL is a subclass of the bigger grammar class LR . That means with a LR grammar it is possible to describe aspects of a language that might not be possible in a LL grammar [Fischer et al., 2009] [Sebesta, 2009].

A CFG is a part of the $LR(k)$ grammar classes if it is possible to produce the rightmost derivation in reverse of a string by looking at most k tokens ahead in the string. LR grammars are a superset for the LL grammars meaning that LR covers a larger variety of programming language than LL . LR parser is a bottom-up parser meaning that it starts constructing the abstract trees from its leaf and works its way to the root. LR parsers are generally harder to implement than LL parsers by hand but there exist tools that automatically can generate LR parsers. $LR(k)$ grammars allow left recursion which means that the LR grammars are a bigger grammar class than LL . $LALR$ and $SLAR$ are subclasses of the $LR(k)$ grammars which means that $LR(k)$ describes a larger language at the cost of a bigger parser table in comparison to $SLAR$ and $LALR$. The balance of power and efficiency makes the $LALR(1)$ a popular table building method compared to LR building method [Fischer et al., 2009] [Sebesta, 2009].

Based on these understandings of grammars there will be a section where there will be looked into which grammar that will be used in this project.

4.3.2 Lexical analyzer

A lexical analyzer reads the input file, and returns a series of tokens based on the input [Fischer et al., 2009]. More specifically it is the scanner in the lexical analyzer which does this. These tokens are matched by rules, usually described by regular expressions. An example of such grammar rules can be seen on table 4.3. Formally a token consists of two parts: The token type, and the token value [Fischer et al., 2009]. As an example the IDENT token seen on 4.4 has the token type IDENT and the value 'c'.

Terminal	Regular expression
dcl	"[a - z]"
assign	"="
digit	"[0 - 9]+"
endassign	","
blank	" "+

Table 4.3: Sample token specification

This specification of tokens, would be used by the scanner to determine how tokens looks, and thereby which text-elements are tokens.

```
1 c = 42;
```

Listing 4.3: Simple example of code

As an example the lines of code seen on listing 4.3 might be read as the tokens seen on table 4.4.

Token	Lexeme
IDENT	c
ASSIGN	=
DIGIT	42
SEMICOLON	;

Table 4.4: Example of tokens

The scanner produces a stream of tokens, which is returned to the parser. The parser checks if the tokens conforms to the language-specification [Fischer et al., 2009].

4.4 Semantics

4.5 Contextual analysis

4.6 Code generation

Design 5

5.1 Syntax design

5.2 Choice of grammar

The programmer, using this projects language, could be a hobby programmer, who would want to program a custom drink machine, but does not possess a high level of education in programming. Therefore it was decided that the grammar should have a high level of readability because this will ensure that it is easier for the person to read and understand their program - also useful if the code has to be edited later on. This on the other hand can decrease the level of write-ability because it has to be written in a specific way and will need to contain some extra words or symbols to mimic a language closer to human language rather than a computer language.

The method to assign a value to a variable is by typing "*variable* <- *valuetoassign*" this approach have been chosen, instead of the more commonly used "=" symbol, because a person not accustomed to programming might confuse which side of the "=" is assigned to the other. Thus by using the arrow, it is more clearly indicated that the value is assigned to the variable, and therefore ensuring readability - especially for the hobby programmer.

To get a more symmetrical structure in the code the functions must always return something, but it can return the value "nothing". This will ensure a better understanding and readability of the code when the programmer can see what it returns, even if no value is parsed. To indicate that *return* is the last thing that will be executed in a function, the *return* must always be at the end of the function. To indicate that a program is called "call *functionname*" must be written. Words are used instead of symbols, when suitable, to improve the understanding of the program(compared to most other programming languages). "begin" and "end" are used to indicate a block (eg. an "if" statement). To combine logical operators the words "AND" and "OR" are used. The ";" symbol is used to improve readability by making it easier to see when the end of a line has been reached.

It would be appropriate to design a grammar that is a subset of $LL(1)$ grammars. This is based on the idea that it easier to implement a parser for $LL(1)$ grammars by hand compared to LR grammars. This approach means it would be possible to both implement a parser by hand or use some of the already existing tools. This way both approaches are possible which are a suited solution for the project because it allows the project group to later go back and make the parser by hand instead of using a tool if so desired.

If the purpose was to create an efficient compiler it would be more appropriate to design the grammar as a subset of the *LALR* grammar class. A parser for *LALR* is balanced between power and efficiency which makes it more desirable than *LL* and other *LR* grammars, see section 4.3.1 for more on the grammars. *LR* parsers can be made by hand but it is much more difficult than the *LL* parsers.

$$\langle \text{program} \rangle \rightarrow \langle \text{roots} \rangle$$

$$\begin{aligned} \langle \text{roots} \rangle &\rightarrow \varepsilon \\ &| \quad \langle \text{root} \rangle \langle \text{roots} \rangle \end{aligned}$$

$$\begin{aligned} \langle \text{root} \rangle &\rightarrow \langle \text{dcl} \rangle; \\ &| \quad \langle \text{function} \rangle \\ &| \quad \langle \text{comment} \rangle \end{aligned}$$

$$\langle \text{dcl} \rangle \rightarrow \langle \text{type} \rangle \langle \text{id} \rangle \langle \text{dclend} \rangle$$

$$\langle \text{type} \rangle \rightarrow \langle \text{primitivetype} \rangle \langle \text{arraytype} \rangle$$

$$\begin{aligned} \langle \text{primitivetype} \rangle &\rightarrow \text{bool} \\ &| \quad \text{double} \\ &| \quad \text{int} \\ &| \quad \text{char} \\ &| \quad \text{container} \\ &| \quad \text{string} \end{aligned}$$

$$\begin{aligned} \langle \text{arraytype} \rangle &\rightarrow \langle \text{type} \rangle [] \\ &| \quad \varepsilon \end{aligned}$$

$$\langle \text{id} \rangle \rightarrow \langle \text{letter} \rangle \langle \text{idend} \rangle$$

$$\langle \text{letter} \rangle \rightarrow [\text{a} - \text{zA} - \text{Z}]$$

$$\begin{aligned} \langle \text{idend} \rangle &\rightarrow \langle \text{letter} \rangle \langle \text{idend} \rangle \\ &| \quad \langle \text{digit} \rangle \langle \text{idend} \rangle \\ &| \quad \varepsilon \end{aligned}$$

$$\begin{aligned} \langle \text{dclend} \rangle &\rightarrow \varepsilon \\ &| \quad \langle \text{assign} \rangle \end{aligned}$$

$$\langle \text{assign} \rangle \rightarrow <-- \langle \text{expr} \rangle$$

$$\langle \text{expr} \rangle \rightarrow \langle \text{term} \rangle \langle \text{exprlend} \rangle$$

$$\langle \text{term} \rangle \rightarrow \langle \text{comp} \rangle \langle \text{termend} \rangle$$

$$\langle \text{comp} \rangle \rightarrow \langle \text{factor} \rangle \langle \text{compend} \rangle$$

$$\begin{aligned} \langle \text{factor} \rangle &\rightarrow (\langle \text{expr} \rangle) \\ &| \quad !(\langle \text{expr} \rangle) \\ &| \quad \langle \text{callid} \rangle \\ &| \quad \langle \text{numeric} \rangle \\ &| \quad \langle \text{string} \rangle \end{aligned}$$

$\mid \langle functioncall \rangle$
 $\mid \langle cast \rangle$
 $\mid \text{LOW}$
 $\mid \text{HIGH}$
 $\mid \text{true}$
 $\mid \text{false}$

$\langle callid \rangle \rightarrow \langle id \rangle \langle arraycall \rangle$

$\langle arraycall \rangle \rightarrow [\langle notnulldigits \rangle]$
 $\mid \varepsilon$

$\langle notnulldigits \rangle \rightarrow \langle notnulldigit \rangle \langle digits \rangle$

$\langle notnulldigit \rangle \rightarrow [1 - 9]$

$\langle digits \rangle \rightarrow \varepsilon$
 $\mid \langle digit \rangle \langle digits \rangle$

$\langle digit \rangle \rightarrow [0 - 9]$

$\langle numeric \rangle \rightarrow \langle plusminus \rangle \langle digitsnotempty \rangle \langle numericend \rangle$

$\langle plusminus \rangle \rightarrow \varepsilon$
 $\mid -$

$\langle digitsnotempty \rangle \rightarrow \langle digit \rangle \langle digits \rangle$

$\langle numericend \rangle \rightarrow \varepsilon$
 $\mid . \langle digitsnotempty \rangle$

$\langle string \rangle \rightarrow " \langle stringmidt \rangle "$

$\langle stringmidt \rangle \rightarrow \langle letter \rangle \langle stringmidt \rangle$
 $\mid \langle symbol \rangle \langle stringmidt \rangle$
 $\mid \langle digit \rangle \langle stringmidt \rangle$
 $\mid \varepsilon$

$\langle symbol \rangle \rightarrow !$

$\mid \%$
 $\mid ^$
 $\mid \&$
 $\mid ($
 $\mid)$
 $\mid _$
 $\mid +$
 $\mid |$
 $\mid \sim$
 $\mid -$
 $\mid =$
 $\mid ,$
 $\mid \{$
 $\mid \}$

| [
 |]
 | :
 | ;
 | ?
 | ,
 | .
 | /
 | ' ,

$\langle functioncall \rangle \rightarrow \text{call } \langle id \rangle (\langle callexpr \rangle)$

$\langle callexpr \rangle \rightarrow \langle subcallexpr \rangle$
 | ε

$\langle subcallexpr \rangle \rightarrow \langle expr \rangle \langle subcallexprend \rangle$

$\langle subcallexprend \rangle \rightarrow , \langle subcallexpr \rangle$
 | ε

$\langle cast \rangle \rightarrow \langle type \rangle (\langle expr \rangle)$

$\langle compend \rangle \rightarrow \langle comparisonoperator \rangle \langle comp \rangle$
 | ε

$\langle comparisonoperator \rangle \rightarrow >$
 | $<$
 | $<=$
 | $>=$
 | $!=$
 | $=$

$\langle termend \rangle \rightarrow * \langle term \rangle$
 | $/ \langle term \rangle$
 | AND $\langle term \rangle$
 | ε

$\langle exprend \rangle \rightarrow + \langle expr \rangle$
 | $- \langle expr \rangle$
 | OR $\langle expr \rangle$
 | ε

$\langle function \rangle \rightarrow \langle functionstart \rangle \langle functionmidt \rangle$

$\langle functionstart \rangle \rightarrow \text{function } \langle id \rangle \text{ return}$

$\langle functionmidt \rangle \rightarrow \langle type \rangle \langle functionend \rangle \langle expr \rangle; \text{end}$
 | nothing $\langle functionend \rangle$ nothing; end

$\langle functionend \rangle \rightarrow \text{using } (\langle params \rangle) \text{ begin } \langle stmts \rangle \text{ return}$

$\langle params \rangle \rightarrow \langle subparams \rangle$
 | ε

$$\begin{aligned}
\langle subparams \rangle &\rightarrow \langle type \rangle \langle id \rangle \langle subparamsend \rangle \\
\langle subparamsend \rangle &\rightarrow , \langle subparams \rangle \\
&| \quad \varepsilon \\
\langle stmts \rangle &\rightarrow \varepsilon \\
&| \quad \langle stmt \rangle \langle stmts \rangle \\
\langle stmt \rangle &\rightarrow \langle callid \rangle \langle assign \rangle; \\
&| \quad \langle nontermif \rangle \\
&| \quad \langle nontermwhile \rangle \\
&| \quad \langle from \rangle \\
&| \quad \langle dcl \rangle; \\
&| \quad \langle functioncall \rangle; \\
&| \quad \langle nontermswitch \rangle \\
&| \quad \langle comment \rangle \\
\langle nontermif \rangle &\rightarrow \text{if}(\langle expr \rangle) \text{begin} \langle stmts \rangle \text{end} \langle endif \rangle \\
\langle endif \rangle &\rightarrow \text{else} \langle nontermelse \rangle \\
&| \quad \varepsilon \\
\langle nontermelse \rangle &\rightarrow \langle nontermif \rangle \\
&| \quad \text{begin} \langle stmts \rangle \text{end} \\
\langle nontermwhile \rangle &\rightarrow \text{while}(\langle expr \rangle) \text{begin} \langle stmts \rangle \text{end} \\
\langle from \rangle &\rightarrow \text{from} \langle expr \rangle \text{to} \langle expr \rangle \text{step} \langle assign \rangle \text{begin} \langle stmts \rangle \text{end} \\
\langle nontermswitch \rangle &\rightarrow \text{switch} (\langle expr \rangle) \text{begin} \langle cases \rangle \text{end} \\
\langle cases \rangle &\rightarrow \text{case} \langle expr \rangle: \langle stmts \rangle \langle endcase \rangle \\
\langle endcase \rangle &\rightarrow \langle cases \rangle \\
&| \quad \text{break}; \langle breakend \rangle \\
&| \quad \text{default:} \langle stmts \rangle \text{break}; \\
\langle breakend \rangle &\rightarrow \langle cases \rangle \\
&| \quad \text{default:} \langle stmts \rangle \text{break}; \\
&| \quad \varepsilon \\
\langle comment \rangle &\rightarrow /* \langle stringmid \rangle */
\end{aligned}$$

5.3 Semantics of SPLAT

In this section the semantics of SPLAT will be described.

5.3.1 Scoping

The scope of a variable is the block of the program, in which it is accessible. A variable is local to a block, if it is declared in that block. A variable is non-local to a block if it is not declared in that block, but is still visible in that block (ex. global variables).

In SPLAT static scoping is used. This means that scopes are computed at compile time, based on the inputted program text. Static scoping means that a hierarchy of scopes are maintained during compilation. To determine the name of used variables, the compiler must first check if the variable is in the current scope. If it is, the value of the variable is found, and the compiler can proceed. Else it must recursively search the scope hierarchy for the variable. When done, if the variable is still not found, the compiler returns an error, because an undeclared variable is used.

Symbol tables

Generally there are two approaches to symbol tables: One symbol table for each scope, or one global symbol table.

Multiple Symbol Tables

In each scope, a symbol table exists, which is an ADT (Abstract Data Type), that stores identifier names and relate each identifier to its attributes. The general operations of a symbol table is: Empty the table, add entry, find entry, open and close scope.

It can be useful to think of this structure of static scoping and nested symbol tables as a kind of tree structure. Then when the compiler analyzes the tree, only one branch/path is available at a time. This exactly creates these features of e.g. local variables.

A stack might intuitively make sense because of the way scopes are defined by begin and end. A begin scope would simply push a symbol table scope to the stack, and when the scope ends, the symbol table is popped from the stack. This also accounts for nested scopes. But searching for a non-local variable would require searching the entire stack.

One Symbol Table

To maintain one symbol table for a whole program, each name will be in the same table. The names must therefore be named appropriately by the compiler, so that each name also contain information about nesting level. Various approaches to maintain one symbol table exists, for example maintaining a binary search tree might seem like a good idea, because it is generally searchable in $O(\lg(n))$. But the fact that programmers generally does not name variables and functions at random, causes the search to take as long as linear search. Therefore hash-tables are generally used. This is because of hash-tables perform excellent, with insertion and searching in $O(1)$, if a good hash function and a good collision-handling technique is used.

5.3.2 Type Checking

5.4 Code examples

Implementation 6

6.1 Known lexers and parsers

In this section some of the different lexers and parsers, that are available on the internet, will be described.

6.1.1 Lexer

These programs generate a lexical analyzer also known as a scanner, that turns code into tokens which a parser uses.

Lex

Files are divided into three sections separated by lines containing two percent signs. The first is the "definition section" this is where macros can be defined and where headerfiles are imported. The second is the "Rules section" where regular expressions are read in terms of C statements. The third is the "C code section" which contains C statements and functions that are copied verbatim to the generated source file. Lex is not open source, but there are versions of Lex that are open source such as Flex, Jflex and Jlex. [Lex]

Flex

Alternativ to lex [Flex]

An optional feature to flex is the REJECT macro, which enables non-linear performance that allows it to match extremely long tokens. The use of REJECT is discouraged by Flex manual and thus not enabled by default.

The scanner flex generates does not by default allow reentrancy, which means that the program can not safely be interrupted and then resumed later on.

Jflex

Jflex is based on Flex that focuses on speed and full Unicode support. It can be used as a standalone tool or together with the LALR parser generators Cup and BYacc/J [Jflex]

Jlex

Based on lex but used for java. [Jlex]

6.1.2 Parser

Parsertools generates a parser, based on a formal grammar from a lexer, checks for correct syntax and builds a data structure (Often in the form of a parse tree, abstract syntax tree or other hierarchical structure).

Yacc

Generates a LALR parser that checks the syntax based on an analytic grammar, written in a similar fashion to BNF. Requires an external lexical analyser, such as those generated by Lex or Flex. The output language is C. [Yacc]

Cup

More or less like Yacc, output language is in java instead. [Cup]

6.1.3 Lexer and parser

Combines the lexer and parser in one tool.

SableCC

Using the CFG(Context Free Grammar) written in Extended Backus-Naur Form SableCC generates a LALR(1) parser, the output languages are: C, C++, C#, Java, OCaml, Python [SableCC].

ANTLR

ANother Tool for Language Recognition uses the CFG(Context Free Grammar) written in Extended Backus-Naur Form to generate an LL(*) parser. It has a wide variety of output languages, including, C, C++ and Java. ANTLR can also make a tree parsers and combined lexer-parsers. It can automatically generate abstract syntax trees with a parser. [Antlr]

JavaCC

Javacc generate a parser from a formal grammar written in EBNF notation. The output is Java source code. JavaCC generates top-down parsers, which limits it to the LL(k) class of grammars (in particular, left recursion cannot be used). JavaCC also generates lexical analyzers in a fashion similar to lex[Norvell]. The tree builder that accompanies it, JJTree, constructs its trees from the bottom uplex[JJTree].

6.1.4 Comparison table

Name	Parsing algorithm	Input notation	Output language
Yacc	LALR(1)	YACC	C
Cup	LALR(1)	EBNF	java
SableCC	LALR(1)	EBNF	C, C++, C#, java, OCaml, Python
ANTLR	LL(*)	EBNF	ActionScript, Ada95, C, C++, C#, Java, JavaScript, Objective-C, Perl, Python, Ruby
JavaCC	LL(k)	EBNF	Java, C++(beta)

Based on the different lexers and parsers attributes, compared to the expectations of this project, it has been decided that ANTLR best fit the project. The reason behind this is that ANTLR uses the LL(*) parser algorithm, this fits the structure of the CFG grammar for this project. Furthermore ANTLR's output language can be in Java, C or C++, this makes it easier to work on an Arduino. Another possibility could be to write the lexer and parser by hand, but many typing errors are avoided by using a tool like ANTLR. Furthermore, it is easier to maintain the lexer and parser with a tool. When the grammar is changed, you can just generate a new lexer and parser with the tool. It has therefore been decided to use ANTLR for generating the lexer and parser in this project.

6.2 Scanner class creation

6.3 Parser generation

6.4 Class generator classgen

6.5 Scope and type checking

6.6 Code generation

6.7 Test/evaluation

Conclusion 7

FiXme Fatal: ko
mangler

Bibliography

- Antlr.** Antlr. *Theory behind Antlr*. URL <http://www.antlr.org/about.html>.
- Cup.** Cup. *Theory behind CUP*. URL <http://www2.cs.tum.edu/projects/cup/manual.html>.
- Fischer et al., 2009.** Charles N. Fischer, K. Cyton Ron og J. LeBlanc. Jr. Richard. *Crafting a Compiler*. Pearson, 2009.
- Flex.** Flex. *Theory behind Flex*. URL <http://flex.sourceforge.net/manual/>.
- Jflex.** Jflex. *Theory behind Jflex*. URL <http://jflex.de/manual.html#SECTION00040000000000000000>.
- JJTree.** JJTree. *JJTree to JavaCC*. URL <http://tomcopeland.blogspot.com/juniordeveloper/2007/10/better-jjtree-v.html>.
- Jlex.** Jlex. *Theory behind Jlex*. URL <http://www.cs.princeton.edu/~appel/modern/java/JLex/current/manual.html>.
- Lex.** Lex. *Theory behind Lex*. URL <http://dinosaur.compilertools.net/lex/index.html>.
- Norvell.** Theodore S. Norvell. *The JavaCC FAQ*. URL <http://www.engr.mun.ca/~theo/JavaCC-FAQ/javacc-faq-moz.htm#jjtree-and-jtb>.
- Nørmark, July 7 2010.** Kurt Nørmark. *Overview of the four main programming paradigms*. web, 2010. URL http://people.cs.aau.dk/~normark/prog3-03/html/notes/paradigms_themes-paradigm-overview-section.html.
- SableCC.** SableCC. *SableCC homesite*. URL <http://sablecc.org/wiki>.
- Sebesta, 2009.** Robert W. Sebesta. *Concepts of Programming Languages*. Pearson, 9 udgave, 2009.
- Sipser, 2013.** Michael Sipser. *Introduction to the Theory of Computation*. PWS Publishing, 3 udgave, 2013.
- Yacc.** Yacc. *Theory behind Yacc*. URL <http://dinosaur.compilertools.net/yacc/index.html>.

List of Corrections

Fatal: synopsis mangler	iii
Fatal: prolog mangler	v
Fatal: indledning mangler	1
Fatal: problem statement mangler	3
Fatal: Er det compileren???	18
Fatal: konklusion mangler	23

Appendix 8
