

Design and implementation of a programming language with LLVM

Josep Maria Domingo Catafal

Abstract— This thesis presents the design and development of a new programming language called *Craft*, using Rust and LLVM. The goal of the project is to gain a deeper understanding of how compilers work by creating one from scratch using these tools. The language is designed to be simple and easy to understand, but at the same time, it aims to be fast and efficient, so some sacrifices have to be made. The thesis covers the design and implementation of the language, including the lexer, parser, and code generation. The project also includes a discussion of the challenges encountered during development and suggestions for future work. Overall, the project serves as a valuable learning experience for understanding the inner workings of compilers and the capabilities of Rust and LLVM.

Keywords— Programming Language, LLVM, SSA, Strongly Typed, Compiled

Resum— Aquesta tesi presenta el disseny i desenvolupament d'un nou llenguatge de programació anomenat *Craft*, utilitzant Rust i LLVM. L'objectiu del projecte és aconseguir una comprensió més profunda de com funcionen els compiladors creant-ne un des de zero amb aquestes eines. El llenguatge està dissenyat per ser senzill i fàcil d'entendre, però al mateix temps, pretén ser ràpid i eficient, per la qual cosa s'han de fer alguns sacrificis. La tesi cobreix el disseny i la implementació del llenguatge, incloent-hi el lexer, el parser i la generació de codi. El projecte també inclou una discussió sobre els reptes trobats durant el desenvolupament i suggeriments per a treballs futurs. En general, el projecte serveix com una valuosa experiència d'aprenentatge per comprendre el funcionament intern dels compiladors i les capacitats de Rust i LLVM.

Paraules clau— Llenguatge de programació, LLVM, SSA, Fortament tipat, Compilat

1 INTRODUCTION

HISTORICALLY, there has always been a dilemma between the speed of execution, and speed of development. Some languages are easy to program: they allow the programmer to not worry about low-level concepts such as memory management, and create abstractions that streamline the development. The problem is that these abstractions limit language efficiency, and create slower programs. Another reason that allows speeding up development is dynamic typing, as it frees the user from the mental overhead that comes with deciding the type that should be used. But this also has its own disadvantages, since it is very likely that you will encounter errors in timing of execution. These languages tend to be interpreted in order to save money it waits for the programmer, but it has an impact on the execution performance of the program if we compare it to compiled languages. Python, for example,

would be one of the largest representatives of this group of languages. On the other side of the currency and we have languages like C, which have almost no abstraction and the programmer must be aware of what he is doing in every line of code he writes. they are languages that provide very good performance, but slow down the development, as the programmer must take into account many concepts of low level. An additional problem is that when managing shape memory manual, it opens the door to a lot of runtime errors in the form of memory leaks. Currently, there are languages like Rust that solve this these memory management issues without losing performance, but the development is still slow and the compilation time long. these languages tend to be strongly typed, which reduces errors in runtime, but slows down development.

• E-mail de contacte: jdomingocatafal@gmail.com

• Menció realitzada: Computació

• Treball tutoritzat per: Javier Sanchez Pujadas (Ciències de la Computació)

• Curs 2022/23

2 GOALS

3 STATE OF THE ART

3.1 Programming Languages

3.2 Compilers

4 METHODOLOGY

5 TECH STACK

Writing a compiler does not require many tools, but some of them can help a lot in paving the road. For this compiler we are going to use two tools: The Rust Programming Language and LLVM.

We are also going to use two tools for managing the source code: Git with GitHub for hosting.

5.1 Rust

Rust is a systems programming language that was first released in 2010. It was developed by the Mozilla Foundation with the goal of creating a safe and concurrent language that would be suitable for low-level systems programming tasks, such as operating systems, and performance critical programs, like a browser engine. One of the major features of the language is its guaranteed memory-safety (and thread-safety) without requiring the use of a garbage collector or reference counting (and thus not compromising on performance). This combined with its powerful type system, help catch a lot of bugs at compile-time.

For this reason we are going to use this language. It's also really comfortable to use and comes with a lot of great tooling like cargo, which is the command line tool used for compiling, managing dependencies, etc. and clippy, a linter that comes built in and gives great hints on how to improve the code.

5.2 LLVM

LLVM is a tool chain for building compilers, i.e. a set of tools different ones that help us in implementing compilers. It was created in 2003 by Chris Lattner (also creator of the Swift programming language) and has the support of companies like Apple (LLVM is a part integral part of XCode and Swift for iOS application development), Google, IBM or Intel. Currently, there are several mainstream programming languages that use it, such as C/C++ (via the Clang compiler, a alternative to GCC), Rust, Swift, Crystal...

As we said LLVM has a lot of tools, but among all them, the LLVM Core libraries are the most important and particularly relevant for us.

We will be referring to them as LLVM from now on for simplicity.

LLVM will allow us to generate assembly for a lot of different architectures without any extra effort. It can even generate Web Assembly, which allows us to run the language in modern web browsers. The Craft compiler will generate LLVM IR (Intermediete Representation), and then it will be piped to LLVM, which will take the IR, apply transformations to it, in order to optimize it, and then the assembly for the target architecture will be generated.

Generating LLVM IR instead of assembly, also frees us from some headaches. For example, since LLVM is architecture independent, when we generate the code, we don't need to worry about the number of registers, as we have an unlimited number of virtual registers, which LLVM will later map to the registers of the corresponding architecture.

As we discussed, LLVM also applies optimizations to the generated code, like dead code elimination, constant folding, loop unrolling, etc. However, in order for LLVM to perform these optimizations, we will have to generate the code in SSA (Static single-assignment). This means that we can only assign a value to a variable once. If we need to reassign a value, a new variable has to be created that replaces the other one. The reason SSA is used is because it makes applying optimizations a lot easier.

6 GIT AND GITHUB

For version control Git was chosen since it's the industry standard and one of the most powerful tools out there. The repo is hosted on GitHub which is great for open source projects and allows us to use GitHub Actions. The repository is configured such as that on every Pull Request or commit to the master branch, a GitHub action runs that compiles the code, runs Clippy, runs the tests and checks the formatting of the code. This way we prevent broken code to enter the stable branch.

7 DEVELOPMENT

8 ARCHITECTURE

Most compilers are divided into two parts: the front-end and the back-end. The front-end is the part of the compilers that takes the source code and transforms it into an intermediate representation than will later be transformed into the actual machine code by the back-end. In our case, since we are using LLVM, we don't need to worry too much about the back-end, since LLVM will be in charge of generating the machine code. Our job will be to go from the source code to the LLVM intermediate representation (we will call it **IR** from now on). The front-end of the compiler is typically composed of three main steps.

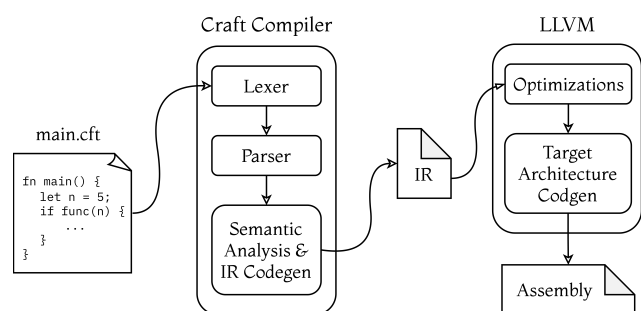


Fig. 1: Diagram showing the compilation workflow of a Craft program

8.1 Lexing

This step consists of breaking the source code into a sequence of tokens. A token is a basic building block of the languages, such as a keyword or an identifier.

A lexer can be implemented rather easily, by using a state machine. An approach to do it programmatically could be the following:

1. Start by reading the source code character by character until we reach the end.
2. If the character by itself forms a valid sequence (e.g. a parenthesis) we create a token from it. If it doesn't, we continue reading characters until we find a valid sequence.

Note that sometimes we may find a valid sequence, but that's not enough to create a token, since it may be the start of another longer and valid sequence. We may need to check the following character/s to check whether it continues. An example of such case would be the '>' operator, since, by itself is a valid sequence, but it may be the start of the '>=' operator. So we need to check if the next character is an '=' or something else.

The lexer requires a bit of work to set up, but after that, expanding it is trivial, since we only need to add a new word to the list of reserved words, in the case we want to add a reserved word, or add a new rule that detects a new symbol for example.

8.2 Parsing

The lexer allowed us to identify the symbols of the program, but it does not allow us to determine if their order is correct, or if they follow the rules of the language (i.e. it's grammatically correct). That is the job of the parser.

The parser takes the sequence of tokens obtained from the lexer and transforms it into an Abstract Syntax Tree (AST). This tree represents the structure of the program and determines its syntactic structure. It will tell us the order in which we need to execute the instruction.

To generate the AST compilers use a context free grammar. A grammar is a set of rules that tells us how to form valid strings of tokens in a specific language. They are formed of a set of symbols, which can be divided into terminal and non-terminal symbols, and a set of production rules that specify how the non-terminal symbols can be replaced by sequences of terminal and non-terminal symbols. A context free grammar is a type of grammar that, the rules of the grammar do not depend on the context in which the symbols appear.

The goal of the parser is to make the program obey the rules of the grammar.

Here's an example of a simple grammar for parsing function prototypes:

```
<proto> ::= fn <id> "(" <params> ")"
<id> ::= letter {letter | digit | "_"}
<params> ::= "(" { <param> {, <param> } } ")"
<param> ::= <id>: <type>
```

If we were to translate the proto rule to code, it would look something like this:

```
1 fn parse_prototype() -> (Prototype, Err) {
2   // we expect to find the fn keyword,
3   // else it's an error
4   match current_token().kind {
5     // The advance function moves to the next token
6     TokenKind::Fn => advance(),
7     _ => return Err("Expected fn keyword"),
8   };
9
10  // we expect to find the function name,
11  // else it's an error
12  let name = match current_token().kind {
13    TokenKind::Identifier => current_token().lexeme,
14    _ => return Err("Expected an identifier"),
15  };
16
17  advance();
18
19  // we call the params rule
20  let params = parse_params();
21
22  // we are done, we return a struct
23  // with the info of the prototype
24  return Prototype {
25    name,
26    params,
27  };
28 }
```

8.3 Semantic Analysis and [IR] Code Generation

Once we have an AST, we can traverse it to generate the LLVM IR. Since this is a simple compiler, we are going to do the semantic analysis in this step. With more complex compilers, we may want to create a specific step of semantic analysis, but in our case it is not necessary.

LLVM has the concept of modules. A module contains all the information associated with one code file. If we have multiple files, we simply have to create different modules and link them. Modules contain functions and functions are made up of instructions, similar to the instructions we find in assembly.

Let's see an example of a small piece of code and what it would look like in the intermediate representation of LLVM. We have the following function that receives two integers and returns the maximum:

```
1 fn max(int a, int b) int {
2   if a > b { a } else { b }
3 }
```

If we translate it to LLVM IR we have the following code:

```
1 define i32 @max(i32 %a, i32 %b) {
2   entry:
3     %0 = icmp sgt i32 %a, %b
4     br i1 %0, label %btrue, label %bfalse
5
6   btrue:
7     br label %end
8
9   bfalse:
10    br label %end
11
12   end:
13    %retval = phi i32 [%a, %btrue], [%b, %bfalse]
14    ret i32 %retval
15 }
```

The first line of the previous code defines a function, which receives two 32 bits integers. A label entry is defined below. These tags are like assembly tags, and we can jump to them.

The first thing it does, on the 'entry' tag, is comparing

both integers (i.e. the if condition). The 'sgt' keyword, means 'signed greater than', which, as the names says, does a greater than signed comparison. The result of the comparison is a one bit integer which acts like a boolean. If it's true, it will jump to the label btrue, otherwise to the label bfalse.

In this case the two branches do the same thing: a jump to the label end. There we come across a concept called phi nodes. The phi nodes are kind of an inverted if. Depending on where we did the jump we will assign one value or another to the variable retval. If we come from btrue, retval is assigned %a, and if we come from bfalse it is assigned to %b.

Now, why do the two branches jump to the end tag and then do a conditional again in the end block? Couldn't we assign the value of retval directly inside the branch btrue or bfalse and spare us that third conditional? Well the answer is no, because then we would be generating code that is not in SSA form. And as we mentioned earlier, LLVM requires the generated code to be in SSA form. And that's why phi nodes exist, to be able to solve this types of problems.

8.3.1 LLVM API

It's important to understand the LLVM IR, but generating all of that code by hand is a lot of work and it would require us to put a lot of infrastructure in place. Fortunately for us, LLVM comes with a C++ API, that makes generating the IR a lot easier. There's the little problem that we are not using C++, but thankfully, there's a Rust library called Inkwell, that offers a comfortable and idiomatic Rust interface to the LLVM API. Under the hood it calls the C++ API by using FFI (Foreign Function Interface).

The API is made of several entities. The most fundamental one is the BasicBlock, which represents the blocks that we find inside a function, i.e. the set of instructions we find within a label. If we group a bunch of BasicBlocks, we have an entity call Function, which, you guess it represents a function in the IR. And by grouping a bunch of functions we get the Module entity. As you may notice these entities represent the different parts that we saw when we talked about the IR.

Una altra classe important és Context. Necessitem una sola instància d'aquesta classe, la qual conte un conjunt d'estructures de dades de LLVM, i que contindrà l'estat de la compilació.

Finally we have the Builder entity, which is the one in charge of actually generating the code. For example, if we want to generate a comparison, like the one we saw on the IR example, we would call the Builder, telling it which instruction to generate:

```
Builder.CreateICmpSGT(a, b, "some-name")
```

The first two parameters are the two values being compared. The third indicates the name of the generated variable (we can give it any name we want, it's only useful for debugging).

The previous call will generate the following code (which is the third line of the IR example we saw):

```
%nom = icmp sgt i32 %a, %b
```

9 RESULTS

10 CONCLUSIONS

REFERENCES

- [1] <http://en.wikibooks.org/wiki/LaTeX>
- [2] Referència 2
- [3] Etc.

APÈNDIX

A.1 Secció d'Apèndix

.....

A.2 Secció d'Apèndix

.....