

P4 Project

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



Titel:

SPLAD - Special Programming Language for Arduino Drinks-mixer

AALBORG UNIVERSITY
STUDENT REPORT

Project period:

P4, spring 2013

Department of Computer Science

Selma Lagerlöfs Vej 300

DK-9220 Aalborg East

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

Project group:

SW407F13

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:

36

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 April 12, 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
1.1 Environment for this project	1
1.1.1 Solution in bars	1
1.2 Problem statement	1
1.2.1 Sub Statements	2
1.3 Report Structure	2
2 Language Specification	3
2.1 Paradigms of Programming Language	3
2.1.1 Choice of Paradigm in This Project	4
2.2 Design Criteria	4
2.2.1 Design Criteria in This Project	4
2.3 Syntax	6
2.3.1 Grammartypes	6
2.3.2 Grammar	6
2.3.3 Choice of Grammar	8
2.3.4 The BNF of SPLAD	9
2.3.5 EBNF?	12
2.3.6 Lexicon	12
2.4 Semantics	12
2.4.1 Scoping	12
2.4.2 Transition Rules	13
2.4.3 Type Rules	13
2.5 Code Examples	14
3 Implementation	15
3.1 Design Criteria	15
3.2 Architecture	15
3.2.1 Hardware	15
3.2.2 Overview of the Compiler	17
3.2.3 Language Processing Strategy	18
3.2.4 Compilation Passes	18
3.2.5 Abstract Syntax Trees	18
3.2.6 The Visitor Pattern	19
3.3 Syntactic Analysis	19
3.3.1 Semantic Analysis	19
3.3.2 Known lexers and parsers	19
3.3.3 ANTLR	20
3.3.4 Lexical Analyzer	22

3.3.5	Tokens	23
3.3.6	Parser	24
3.4	Contextual Analysis	25
3.4.1	Scope Checking	25
3.4.2	Type Checking	25
3.5	Code Generation	25
3.5.1	Declaration	25
3.5.2	Commands	25
3.5.3	Expressions	25
4	Conclusion	27

Introduction

1

1.1 Environment for this project

The environment in this project is an Arduino-based drinks-machine. The basic idea behind this environment, is that the drinks-machine resides in a bar or to a party, where customers will buy RFID-tags, which contains information about a specific drink. These RFID-tags can be read by the drinks-machine by a RFID-reader. The machine will then automatically mix the drink on the RFID-tag.

1.1.1 Solution in bars

Currently bars and clubs often have multiple bartenders who mixes the drinks and serves the customers. The bartender handles both receiving the order, mixing the drink and handles the payment of the drink. This process can be done more efficiently. If a bar had a drinks-machine like the one described above, the bar would require only one cashier instead of multiple bartenders. The cashier would handle selling and programming the RFID-tags. The customers themselves then places the RFID-tag on the RFID-reader on the drinks-machine, and the machine mixes the appropriate drink, and either decreases the count on the RFID-tag, or disables the tag, and display an appropriate message to the customer on the display of the machine. The RFID-tags are programmed by the cashier which encodes in the tag a drink id, and a drink count. This allows the customer to buy for example 10 cosmopolitans on the same tag. The machine will simply check if the drink-count on the tag is above zero, and display an error if it is not.

1.2 Problem statement

In this section a problem statement will be presented, which will be used as a basis for this project. In this project it has been chosen to examine how drink machine could be programmed using Arduino as platform for the processing. As mentioned in section 3.2.1, the programming language usually used for Arduino is based on C and C++, which is not aimed at programming drink machines as programming purpose. It could be useful to have a niche programming language aimed directly at programming drink machines on a Arduino platform. This will be the goal of this project.

The programming language in this project is aimed at the hobby programmer who wants to program his own drink machine. Because of this, the programs written in this language must be easy to understand and maintain. This however sacrifices some write-ability of the programs, because of constraints imposed to make sure programs are easily

Fixme Fatal: in
mangler

understandable. These trade-offs and will be further discussed in section 2.3.3. A hobby programmer is defined as a programmer who knows the basic structure of programming, but does not have an education in programming or work with software development.

Based on the above, the following problem statement comes to light:

- **How can a programming language be developed, which makes it easy for the hobby programmer to program drink machines based on Arduino platforms?**

The purpose of this problem statement is to guide the programming language for this project, so when the programming language reaches a final state, it is easy for hobby programmers to program using the language.

1.2.1 Sub Statements

On the basis of the problem statement, a number of sub-statements arises:

- **How can a programming language be specified, which makes it easy for novice programmers to learn it?** Because the language of this project is aimed at hobby programmers, the programming language should be specified in a way which is suited for the programmer.
- **How can a compiler be developed, which recognizes the language, and translates the source program into Arduino suitable code?** Of course it is not enough to have an easy-to-understand language, if it does not have a compiler for that language. The language would then render useless. This is the reason why a compiler must be developed, either by compiling the program code directly to Arduino machine code, or by first compiling the program code to c code, and then use the Arduino compiler to compile that code further.

1.3 Report Structure

Language Specification 2

2.1 Paradigms of Programming Language

In computer science, four main paradigms of programming languages exists [?]. In this section these paradigms will be briefly 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 state of the program. 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 is 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 the given time. Some examples of functional programming languages are Haskell and OCaml.

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 reach a certain size. In object-oriented programming, each class of objects can be given methods, which is a kind of function 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 relatively 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.

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.

2.1.1 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 desirable.

2.2 Design Criteria

To make it easier to program the machine described above, it would be nice to have a programming language aimed specifically at programming drinks-machines. This would make it easier for programmers with little or no experience to program the machine, and thereby making it easier for e.g. bar-owners to program and install their own drinks-machines in their bar. Therefore it has been decided to make a programming language aimed at this problem. The SPLAD language: **S**imple **P**rogramming **L**anguage for **A**rduino **D**rink-mixer. The SPLAD will be described more thoroughly in section 1.2

2.2.1 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 account. The different characteristics of a programming language, that will be used to evaluate trade-offs can be seen on table 2.1.

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 2.1: Brief explanation of language characteristics [?]

These characteristics are used to evaluate the the trade-offs of programming language. An overview of these can be seen on table 2.2.

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

Table 2.2: Overview of trade-offs [?]

In table 2.2 • means that the characteristic affects the feature of the programming language where the • is placed. If there is no • in front of a feature, it means that this particular characteristic is not affected by the feature.

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 2.1 with listing 2.2, 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.

```

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

Listing 2.1: Simple example of if-statement in C using the 'if'-keyword

```

1      res = x > y ? 1 : 0;
```

Listing 2.2: 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 hobbyist-programmer, who might not have much experience in programming.

2.3 Syntax

2.3.1 Grammartypes

The Chomsky hierarchy is a hierarchy of formal grammars. There are 4 types of grammar in the Chomsky hierarchy, where type 0 is the most unrestricted grammar, and type 3 is the most restricted grammar. These types of grammar are described below. The Chomsky hierarchy is used to divide the grammars into different types. All the grammars in a given type is a subset of the less restricted types. This means that if a language is a type 2, it is a subset of the grammars of type 1 and 0, but it is not equal to any of these grammars [?].

Type - 3: Regular Grammar

Regular grammars can be described by finite automata or regular expressions. Regular grammars are meant to be used on computers with an extremely limited amount of memory, because regular languages do not need to use a lot of memory to recognize a language[?].

Type - 2: Context-Free Grammar

Context-free grammars are described by substitution rules, also called productions. Substitution rules for context-free grammars can make the grammar ambiguous. This is a problem since different computers might yield different output for the same grammar. Context-Free Grammars can be described in Backus Naur form or by a Pushdown automata (PDA). PDA's works almost in the same way as finite automata. The difference is that a PDA uses a stack as memory to help create the output[?].

Type - 1: Context-Sensitive Grammar

Context-sensitive grammars substitution rules have nearly the same rules as those used in Context-free grammar. But in context-sensitive the right side of the production can have more then one terminal and there can be non-terminals on the right side of the production. A context-sensitive grammar can be recognized by a linear-bounded automaton.

Type - 0: Recursively Enumerable

Recursively enumerable or unrestricted grammar is a type of grammar, where there is no restrictions on the left and right sides of the grammars productions. On top of that, a language is recursively enumerable if it is recognized by some Turing machine [?].

2.3.2 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 [?]:

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 or 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 [?]. By using the BNF-notation it is possible to describe a CFG. It is preferred to have a unambiguously grammar. A CFG is ambiguous if a string derived in the grammar has two or more different leftmost derivations [?]. An unambiguously grammar will ensure that a program reading a string using CFG can only read the string in one way. It is worth spending time making a grammar unambiguous, because if the grammar is ambiguous, multiple compilations of a program using that grammar can result in different programs with different meanings and programs yielding different results [?]. But it also worth noting that some grammars are inherently ambiguous, meaning that no matter what it can not become unambiguous.

A CFG is a part of the $LL(k)$ grammar class if it is possible to produce the leftmost derivation of a string by looking at most k tokens ahead in the string. LL algorithms works on the same subset of free grammars which means that LL parsers works on $LL(k)$ grammars. $LL(k)$ means that the grammar needs to be free of left-recursion 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. One property is that $LL(1)$ grammars are fairly easy compared to $LL(k)$ where $k > 1$ to implement because the parser analyzer only has to look one element ahead in order to determine the appropriate parser action. $LL(1)$ is also relatively faster than $LL(k)$ where $k > 1$ because of the same reason: 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 detects the syntax errors faster. LL is also inferior compared to LR in terms of describing a languages 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 been possible in a LL grammar [?] [?].

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 parsers are bottom-up parsers meaning that they begin constructing the abstract tree from its leaf and works its way to the root. LR parsers are generally harder to implement by hand than LL parsers but there exists tools which automatic generates LR parsers for a given grammar. $LR(k)$ grammars allows left recursion which means that the LR grammars are a bigger grammar class than LL . $LALR$ and $SLAR$ is subclasses of the $LR(k)$ grammars which means that $LR(k)$ describes a larger class of languages 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 compare to LR building method [?] [?].

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

2.3.3 Choice of Grammar

The programmer, using the language of this project, could be a hobby programmer, who wants to program a custom drinks machine, but does not possess a high level of experience in programming. This is the reason why it was decided that the grammar should have a high level of readability because this in turn will ensure that it is easier for the programmers to read and understand their programs - this is also useful if the code has to be maintained later on. This on the other hand can decrease the level of write-ability because the programs have to be written in a specific way, and will need to contain some overhead in form of extra words and symbols to mimic a language easier for humans to comprehend.

The method to assign a value to a variable is by typing "*variable* <- "value to assign", without the quotes. This approach has 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 clearly indicated that the value is assigned to the variable, and therefore ensuring readability - especially for the hobby programmer.

When declaring a function it has to be written on the form "function *functionname* return *type* using (*parameter(s)*) begin *statements* return *expresion* end". Where *functionname* is the name of the function that is about to be declared, *type* is the type of value that is returned by the function. *Parameter(s)* are used to parse the function some input values from its call(s). *Statements* is were the function can call other functions, declare variables, calculate and assign values. *Expresion* is were the value of the right type is returned or an expression which result is of the correct type. An example of this can be seen on listing 2.3.

```
1      function DoSomething return int using (int x)
2      begin
3          x <-- x + 1;
4          return x;
5      end
```

Listing 2.3: Example of function declaration in SPLAD

To get a more continuity 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 because the programmer can see what it returns, even if no value was 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 a function. To indicate that a function 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 statement 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 is a suitable solution for the project because it allows the project group to later go back and make the parser by hand instead of using a parser generator 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 2.3.2 for more on the grammars.

2.3.4 The BNF of SPLAD

This section contains the BNF for the programming language of this project; SPLAD. It is clear that the BNF begins with the non-terminal "program". The "program" non-terminal can then be derived in a number of ways, to represent a specific program. The grammar for SPLAD is:

$$\begin{aligned} \langle program \rangle &\rightarrow \langle roots \rangle \\ \langle roots \rangle &\rightarrow \varepsilon \\ &\quad | \quad \langle root \rangle \langle roots \rangle \\ \langle root \rangle &\rightarrow \langle dcl \rangle; \\ &\quad | \quad \langle callid \rangle \langle assign \rangle; \\ &\quad | \quad \langle function \rangle \\ &\quad | \quad \langle COMMENT \rangle \\ \langle dcl \rangle &\rightarrow \langle type \rangle \langle id \rangle \langle dclend \rangle \\ \langle type \rangle &\rightarrow \langle primitivetype \rangle \langle arraytype \rangle \\ \langle primitivetype \rangle &\rightarrow \text{bool} \\ &\quad | \quad \text{double} \\ &\quad | \quad \text{int} \\ &\quad | \quad \text{char} \\ &\quad | \quad \text{container} \\ &\quad | \quad \text{string} \\ \langle arraytype \rangle &\rightarrow [\langle NOTZERODIGIT \rangle] \\ &\quad | \quad \varepsilon \\ \langle id \rangle &\rightarrow \langle LETTER \rangle \\ \langle dclend \rangle &\rightarrow \varepsilon \\ &\quad | \quad \langle assign \rangle \\ \langle assign \rangle &\rightarrow <-- \langle expr \rangle \\ \langle expr \rangle &\rightarrow \langle term \rangle \langle exprend \rangle \\ \langle term \rangle &\rightarrow \langle comp \rangle \langle termend \rangle \\ \langle comp \rangle &\rightarrow \langle factor \rangle \langle compend \rangle \end{aligned}$$

$$\begin{aligned} \langle factor \rangle \rightarrow & (\langle expr \rangle) \\ & | \quad !(\langle expr \rangle) \\ & | \quad \langle callid \rangle \\ & | \quad \langle numeric \rangle \\ & | \quad \langle string \rangle \\ & | \quad \langle functioncall \rangle \\ & | \quad \langle cast \rangle \\ & | \quad \text{LOW} \\ & | \quad \text{HIGH} \\ & | \quad \text{true} \\ & | \quad \text{false} \end{aligned}$$

$$\langle callid \rangle \rightarrow \langle id \rangle \langle arrayiden \rangle$$

$$\begin{aligned} \langle arrayiden \rangle \rightarrow & \langle ARRAYCALL \rangle \\ & | \quad \varepsilon \end{aligned}$$

$$\langle ARRAYCALL \rangle \rightarrow [\langle NOTZERODIGIT \rangle]$$

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

$$\begin{aligned} \langle plusminus \rangle \rightarrow & \varepsilon \\ & | \quad - \\ & | \quad + \end{aligned}$$

$$\begin{aligned} \langle numericend \rangle \rightarrow & \varepsilon \\ & | \quad . \langle DIGIT \rangle \end{aligned}$$

$$\langle string \rangle \rightarrow \langle STRINGTOKEN \rangle$$

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

$$\begin{aligned} \langle callexpr \rangle \rightarrow & \langle subcallexpr \rangle \\ & | \quad \varepsilon \end{aligned}$$

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

$$\begin{aligned} \langle subcallexprend \rangle \rightarrow & , \langle subcallexpr \rangle \\ & | \quad \varepsilon \end{aligned}$$

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

$$\begin{aligned} \langle compend \rangle \rightarrow & \langle comparisonoperator \rangle \langle comp \rangle \\ & | \quad \varepsilon \end{aligned}$$

$$\begin{aligned} \langle comparisonoperator \rangle \rightarrow & > \\ & | \quad < \\ & | \quad <= \\ & | \quad >= \\ & | \quad != \\ & | \quad = \end{aligned}$$

$$\begin{aligned} \langle termend \rangle \rightarrow & \langle termsymbol \rangle \langle term \rangle \\ & | \quad \varepsilon \end{aligned}$$

$\langle \text{termsymbol} \rangle \rightarrow *$
 $\quad |$
 $\quad | \text{ AND}$

$\langle \text{exprend} \rangle \rightarrow \langle \text{exprsymbol} \rangle \langle \text{expr} \rangle$
 $\quad | \quad \varepsilon$

$\langle \text{exprsymbol} \rangle \rightarrow +$
 $\quad | \quad -$
 $\quad | \quad \text{OR}$

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

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

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

$\langle \text{params} \rangle \rightarrow \langle \text{subparams} \rangle$
 $\quad | \quad \varepsilon$

$\langle \text{subparams} \rangle \rightarrow \langle \text{type} \rangle \langle \text{id} \rangle \langle \text{subparamsend} \rangle$

$\langle \text{subparamsend} \rangle \rightarrow , \langle \text{subparams} \rangle$
 $\quad | \quad \varepsilon$

$\langle \text{stmts} \rangle \rightarrow \varepsilon$
 $\quad | \quad \langle \text{stmt} \rangle \langle \text{stmts} \rangle$

$\langle \text{stmt} \rangle \rightarrow \langle \text{callid} \rangle \langle \text{assign} \rangle;$
 $\quad | \quad \langle \text{nontermif} \rangle$
 $\quad | \quad \langle \text{nontermwhile} \rangle$
 $\quad | \quad \langle \text{from} \rangle$
 $\quad | \quad \langle \text{dcl} \rangle;$
 $\quad | \quad \langle \text{functioncall} \rangle;$
 $\quad | \quad \langle \text{nontermswitch} \rangle$
 $\quad | \quad \langle \text{COMMENT} \rangle$

$\langle \text{nontermif} \rangle \rightarrow \text{if}(\langle \text{expr} \rangle) \text{ begin } \langle \text{stmts} \rangle \text{ end } \langle \text{endif} \rangle$

$\langle \text{endif} \rangle \rightarrow \text{else } \langle \text{nontermelse} \rangle$
 $\quad | \quad \varepsilon$

$\langle \text{nontermelse} \rangle \rightarrow \langle \text{nontermif} \rangle$
 $\quad | \quad \text{begin } \langle \text{stmts} \rangle \text{ end}$

$\langle \text{nontermwhile} \rangle \rightarrow \text{while}(\langle \text{expr} \rangle) \text{ begin } \langle \text{stmts} \rangle \text{ end}$

$\langle \text{from} \rangle \rightarrow \text{from } \langle \text{callid} \rangle \langle \text{assign} \rangle \text{ to } \langle \text{expr} \rangle \text{ step } \langle \text{plusminus} \rangle \langle \text{DIGIT} \rangle \text{ begin } \langle \text{stmts} \rangle \text{ end}$

$\langle \text{nontermswitch} \rangle \rightarrow \text{switch } (\langle \text{expr} \rangle) \text{ begin } \langle \text{cases} \rangle \text{ end}$

$\langle \text{cases} \rangle \rightarrow \text{case } \langle \text{expr} \rangle: \langle \text{stmts} \rangle \langle \text{endcase} \rangle$

$$\langle \text{endcase} \rangle \rightarrow \langle \text{cases} \rangle$$

$$\quad | \quad \text{break}; \langle \text{breakend} \rangle$$

$$\quad | \quad \text{default}: \langle \text{stmts} \rangle \text{ break};$$

$$\langle \text{breakend} \rangle \rightarrow \langle \text{cases} \rangle$$

$$\quad | \quad \text{default}: \langle \text{stmts} \rangle \text{ break};$$

$$\quad | \quad \varepsilon$$

$$\langle \text{STRINGTOKEN} \rangle \rightarrow " . * ? "$$

$$\langle \text{LETTER} \rangle \rightarrow [\text{a} - \text{zA} - \text{Z}] ^+$$

$$\langle \text{DIGIT} \rangle \rightarrow [0 - 9] ^+$$

$$\langle \text{NOTZERODIGIT} \rangle \rightarrow [1-9][0-9]^*$$

$$\langle \text{COMMENT} \rangle \rightarrow /* . * ? */$$

2.3.5 EBNF?

2.3.6 Lexicon

2.4 Semantics

In this section the semantics of SPLAD will be described.

2.4.1 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 SPLAD static scoping is used. This means that scopes are computed at compile time, based on the program text input. The main reason for this, is that programs for the Arduino platform is mainly written in C, which also uses static scoping. This makes the compilation from SPLAD to C simpler for the compiler [?]. 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 [?].

2.4.2 Transition Rules

2.4.3 Type Rules

This section contains the type rules for the comparison operator.

Type rule for $<$, $>$, $<=$, $>=$:

" $E_1 (<, > <=, >=) E_2$ " is type correct and of type boolean if E_1 and E_2 are type correct and of type integer, double.

Type rule for $!=$, $=$:

" $E_1 (!=, =) E_2$ " is type correct and of type boolean if E_1 and E_2 are type correct and of type integer, double, or if E_1 and E_2 are of the same type of either char or string.

Type rule for $+$, $-$, $*$: " $E_1 (+, -, *) E_2$ " is type correct and of type integer or double if E_1 and E_2 are type correct and of type integer or double.

Type rule for $/$: " $E_1 (/) E_2$ " is type correct and of type integer or double if E_1 and E_2 are type correct and of type integer or double and E_2 does not have the value of zero.

Here the type rules of assign will be described:

" $E_1 <- E_2$ " is type correct if E_1 and E_2 are of the same of type integer, double, char or string.

Here the type rules of loops will be described.

Type rule of 'while'-statement: "while E begin C end" is type correct if E of type boolean and C are type correct.

Type rule of 'from to'-statement: "from E_1 to E_2 begin C end" are type correct if E_1 and E_2 are type correct and of type integer, and C are type correct.

This is the type rules for 'if'-statement: "if(E) begin C end" is type correct if E are type correct and of type boolean, and C are type correct.

Here the type rules for switch/case will be described:

"switch (E) begin case E_1 : C_1 break; ... case E_n : C_n break; default: C_d break; end" is type correct if E, $E_1...E_n$ are type correct and of type integer, double, char or string and are the same type, and $C_1...C_n$ and C_d are type correct.

2.5 Code Examples

Implementation 3

3.1 Design Criteria

3.2 Architecture

3.2.1 Hardware

This section will be about the hardware components used in this project, describing them and the reasons they are used in this project. In the description it will be looked at the more basic technical specification that will be relevant for this project.

Hardware platform

Arduino UNO is a powerful micro controller board which provides the user with ways to communicate with other components such as LCDs, diodes, sensors and other electronic bricks(building blocks) which is a desirable feature in this project. Arduino uses the ATmega328 chip which provides more memory than its predecessors [?].

There exists alternatives to Arduino product which could be considered for this project. Teensy is similar to Arduino in many ways but Teensy differences in the actual size of the product. Teensy is also cheaper than Arduino but does require soldering for simple set-ups where Arduino comes with a board and pin-ports which means that it requires little pre work before using it [?]. Seeeduino is a near replica of Arduino, an example could be Seeeduino Stalker which offers features as SD-card slot, flat-coin battery holder and X-bee module-headers. X-bee is a module for radio communication between one or more of these modules. Seeeduino is compatible with the same components as Arduino that makes it suited for acting as a replacement [?]. Netduino is a faster version of Arduino but it comes at higher cost. Netduino also require the .net framework so it will only work together with windows operating systems. Netduino uses a micro-USB instead of regular USB as Arduino does [?].

Arduino is more accessible because the Aalborg university already has some in stock that could be used where the other alternatives have to be bought first. Arduino, Teensy and Seeeduino are all compatible with the other equipments that will be used in this project. Netduino is limited to the .net framework where Arduino and the other alternatives are more flexible and therefore more ideal because they work with more platforms. So it comes down to that Arduino have all the necessary features and is the most convenient one to obtain. The project group has found this platform suited for this project based on these reflections.

Arduino UNO board has 14 digital input/output pins where six of them can emulate an analogy output through PWM (Pulse-Width modulation) which are available on the Arduino board. The Arduino board also provides the user with six analogy inputs which enables the reading of a alternating current and provides the user with the currents voltages. These pins can be used to control or perform readings on other components and in that way provides interaction with the environment around the board. The Arduino board is also mounted with an USB-port and jack socket. The board can be hooked up with a USB cable or an AC-to-DC (Alternating Current to Direct Current) adapter through the jack socket to power the unit. Arduino UNO operates at 5v (volts) but the recommend range is 7-12v because lower current than 7v may cause instability if the unit needs to provide a lot of power to the attached electronic brick. The USB is also used to program the unit with the desired program through a computer [?].

Programs for Arduino are commonly made in Arduino's own language that are based on C and C++. The produces of the Arduino platform provides a development environment (Arduino IDE) that makes it possible to write and then simply upload the code to the connected Arduino platform. This process also provides a library with functions to communicate with the platform and compatible components [?]. Arduino is suited for this project because it makes it possible to demonstrate the language and illustrate that the translation works.

RFID

To administrate the users collection of purchased drinks the plan is to store the number and kind of drinks on an RFID tag that the customer then can use at the drink machine to get their drinks served.

RFID (Radio Frequency IDentification) is used to identify individual objects using radio waves. The communication between the reader and the RFID tag can go both ways, and it is possible to both read and write to most tag types. The objects that are able to be read differs a lot. It can be clothes, food, documents, pets, packaging and a lot of other kinds. All tags contains a unique ID that can in no way be changed once made. This ID is used to identify an individual tag. Tags can be either passive or active. Passive tags do not do anything until a signal from a reader transfers energy to the tag once activated it sends a signal back in return. Active tags have a power source and therefore is able to send a signal on their own, making the read-distance greater. The tags can also be either *read only tag* or *read/write tag*. A *read only tag* only sends its ID back when it connects with a reader, while a *read/write tag* have a memory for storing additional information it then sends with the ID [?].

Other components

The demonstration situation will require something to illustrate more advanced parts of theoretical machine. The plan is to use LEDs (light emitting diode) to illustrate the different function of the machine, when they are active or inactive. The LED is made of a semiconductor which produces a light when a current runs through the unit.

LEDs are normally easy to use by simply running a current the correct way through the LED. The reason why LEDs are being using instead of making the machine is that there are neither time for it nor is it the main focus of this project.

It would also be good to be able to print a form of text to the customer. To do this there will be used an LCD 16-pin (Liquid Crystal Display). Arduino's Liquid Crystal library

provides the functions to write to LCD so no low level code is needed to communicate with the LCD [?].

As input switches/buttons will be used that will allow interaction with the program at runtime. The switches will illustrate a more advanced control unit but in the project switches will be sufficient.

3.2.2 Overview of the Compiler

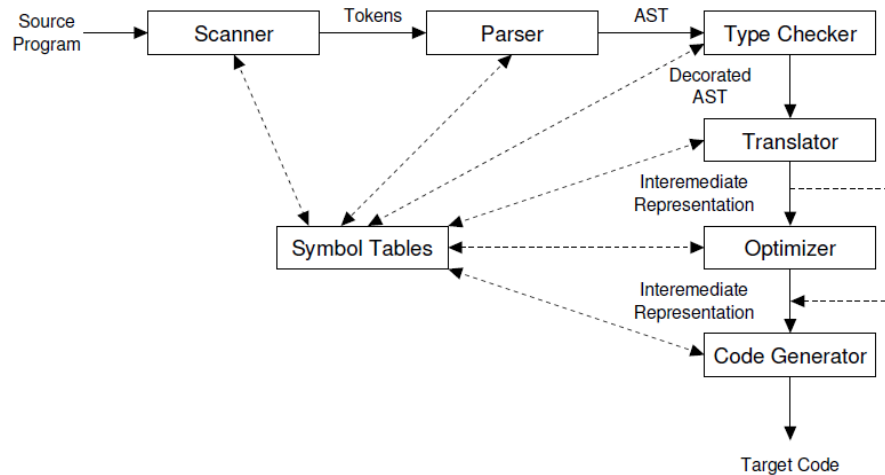


Figure 3.1: This is an overview of how the compiler is structured. The figure is from [?].

Figure 3.1 shows an overview of each of the different phases in the compiler, what each phase requires as input, and what each step returns to the next phase.

A compiler is a fundamental part of modern computing. Their job is to translate programming language into machine language. A compiler allows programs to make a virtual computer to ignore the machine-dependency details of machine language and therefor be portable across different computers [?].

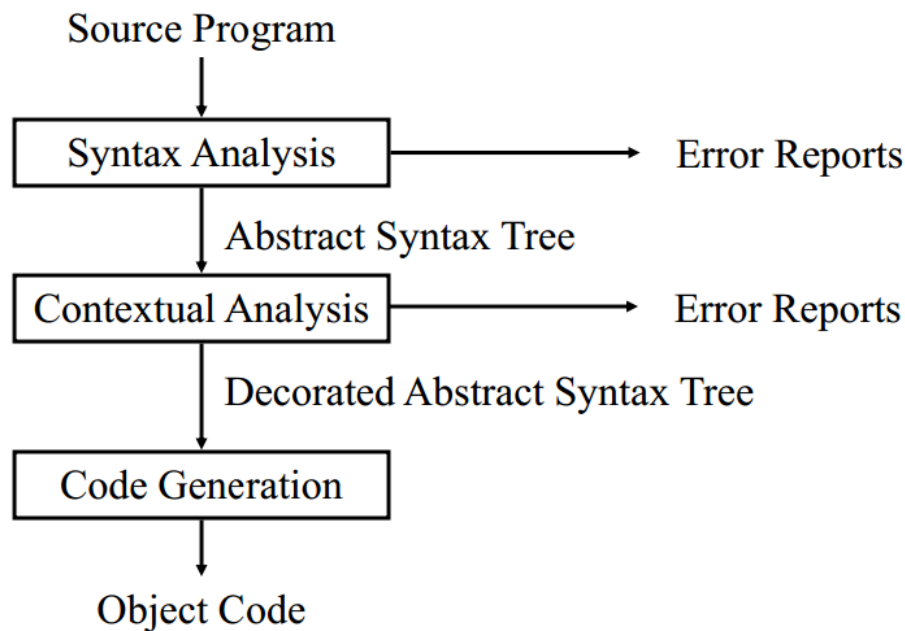


Figure 3.2: This is another overview of how the compiler is structured.

A compiler consists of 3 different phases. The different phases roughly correspond to the different parts in a language specification which can be seen on figure 3.2. The syntax analysis correspond to the syntax, the contextual analysis to the contextual constraints and the code generation phase roughly corresponds to the semantics.

Given a simple compiler it will go through more than three phases. This can be seen on figure 3.1. In the syntax analysis phase the compiler consists of a scanner and a parser. The scanner takes the source program and transforms it into a stream of tokens. The parser then uses the tokens to create an abstract syntax tree (AST). In the contextual analysis a symbol table is created from the abstract syntax tree. At the end, the semantic analysis decorate the AST, and translates this into the target language.

The different phases will be described more thorough later in the rapport.

3.2.3 Language Processing Strategy

3.2.4 Compilation Passes

3.2.5 Abstract Syntax Trees

The parser generates an abstract syntax tree (AST) [?], which is an abstract data type describing the structure of the source program. This means that the AST contains information about which constructs the source program contains. More specifically, each node in the AST represent a construct in the source language, for example an 'if'-block.

When the AST has been generated, it is decorated with types by the type checker. The type checker traverses the AST, and checks the static semantics of each node, which means that it verifies that the node represent valid constructs. If each node is correct it is returned to the translator [?]. The translator then uses the AST to an intermediate representation (IR code), which is used in the later phases of the compiler.

3.2.6 The Visitor Pattern

3.3 Syntactic Analysis

3.3.1 Semantic Analysis

3.3.2 Known lexers and parsers

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

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. [?]

Flex Alternativ to lex [?]

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 [?]

Jlex Based on lex but used for java. [?]

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. [?]

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

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 [?].

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. [?]

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[?]. The tree builder that accompanies it, JJTree, constructs its trees from the bottom uplex[?].

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.

3.3.3 ANTLR

ANTLR (ANother Tool for Language Recognition) is a tool for generating a parser or lexer-parser from a given grammar. The ANTLR starts by generating the lexer based on the lexer rules that are define in the grammar. Lexer rules is written with an upper-case beginnings letter so that ANTLR can distinguish between lexer rules and parser rules [?]. In the SPLAD language it has been decided that a lexer rule should all be written in upper-case. This is done in order to better distinguish between lexer parser rules when writing or reading the grammar. Lexer starts from the top of the rules and work its way down through the rules, meaning that it will try to generate tokens from the very first rule and work its way down until it meets a possible match between a given input and a rule. Because of this the most complex rules should be placed first in the grammar in

order for the lexer to generate the correct tokens. The token stream from the lexer are then parsed following the parser rules that have been defined in the grammar. Parser rules are all written with lower-case in contrast to the lexer rules. ANTLR works with *LL(*)* grammars which means that the parser uses left-most derivation to parse the token stream. ANTLR can generate a abstract syntax tree for the grammar by incorporating specific operators in the grammar that tells if an element should be a root node of a subtree with its children or if an element should be left out of the tree construction. ANTLR also allows the use of rewritten rules to generate a tree from the given grammar [?].

```
1 function setup return nothing using()
2 begin
3   /*Do something*/
4   return nothing;
5 end
6
7 function LCDPrint return nothing using(string text)
8 begin
9   /*Function to write a string to the LCD connected to the arduino
10      */
11   return nothing;
12 end
13
14 function makedrink return nothing using()
15 begin
16   string message <-- "Hello World!";
17   call LCDPrint(message);
18   return nothing;
19 end
```

Listing 3.1: Here the code for the simple "Hello world" program can be seen.

The ANTLR library comes with a tool for testing the generated lexer and parser. The tool, (test rig), allows for parsing some code and get it represented in a GUI or tree representation. This have been done for a simple "Hello world" program, see listing ??, to show how a program is parsed and represented using the GUI option.

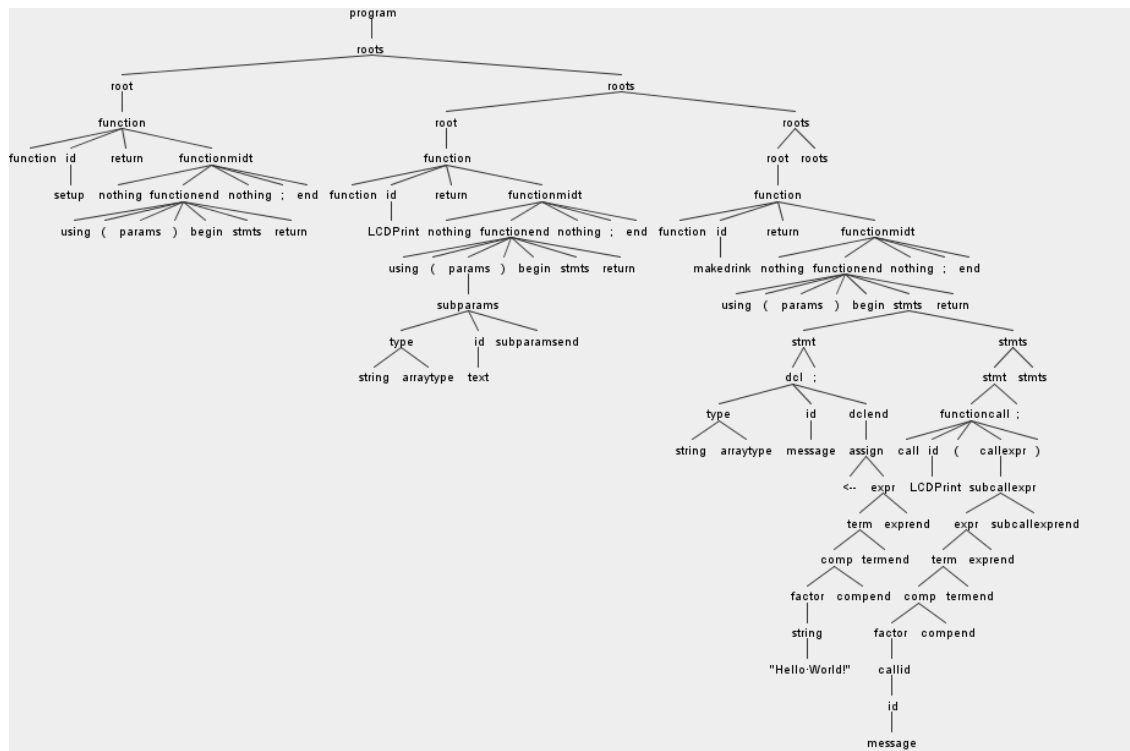


Figure 3.3: On this figure the parse tree for the parsed program "Hello world" can be seen.

On figure 3.3 a parse tree for the "Hello world" program can be seen. Derivations are illustrated by being children of the parent node, see section for more about parse trees.

3.3.4 Lexical Analyzer

A lexical analyzer reads the input file, and returns a series of tokens based on the input [?]. 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 3.1. Formally a token consists of two parts: The token type, and the token value [?]. As an example the IDENT token seen on 3.2 has the token type IDENT and the value 'c'.

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

Table 3.1: 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 3.2: Simple example of code

As an example the lines of code seen on listing 3.2 might be read as the tokens seen on table 3.2.

Token	Lexeme
IDENT	c
ASSIGN	=
DIGIT	42
SEMICOLON	;

Table 3.2: 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 [?].

Scanner Class Generation

3.3.5 Tokens

For a compiler to able to distinguish between variables names and types the compiler will need some rules to describe the difference between them. This is done by reserving the words, called keywords, which are used to describe types, the beginnings and endings of blocks, and declaration of statements. A variable may not be named the same as any of the keywords since the compiler can not distinguish if it is a variable name or a reserved keyword.

Reserved Keywords

The reserved keywords for SPLAD can be seen on table 3.3.

bool	container	while
int	from	to
double	step	default
char	break	case
string	switch	nothing
OR	HIGH	LOW
AND	return	using
true	function	else
false	if	end
begin		

Table 3.3: The reserved keywords in SPLAD

This list is used to keep track of which words are going to be reserved and in that way provide an overview for the programmer.

Token Specification

A parser needs a stream of tokens to parse a program correctly. These tokens are generated by a lexer which reads a stream of input symbols and from a given set of rules, makes the corresponding tokens. A token specification is used to describe the rules the lexer need in the construction of tokens. Token specification are expressed in way related to regular expressions [?]. Regular expressions are strong in describing patterns which is the core of token production [?]. The tokens used for this project can be seen on table 3.4.

PRIMITIVE TYPE	'int' 'double' 'bool' 'char' 'container' 'string'
STRING TOKEN \rightarrow " ... "	
DIGIT	$[0 - 9]^+$
NOTZERODIGIT	$[1 - 9][0 - 9]^*$
LETTER	$[A - Z a - z]^+$
COMMENT	$/* \dots */$
WHITESPACE	
r	
n	
t	
OTHER	ε

Table 3.4: The tokens in SPLAD

Further work would be making a lexer to generate a token for the parser. Another options was to find a suited tool for generating a lexer for the given rules. This is a valid option because making a lexer can be automated and therefore already exists a lot of good lexer generators that can be used, see section 3.3.2.

3.3.6 Parser

A parser takes the tokens from the scanner and use them to create an abstract syntax tree. It also checks if the stream of tokens conforms to the syntax specification, usually written formal using context-free grammar (CFG).

The main purpose of the parser is to analyze the tokens and check if the source program is written in the correct syntax. If this is not the case the parser should show a message describing the error. The parser will at the end create an abstract syntax tree.

Generally there are two different approaches to parsing: top-down and bottom-up. Before describing the different approaches to parsing, it is worth to describe derivation shortly. Derivations is how the parser will create the parse tree. Either it will be built leftmost or it will be built rightmost. Leftmost-derivation is where the parser will take the terminal that is most to the left, and create a derivation for that. A rightmost-derivation is the opposite: The parser chooses the first terminal from the right, and creates a derivation for that.

Top-down Parsers

Then top-down parser starts at the root and works its way to the leaves in a depth-first manner, doing a pre-order traversal of the parse tree. This is done by reading tokens from left to right using a leftmost derivation. Furthermore top-down parsers can be split into table-driven LL and recursive descent parse algorithms.

Table-driven LL Parsers Uses a parse table to determine what to do next. The entries in the parse table is determined by the particular LL(k) grammar. The parser then searches the table to see what to do.

Recursive-descent Parser The recursive-descent parsers consists of mutually recursive parsing routines. Each of the non-terminals in the grammar has a parsing procedure that determines if the token stream contains a sequence of tokens derivable from that non-terminal.

Bottom-Up Parsers

A bottom-up parser has to do a post-order traversal of the parse tree, meaning that it starts from the leaves and works towards the root. A bottom-up parser is more powerful and efficient than a top-down parser, but not as simple.

LR A LR parser reads from left to right and because it is a bottom-up parser it uses a reversed rightmost derivation which means it takes terminals and turn them into non-terminals. It is as the LL parser driven from a parse table. The biggest difference is how it is derived and how the parse table is handled.

LALR A LALR(Lookahead Ahead LR) parser is one of the most commonly used algorithms today, because it is a powerful algorithm but do not need a very large parse table. It works like the LR parser.

3.4 Contextual Analysis

3.4.1 Scope Checking

3.4.2 Type Checking

3.5 Code Generation

3.5.1 Declaration

3.5.2 Commands

3.5.3 Expressions

Conclusion 4

FiXme Fatal: ko
mangler

List of Corrections

Fatal: synopsis mangler	iii
Fatal: prolog mangler	v
Fatal: indledning mangler	1
Fatal: ItLatToC	6
Fatal: Skal det med?	12
Fatal: Skal det med?	12
Fatal: Få lavet konkret henvendelse når opbygningen på rapporten er fastsat.	18
Fatal: skal det med?	18
Fatal: Skal det med?	18
Fatal: ref til afsnittet som parsetræer	22
Fatal: konklusion mangler	27