

Compiler Design — Lecture note week 1

- Author: Ruben Schenk
- Date: 23.09.2021
- Contact: ruben.schenk@inf.ethz.ch

1. Introduction

1.1 What is a Compiler?

A **compiler** is what we as developers usually see as a black box. We might use the C compiler `gcc` in the following way:

```
1  #include <stdio.h>
2
3  int main() {
4      printf("Hello world!\n");
5      return 0;
6  }
```

```
1  % gcc -o hello hello.c
2  % ./hello
3  Hello world!
4  %
```

The goal of a **compiler** is to *translate one programming language to another*, typically that is translating a high-level source code to a low-level machine code (**object code**).

Source Code

Source code is optimized for human readability. This means it is:

- Expressive: match human ideas of grammar/syntax/meaning
- Redundant: more information than needed to help catch errors
- Abstract: exact computations possibly not fully determined by code

Following an example of some C source code:

```
1  #include <stdio.h>
2
3  int factorial(int n) {
4      int acc = 1;
5      while(n > 0) {
6          acc = acc * n;
7          n = n - 1;
8      }
```

```
8     }
9     return acc;
10 }
11
12 int main(int argc, char *argv[]) {
13     printf("factorial(6) = %d\n", factorial(6));
14 }
```

Low-level code

Low-level code is optimized for hardware. This means that:

- Machine code is hard to read for humans
- Redundancy and ambiguity is reduced
- Abstractions and information about intent is lost

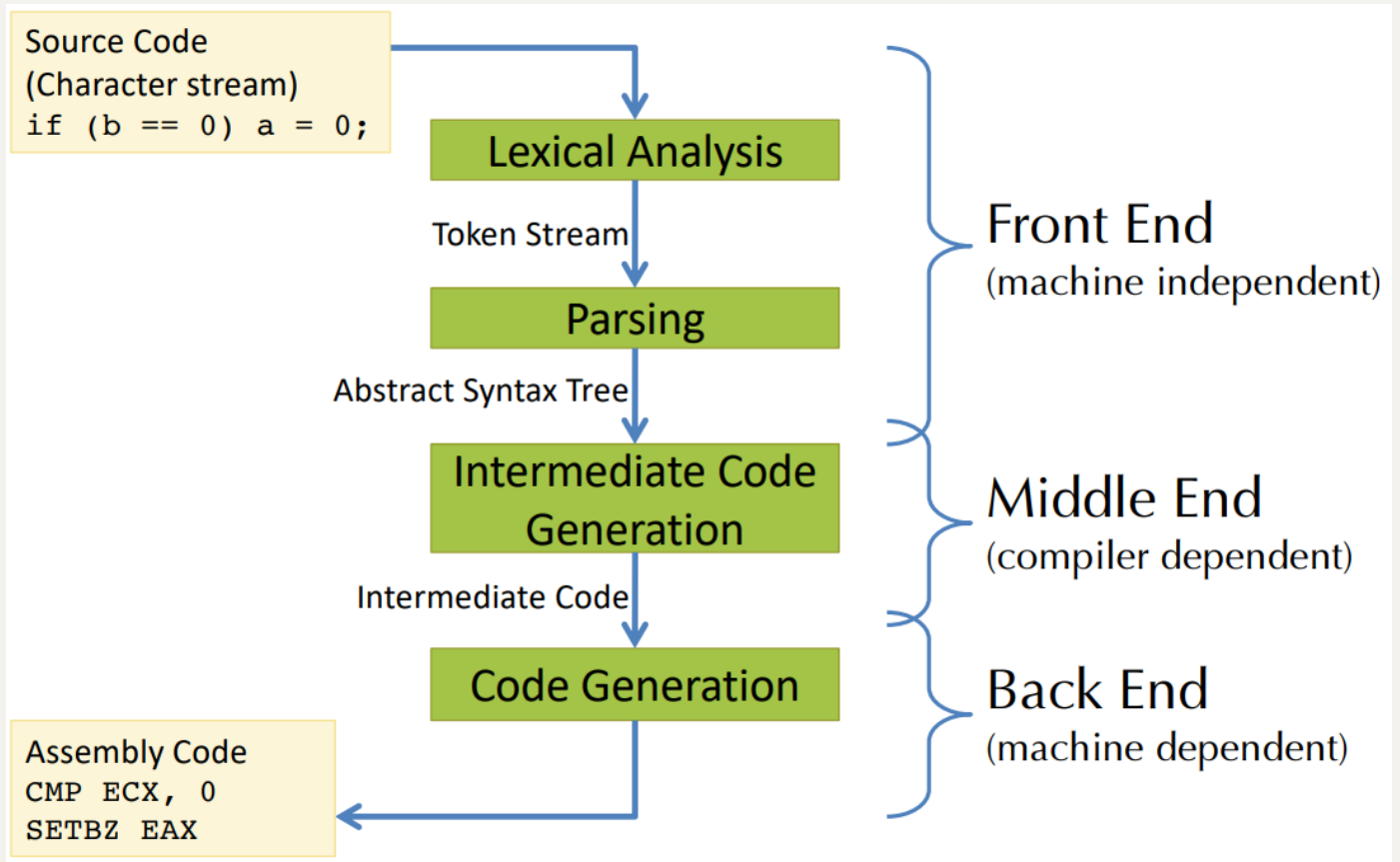
Compiler Bug Types

When compiling code we might encounter different bugs. We can distinguish them into different types:

- **Miscompilation (wrong code bug):** The compiler, maybe after enabling some optimization, gives us back a wrong code.
- **Internal compilation error (ICE):** The compiler crashes on trying to compile some source code.
- **Compiler hang (slow compilation):** The compiler is stuck or takes a very long time trying to compile some simple source code.
- **Missed optimizations:** The compiler might miss optimizations it could do to some given source code.

1.2 Compiler Structure

The following figure shows a simplified view of the **compiler structure**:

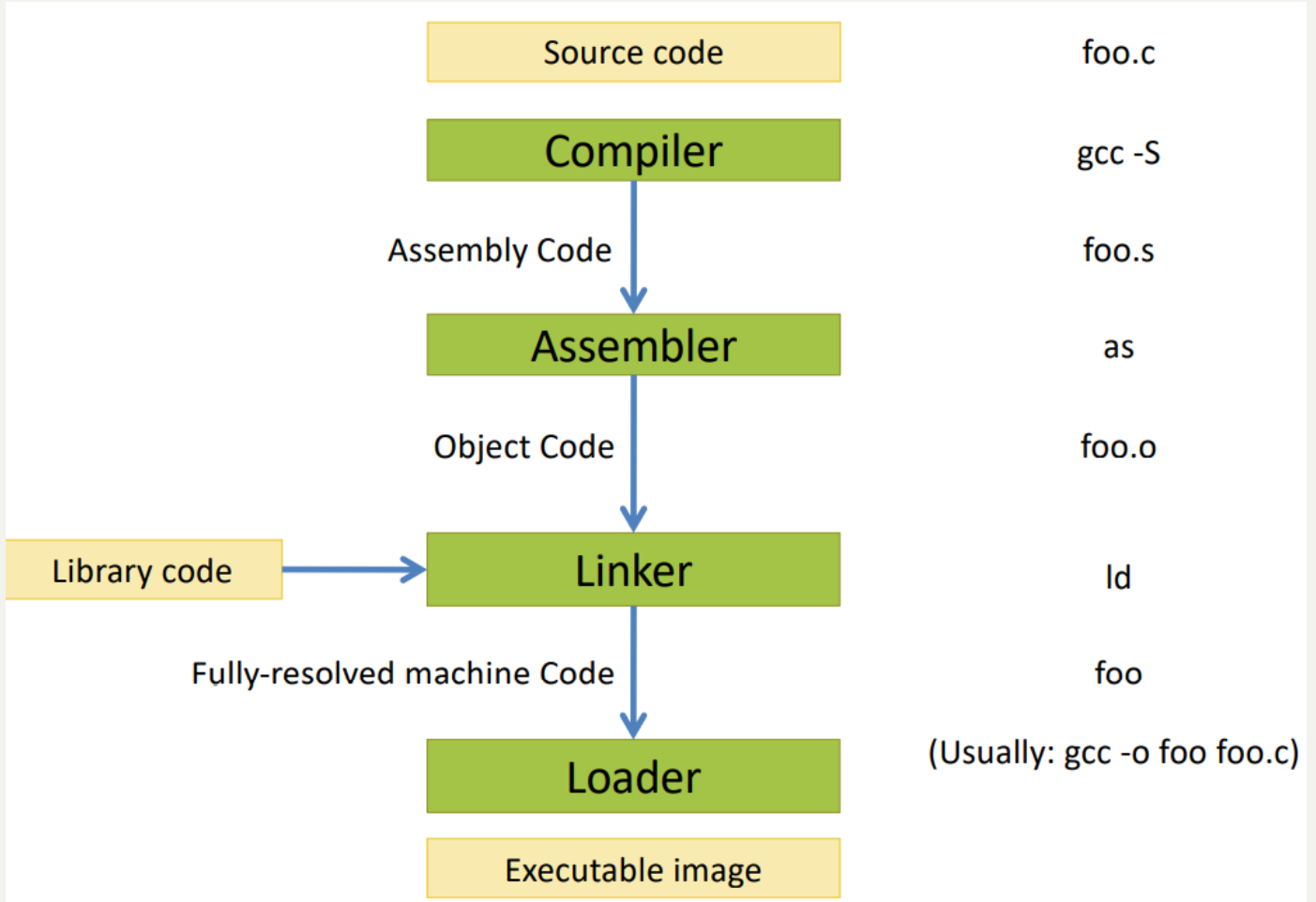


The typical **compiler stages** are as follows:

- Lexing -> token stream
- Parsing -> abstract syntax
- Disambiguation -> abstract syntax
- Semantic analysis -> annotated abstract syntax
- Translation -> intermediate code
- Control-flow analysis -> control-flow graph
- Data-flow analysis -> interference graph
- Register allocation -> assembly
- Code emission

Optimization may be done at *many* of these stages!

Another simplified view on the compilation and execution is given by the following figure:



2. OCaml

2.1 OCaml Tools

The programming language **OCaml** includes the following tools:

- `ocaml` -> The top-level interactive loop
- `ocamlc` -> The byte code compiler
- `ocamlopt` -> The native code compiler
- `ocamldep` -> The dependency analyzer
- `ocamldoc` -> The documentation generator
- `ocamllex` -> The lexer generator
- `ocamlyacc` -> The parser generator

In addition to the above-mentioned tools, one might use the following additional tools:

- `menhir` -> A more modern parser generator
- `ocamlbuild` -> A compilation manager
- `utop` -> A more fully-featured interactive top-level
- `opam` -> Package manager

2.2 OCaml Characteristics

OCaml has the following two main distinguishing characteristics:

Functional and mostly pure

- Programs manipulate values rather than issue commands
- Functions are first-class entities
- Results of computations can be "named" using `let`
- Has relatively few "side effects"

Strongly and statically typed

- Compiler typechecks every expression of the program, issues errors if it can't prove that the program is type safe
- Good support for type inference and generic polymorphic types
- Rich user-defined "algebraic data types" with pervasive use of pattern matching
- Very strong and flexible module system for constructing large projects

2.3 Factorial on OCaml

Consider the following implementation of the factorial function in a hypothetical programming language:

```
1  x = 6;  
2  ANS = 1;  
3  whileNZ (x) {  
4      ANS = ANS * x;  
5      x = x + -1;  
6  }
```

For this hypothetical language, we need to describe the following two constructs:

- **Syntax:** which sequences of characters count as a legal program?
- **Semantics:** what is the meaning of a legal program?

2.4 Grammar for a Simple Language

2.4.1 Grammar and Interpreter

We introduce the following two **nonterminals** for our simple language:

```
1  <exp> ::=  
2      | <X>  
3      | <exp> + <exp>  
4      | <exp> * <exp>  
5      | <exp> < <exp>  
6      | <integer constant>  
7      | (<exp>)  
8  
9  <cmd> ::=  
10     | skip  
11     | <X> = <exp>  
12     | ifNZ <exp> { <cmd> } else { <cmd> }  
13     | whileNZ <exp> { <cmd> }  
14     | <cmd>; <cmd>
```

The above given syntax (or *grammar*) for a simple imperative language has the following properties:

- It is written in *Backus-Naur form*
- The symbols `::=`, `|`, and `<...>` are part of the **meta language**

- Keywords like `skip`, `ifNZ`, and `whileNZ` and symbols like `{` and `+` are part of the **object language**

Example: Define Grammar in OCaml

With the above definition of our grammar in BNF, we can transform this into OCaml. It looks the following way:

```

1  type var = string;
2
3  type exp =
4    | Var of var
5    | Add of (exp * exp)
6    | Mul of (exp * exp)
7    | Lt  of (exp * exp)
8    | Lit of int
9
10 type cmd =
11   | Skip
12   | Assn  of var * exp
13   | IfNZ  of exp * cmd * cmd
14   | WhileNZ of exp * cmd
15   | Seq   of cmd * cmd

```

With the definition of our hypothetical language, we can build a command for the factorial function in OCaml the following way:

```

1  let factorial : cmd =
2    let x = "x" in
3    let ans = "ANS" in
4    Seq (Assn (x, Lit 6)),
5    Seq (Assn (ans, Lit 1)),
6    WhileNZ(Var x,
7      Seq (Assn (ans, Mul (Var and, Var x)),
8        Assn (x, Add (Var x, Lit (-1))))))

```

With the above two examples we can now finally build a simple **interpreter** for our simple language:

```

1  type state = var -> int
2
3  let rec interpret_exp (s:state) (e:exp) : int =
4    match e with

```



```

5   | Var x -> s x
6   | Add (e1, e2) -> (interpret_exp s e1) + (interpret_exp s e2)
7   | Mul (e1, e2) -> (interpret_exp s e1) * (interpret_exp s e2)
8   | Lt  (e1, e2) -> if (interpret_exp s e1) < (interpret_exp s e2) then 1
    else 0
9   | Lit n -> n
10
11 let update s x v =
12     fun y -> if x = y then v else s y
13
14 let rec interpret_cmd (s:state) (c:cmd) : state =
15     match c with
16     | Skip -> s
17     | Assn (x, e1) ->
18         let v = interpret_exp s e1 in
19         update s x v
20     | IfNZ (e1, c1, c2) ->
21         if (interpret_exp s e1) = 0 then interpret_cmd s c2 else
22         interpret_cmd s c1
23     | WhileNZ (e, c) ->
24         if (interpret_exp s e) = 0 then s else interpret_cmd s (Seq (c,
25         WhileNZ (e, c)))
26     | Seq (c1, c2) ->

```

Main Function

We might write and invoke a `main` function in the following way:

```

1  let main() =
2      Printf.printf("Hello world!")
3
4  (* let _ = main () *)
5  ;; main ()

```

2.4.2 Optimizer

We might `optimize` our interpreter. This could be that we evaluate simple expressions ourselves, instead of letting the compiler evaluate it completely. Examples:

```

1  (*
2  e + 0 -> e
3  e * 1 -> e
4  e * 0 -> 0

```

```

5  0 + e -> e
6  e - e -> 0
7  ...
8
9  skip; c -> c
10
11 ifNZ 0 then c1 else c2 -> c2
12 ifNZ 1 then c1 else c2 -> c1
13
14 whileNZ 0 c -> skip
15 *)

```

In general, we want to make our program as simple as possible based on some rewriting rules before interpreting it.

We might realize an *optimizer for commands* in the following way:

```

1  let rec optimize_cmd (c:cmd) : cmd =
2      match c with
3      | Assn(x, Var y) -> if x = y then Skip else c
4      | Assn(_, _) -> c
5      | WhileNZ(Lit 0, c) -> Skip
6      | WhileNZ(Lit _, c) -> loop
7      | WhileNZ(e, c) -> WhileNZ(e, optimize_cmd c)
8      | Skip -> Skip
9      | IfNZ(Lit 0, c1, c2) -> optimize_cmd c2
10     | IfNZ(Lit _, c1, c2) -> optimize_cmd c1
11     | IfNZ(e, c1, c2) -> IfNZ(e, optimize_cmd c1, optimize_cmd c2)
12     | Seq(c1, c2) ->
13         begin match (optimize_cmd c1, optimize_cmd c2) with
14         | (Skip, c2') -> c2'
15         | (c1', Skip) -> c1'
16         | (c1', c2') -> Seq(c1', c2')
17     end

```

2.4.3 Translator

We might imagine trying to build a translator from *Simple* to *OCaml*. This process consists of several different steps.

Set of Variables

In the following code we explore how to get the set of variables from a given expression.

```
1  ;; open Simple
2
3  module OrderedVars = struct
4      type t = var
5      let compare = String.compare
6  end
7
8  module VSet = Set.Make(OrderedVars)
9  let (++) = VSet.union
10
11  (* Calculate the set of variables mentioned in either an expression or a
12     command *)
13
14  let rec vars_of_exp (e:exp) : VSet.t =
15      begin match e with
16      | Var x -> VSet.singleton x
17      | Add(e1, e2)
18      | Mul(e1, e2)
19      | Lt (e1, e2) ->
20          (vars_of_exp e1) ++ (vars_of_exp e2)
21      | Lit _ -> VSet.empty
22      end
23
24  let rec vars_of_cmd (c:cmd) : VSet.t =
25      begin match c with
26      | Skip -> VSet.empty
27      | Assn(x, e) -> (VSet.singleton x) ++ (vars_of_exp e)
28      | IfNZ(e, c1, c2) ->
29          (vars_of_exp e) ++ (vars_of_cmd c1) ++ (vars_of_cmd c2)
30      | WhileNZ(e, c) ->
31          (vars_of_exp e) ++ (vars_of_cmd c)
32      | Seq(c1, c2) ->
33          (vars_of_cmd c1) ++ (vars_of_cmd c2)
34      end
```

Translation

The translation invariants are guided by the *types* of the operations:

- variables are a global state, so they become mutable references
- expressions denote integers
- commands denote imperative actions of type unit

We might build our translator the following way:

```
1  let trans_var (x:var) : string =
2    "V_" ^ x
3
4  let rec trans_exp (e:exp) : string =
5    let trans_op (e1:exp) (e2:exp) (op:string) =
6      Printf.sprintf "(%s %s %s)"
7        (trans_exp e1) op (trans_exp e2)
8    in
9    begin match e with
10     | Var x -> "!" ^ (trans_var x)
11     | Add(e1, e2) -> trans_op e1 e2 "+"
12     | Mul(e1, e2) -> trans_op e1 e2 "*"
13     | Lt (e1, e2) ->
14       Printf.sprintf "(if %s then 1 else 0)"
15         (trans_op e1 e2 "<")
16     | Lit 1 -> string_of_int 1
17   end
18
19  let rec trans_cmd (c:cmd) : string =
20    begin match c with
21     | Skip -> "()"
22     | Ass(x, e) ->
23       Printf.sprintf "%s := %s"
24         (trans_var x) (trans_exp e)
25     | IfNZ(e, c1, c2) ->
26       Printf.sprintf "if %s <> 0 then (%s) else (%s)"
27         (trans_exp e) (trans_cmd c1) (trans_cmd c2)
28     | WhileNZ(e, c) ->
29       Printf.sprintf "while %s <> 0 do \n %s done"
30         (trans_exp e) (trans_cmd c)
31     | Seq(c1, c2) ->
32       Printf.sprintf "%S; \n %s"
33         (trans_cmd c1) (trans_cmd c2)
34   end
```

```
35
36 let trans_prog (c:cmd) : string =
37   let vars = vars_of_cmd c in
38   let decls =
39     VSet.fold (fun x s ->
40       Printf.sprintf " let %s = ref 0 \n %s \n"
41         (trans_var x) d)
42     vars
43     ""
44   in
45     Printf.sprintf "module Program = struct \n %s let run () = \n %s \n
end"
46     decls (trans_cmd c)
47
48 (* Do some testing using the factorial code: Simple.factorial *)
49 let _ =
50   Printf.printf ("%s \n") (trans_prog factorial)
```