

UNIVERSITY OF ZAGREB  
FACULTY OF ELECTRICAL ENGINEERING AND  
COMPUTING

MASTER THESIS No. 2568

# **Design of a strongly-typed programming language**

Petar Mihalj

Zagreb, May 2021.

## MASTER THESIS ASSIGNMENT No. 2568

Student: **Petar Mihalj (0036497900)**

Study: Computing

Profile: Computer Science

Mentor: asst. prof. Ante Đerek

Title: **Design of a strongly-typed programming language**

### Description:

Programming languages play a central role in the development of software. However, there is no silver bullet programming language, which would be efficient, easy to use, portable and consistent. The goal of this graduate thesis is to design a new programming language with the emphasis on how strong typing can address common drawbacks. The language specification should be described in a formal form. Design should be followed-up by a careful analysis of some aspects of the language, for example to prove the desirable properties of the type system. The final goal of the thesis is to implement the tool chain, which primarily consists of a compiler, but can include other tools, such as language servers and formatters. The thesis should be accompanied with the source code of all developed software. All references should be clearly cited. Any assistance received should be clearly acknowledged.

Submission date: 28 June 2021

## **DIPLOMSKI ZADATAK br. 2568**

Pristupnik: **Petar Mihalj (0036497900)**

Studij: Računarstvo

Profil: Računarska znanost

Mentor: doc. dr. sc. Ante Đerek

Zadatak: **Oblikovanje strogo tipiziranog programskog jezika**

### Opis zadatka:

Programski jezici centralna su karika u razvoju programske potpore. Ipak, ne postoji programski jezik najbolji za sve primjene, u isto vrijeme efikasan, jednostavan, prenosiv i konzistentan. Cilj ovog diplomskog rada dizajnirati je novi programski jezik, s naglaskom na prednosti strogih tipova u rješavanju čestih problema. Specifikacija jezika treba biti opisana u formalnom obliku. Nakon dizajna jezika potrebno je provesti pažljivu analizu nekih svojstava jezika, na primjer pokazati poželjna svojstva sustava tipova. Posljednji cilj rada implementacija je skupa alata, kojeg primarno čini prevoditelj, ali koji može uključivati i druge alate, poput jezičnih poslužitelja (language servers) i formatera koda. Rad treba uključivati izvorni kod razvijenih alata. Sve reference trebaju biti jasno citirane. Sva primljena pomoć treba biti jasno obznanjena.

Rok za predaju rada: 28. lipnja 2021.

*I thank everybody...*

# CONTENTS

<b>1. Introduction</b>	<b>1</b>
1.1. Paper organization . . . . .	2
<b>2. The AGT Programming Language - a gentle introduction</b>	<b>3</b>
2.1. Hello world! . . . . .	3
2.2. In depth . . . . .	5
2.2.1. Lexical analysis . . . . .	5
2.2.2. Syntactical analysis . . . . .	8
2.2.3. Semantic analysis . . . . .	8
2.2.4. Type inference . . . . .	8
<b>3. Conclusion</b>	<b>9</b>
<b>Bibliography</b>	<b>10</b>

# 1. Introduction

The programming language developed as a part of this thesis is called AGT. The name AGT stands for an unfortunate fact that the names of all precious stones are taken by other programming languages, hence "All Gems are Taken".

AGT is a statically and strongly typed language, with a highly expressive type system. The type system is used for three main purposes:

1. types determine the size of values in memory
2. types enable polymorphic behaviour (of both builtin and user-defined (struct) types)
3. types allow for compile-time computation

The language also allows for implementation of object lifetime constructs, such as those seen in languages like C++ or Rust; object creation, copying and destruction. These constructs can be used to implement various object ownership semantics, which will be demonstrated in later chapters.

The reference AGT compiler (AGTC) produces executable binaries only. This is in sharp contrast to some other languages, which can produce object files, to later be linked into executables for a particular runtime, or used to augment the runtime environment itself (kernel modules, for example). This decision was made primarily because the goal of the thesis was to study the language from the application programmers perspective. This goal implies that AGT is running on a rich runtime environment, where heap memory management and standard I/O are available. Future updates to AGT specification are meant to allow AGT to be compiled into linkable files compatible with various platform specific ABIs.

The compiler frontend is implemented using Python3 programming language, while the backend uses LLVM (Lattner und Adve, 2004) compiler infrastructure. Allowing for some simplification, AGT is first compiled into LLVM intermediate representation

(LLVM-IR), and then is turned into executable by a chain of tools for compiling LLVM-IR and linking the resulting object files into executables.

## **1.1. Paper organization**

This master thesis will gradually introduce the reader to AGT, starting from the simplified introduction to syntax and semantics, with an emphasis on type system and inference. This introduction will be augmented with simple examples of AGT code.

After the surface-level exploration of the language, we will study the language in depth. This part will have many references to the AGTC reference implementation, rather than to a separate formal specification. Since AGT is still not a finalized product, formal specification does not exist, and its functionality is determined by a reference compiler. We will mostly focus on the type system, since it is unorthodox and is the main contribution of the author. In this part, the author will justify the decisions which have been made on all levels of AGT design process.

In the last part, we will address the issues AGT has both on definition and implementation levels, and comment on the various features that can be improved or added to the language.

## 2. The AGT Programming Language - a gentle introduction

Before diving into

### 2.1. Hello world!

Let us first check out an example AGT program:

```
1 fn outnl(i){
2     out(i);
3     out('\n');
4 }
5
6 fn mul(a, b) -> a{
7     return a*b;
8 }
9
10 fn main()->i32{
11     let a = in<i32>();
12     let b = in<i8>();
13
14     let c = mul(a, cast<i32>(b));
15     outnl(a);
16     outnl(b);
17     outnl(c);
18 }
```

Every AGT program has to have a function called `main`, which returns a 32-bit-wide integer. Some functions, like `outnl`, don't return anything. It is important to note that you can't specify that a function returns `void` like in C for example.



A `let <x> = <y>;` construct is *initialization assignment* statement. This roughly translates to an allocation of memory on a stack (which is allocated at the newly generated *location* of `<x>`), and copying the value of expression `<y>` to the location of identifier `<x>`.

Note that call of the function `in` is *parametrized* by a *type argument* `i32`. Type arguments are in contrast to the *value arguments*, which carry value too, along with a type. The function `in` is a builtin, and the `i32` is used to signal the compiler that we want an instance of this function which returns an `i32`, rather than a `char`, for example. Of course, user-defined functions can also have *type parameters*, along with *value parameters*; we will get to these later.

Notice that we used the `cast<i32>` builtin function to convert `b` to `i32` before passing it to function `mul`. If we didn't do that, the compiler would signal that it can't compute the expression `a*b`, since multiplication is predefined only for integer types of same size. Apart from the `i8` and `i32` we used in this example, `i16` and `i32` are also available.

There is a certain feature that you might have missed while inspecting the code; namely the absence of types in the definition of function `mul`. While some dynamically typed languages (ex. Python) allow for object of any type to be passed to the function, AGT behaves quite differently.

Every time AGTC (AGT compiler) encounters a function call, it tries to infer a *function type*. Function type is inferred from function definitions. Consequently, you can think of the definitions in source code more as *templates* for synthesis of actual code, than the actual representation of code. These definitions lack types, and can't be considered in isolation. In this example, the call of function `mul` causes the compiler to try to infer a function type.

It is important to emphasize that one function definition can be a source of many function types. The `outnl` (output with newline) function is a prime example of this behaviour. It is called three times, with argument of types `i32`, `i8` and `i32`. Thus, the compiler has to infer two function types, one which can be fed with an `i8`, and one with `i32`. The success of this inference corresponds with compilers ability to infer function body of the given function definition. This will be possible if and only if the compiler can also infer the function type `out` function which can take ... The property of `outnl` function definition which allows it to be a source for inference of two distinct function types, given that it *makes sense* for these types (regarding the function body), is called **parametric polymorphism**.

## 2.2. In depth

AGT compilation process consists of 4 main phases:

1. Lexical analysis - tokenization of input source code text
2. Syntactical analysis - grouping of tokens in an abstract syntax tree
3. Semantical analysis - transforming the syntax tree into a context aware semantical syntax tree
4. Type inference - inferring the main function type and all other required types

The process is successful if the type inference step successfully infers the main function type.

We will now describe each phase in great detail, contrasting the newly defined terms against their use in other languages, since the subtleties sometimes are different.

### 2.2.1. Lexical analysis

The AGT Lexical analysis phase is a standard regular grammar parsing process. Lexical analyzer is fed a stream of characters, which get grouped into tokens, according to regular expressions and ambiguity resolution rules.

Lexical analyzer tool used for the AGTC was PLY - Python Lex-Yacc (Beazley, a). Python Lex-Yacc is a tool which tries to bring the functionality of well known lexical and syntactical analyzers Lex and Yacc into Python. The main difference is that PLY does not *generate* the concrete lexer and parser. Instead, a PLY user defines the rules fed to PLY by specifying programming constructs like functions that correspond to these rules. PLY then inspects these constructs using reflective capabilities of python, and acts according to these rules when performing lexical or syntactical analysis.

Lets check out some PLY lexer rules used in AGTC.

```
1 class Lexer():
2     states = (
3         ('mlc', 'exclusive'), ('strchar', 'exclusive'),
4     )
5     #...
6     tokens = (
7         'INTL', 'BOOLL', 'ID',
8         'IF', 'ELSE', 'BREAK', 'FOR', 'WHILE', 'RETURN',
```

```

9          #...
10     )
11     t_ADD = r'\+'
12     t_SUB = r'\-'
13     t_MUL = r'\*'
14     #...
15     reserved = {
16         'type': 'TYPE',
17         'fn': 'FN',
18         'let': 'LET',
19         #...
20     }
21     def t_ID(self, t):
22         r'[a-zA-Z_][a-zA-Z_0-9]*'
23         t.type = self.reserved.get(t.value, 'ID')
24         return t
25
26     def t_INTL(self, t):
27         r'(\d+)(li8|li16|li32|li64)'
28         return t
29     #...

```

Rules for simple operators are given in a string form, in a form of `t_<TOKEN> = <REGEX>`.

Rules for more complex tokens, like int literals or identifiers are defined using functions. Those same functions can alter the behaviour of the lexer when such a token is parsed. For example, whenever an identifier is parsed, the function `t_ID` looks up whether the identifier is in a list of reserved words, and if it is, changes the tokens type. This neat trick greatly simplifies the regular expressions (Beazley, b).

PLY documentation states that priority is given to regular expressions defined by functions, in the order they are given, and then to string defined rules, according to their length.

We also employ the *state* feature of PLY, which allows us to change the lexer's behaviour when it encounters, for example, a beginning of a multiline comment. In that particular case, the lexer parses and caches every symbol until it gets to the end of that comment. The state feature is also used for string literals.

An interested reader will find more details about the implementation in the source code supplied with this paper.

This is a great place to introduce reader to special meaning of arithmetic operators. When an expression such as `a + 1` is evaluated, the function call `__add__(a, 1)` is evaluated instead. Same goes for all operators, so now we give their equivalent *dunder* functions (the name comes from the Python community).

Arithmetic operators and their correspondant functions methods are:

- `+` : `__add__`
- `-` : `__sub__`
- `*` : `__mul__`
- `/` : `__div__`
- `%` : `__mod__`

Comparison operators and their correspondant dunder functions are:

- `==` : `__eq__`
- `!=` : `__ne__`
- `>` : `__gt__`
- `<` : `__lt__`
- `>=` : `__ge__`
- `<=` : `__le__`

The strict less or greater operators include the exclamation mark in order to prevent some common ambiguities in lexing (which occur due to the presence of angle brackets).

Finally, the last type of builtins are boolean operators.

- `&` : `__and__`
- `|` : `__or__`
- `~` : `__not__`

Note that, even though we call the last group the *boolean* operators, those operators can be evaluated on operands which are not booleans, as long as the corresponding dunder functions are defined in the program. The same goes for arithmetic and comparison operators; a programmer can use builtins functions for all integer types, but can define additional dunder functions for `rectangle` objects, for example.

Both type of comments, `/*multiline comment*/` and `//singleline comment` are discarded after being parsed.

String and character literals are given in double and single quotes, respectively. They can include escaped characters: `(\0, \n, \t, \", \')`. Character literal must consist of exactly one character or escaped character.

Identifiers consist of a letter or an underscore, followed by any number of letters, underscores or digits. Reserved words are excluded from identifiers.

Reserved words are:

- 'fn', 'struct',
- 'let', 'type',
- 'if', 'else', 'break', 'for', 'while', 'return',
- 'true', 'True', 'false', 'False',

Int literals are characterized by both the value and the size of integer containing this value. Possible integer sizes are 8, 16, 32 and 64.

Example int literals are:

- 5 (implicitly i32)
- 5i32
- 5i6
- 5i8
- 5i64
- 324243 (implicitly i32)

### **2.2.2. Syntactical analysis**

### **2.2.3. Semantic analysis**

### **2.2.4. Type inference**

## **3. Conclusion**

Zaključak.

# BIBLIOGRAPHY

- [Beazley a] BEAZLEY, David: *PLY (Python Lex-Yacc)*. – URL <http://www.dabeaz.com/ply/>
- [Beazley b] BEAZLEY, David: *PLY (Python Lex-Yacc) documentation*. – URL <https://ply.readthedocs.io/en/latest/index.html>
- [Lattner und Adve 2004] LATTNER, Chris ; ADVE, Vikram: *LLVM: A compilation framework for lifelong program analysis transformation*. 2004

## **Design of a strongly-typed programming language**

### **Sažetak**

Sažetak na hrvatskom jeziku.

**Ključne riječi:** Ključne riječi, odvojene zarezima.

### **Title**

### **Abstract**

Abstract.

**Keywords:** Keywords.