



IMT Atlantique
Bretagne-Pays de la Loire
École Mines-Télécom



Compilation

A crash course

Fabien Dagnat
LaLog – Spring 2026



Plan

1 The structure of a compiler

2 Lexing

3 Parsing

4 Core

1 The structure of a compiler

2 Lexing

3 Parsing

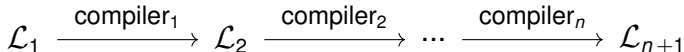
4 Core

What is a compiler?

- ▶ It is a program transformer from \mathcal{L}^I to \mathcal{L}^O

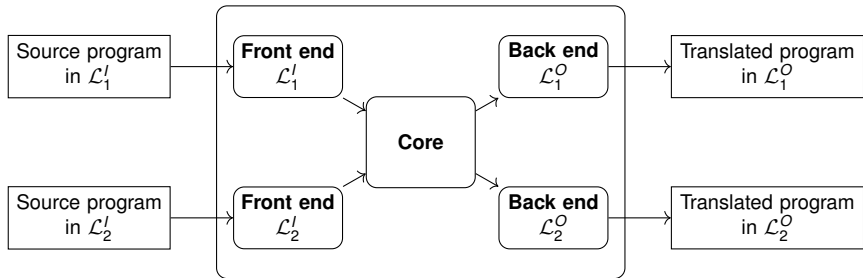


- ▶ \mathcal{L}^O is generally "more" executable than \mathcal{L}^I
- ▶ Most of the time, a program in \mathcal{L}^O is directly executable
 - ▶ either by a machine (e.g. $\mathcal{L}^O = \text{X64}$)
 - ▶ or by an **abstract machine**, a piece of software acting as a machine (e.g. $\mathcal{L}^O = \text{OCaml bytecode}$)
- ▶ Often a compiler is composed of a flow of simpler compilers



Structure of a compiler

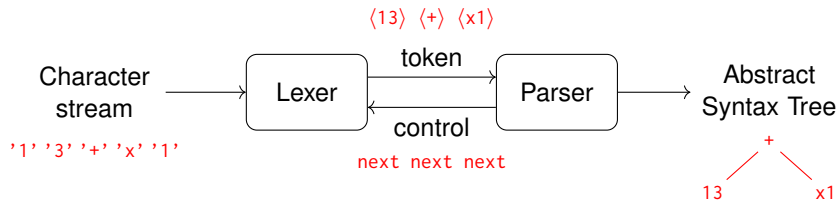
- ▶ It is composed of three stages
 - ▶ a *front end* to recognize \mathcal{L}^I (e.g. LLVM has C, C++, Go)
 - ▶ a *core* doing the hard work
 - ▶ a *back end* in charge of emitting \mathcal{L}^O (e.g. LLVM has X64, ARM)



- ▶ Several paths are possible

Front end

- ▶ It checks whether the program is syntactically correct
 - ▶ The program belongs to the language \mathcal{L}'
 - ▶ It must build an internal representation (IR) of the program
 - ▶ The IR is an internal data structure for the compiler
- ⇒ The front end is highly dependent on the input language
- ▶ It is decomposed in two parts
 - ▶ **lexer** recognizes **tokens** in a character stream, returns a token
 - ▶ **parser** recognizes **sentences** in a token stream, returns an AST



Core

- ▶ Works on internal data structures
- ▶ In charge of
 - ▶ the verification of validity (**typing**) of the program
 - ▶ the main transformation work, for instance
 - ▶ simplify programs by removing useless elements
 - ▶ transform function calls
 - ▶ transform object oriented access
 - ▶ ...
- ▶ Relatively usual software
- ▶ Functional paradigm is adapted for this kind of code
 - ▶ recursive functions for the visiting part
 - ▶ sum types for representing the various elements
 - ▶ product types to add information to the various elements
 - ▶ pattern matching for the recognition of structure
 - ▶ ...

Back end

- ▶ Translates internal data structures in instructions of the target language

...

```
movq $1, %rax
```

```
; some code to get the value of x1
```

```
addq $26, %rax
```

...

- ▶ Implements all optimizations specific to the target
- ▶ Complex requiring to master the target machine
- ▶ Not in the scope of this introduction...

Progress

1 The structure of a compiler

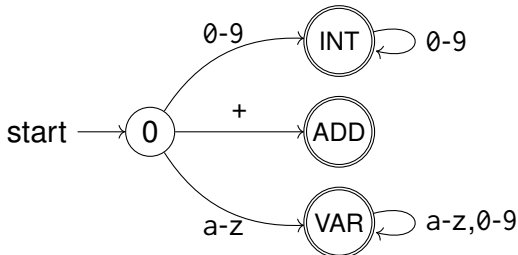
2 Lexing

3 Parsing

4 Core

Lexical analyzer (*a.k.a* lexer)

- ▶ A lexer reads enough characters from an entry stream to produce a token
- ▶ To be efficient it is generally built as an automaton where
 - ▶ transitions correspond to the received characters
 - ▶ final states correspond to the production of the recognized token



- ▶ When reaching a final state, it produces a token (data structure for the parser)

Building a lexer

- ▶ One way to define such an automaton is to define a set of **translation rules**
- ▶ A translation rule is composed of a **pattern** and an **action**
 - ▶ the pattern is a regular expression specifying the accepted input
 - ▶ the action defines what to do when accepting (often just returning the right token)
- ▶ For the previous slide example

```
| ['0'-'9']+ { INT((* input converted in int *)) }  
| '+' { ADD }  
| ['a'-'z']['a'-'z''0'-'9']* { VAR((* input *)) }
```

- ▶ A Domain Specific Language: OCamllex
 - ▶ a compiler `ocamllex` producing OCaml code for the automaton

OCamllex syntax

► File with extension `.ml1`

`{ (* OCaml code: optional prelude *) }` as is in the result

(* useful regular expressions only for regexp part *)

let `ident` = `regexp`

let `ident` = `regexp`

(* a group of rules *)

rule `ident` [`ident1` ... `identn`] = **parse**

| `regexp` { (* OCaml code *) }

| `regexp` { (* OCaml code *) }

(* another group of rules *)

and `ident` [`ident1` ... `identn`] = **parse**

...

`{ (* OCaml code: optional postlude *) }` as is in the result

compiled as an automaton

actions executed

on accepting

translated as a function

OCamllex regexp syntax by example

▶ `[' ' '\014' '\t' '\012']+` ⇒ at least one space

OCamllex regexp syntax by example

- ▶ `[' ' '\014' '\t' '\012']+` \Rightarrow at least one space
- ▶ `(['\n' '\r'] | "\r\n")` \Rightarrow newline

OCamllex regexp syntax by example

- ▶ `[' '\014' '\t' '\012']+` \Rightarrow at least one space
- ▶ `(['\n' '\r'] | "\r\n")` \Rightarrow newline
- ▶ `[^ '\n' '\r']` \Rightarrow any character except newline

OCamllex regexp syntax by example

- ▶ `[' ' '\014' '\t' '\012']+` ⇒ at least one space
- ▶ `(['\n' '\r'] | "\r\n")` ⇒ newline
- ▶ `[^ '\n' '\r']` ⇒ any character except newline
- ▶ `"//[^\n\r]*"`

OCamllex regexp syntax by example

- ▶ `[' ' '\014' '\t' '\012']+` ⇒ at least one space
- ▶ `(['\n' '\r'] | "\r\n")` ⇒ newline
- ▶ `[^ '\n' '\r']` ⇒ any character except newline
- ▶ `"//[^\n\r]*"` ⇒ C like line comment
- ▶ Suppose

```
let digit = ['0'-'9']
let letter = ['a'-'z' 'A'-'Z']
let id_char = (letter | digit | '_')
```
- ▶ `letter id_char*` as `id`

OCamllex regexp syntax by example

- ▶ `[' ' '\014' '\t' '\012']+` ⇒ at least one space
- ▶ `(['\n' '\r'] | "\r\n")` ⇒ newline
- ▶ `[^ '\n' '\r']` ⇒ any character except newline
- ▶ `"//[^\n\r]*"` ⇒ C like line comment
- ▶ Suppose
 - `let digit = ['0'-'9']`
 - `let letter = ['a'-'z' 'A'-'Z']`
 - `let id_char = (letter | digit | '_')`
- ▶ `letter id_char* as id` ⇒ identifiers, `id` contains the result
- ▶ integers

OCamllex regexp syntax by example

- ▶ `[' ' '\014' '\t' '\012']+` ⇒ at least one space
- ▶ `(['\n' '\r'] | "\r\n")` ⇒ newline
- ▶ `[^ '\n' '\r']` ⇒ any character except newline
- ▶ `"//[^\n\r]*"` ⇒ C like line comment
- ▶ Suppose
 - `let digit = ['0'-'9']`
 - `let letter = ['a'-'z' 'A'-'Z']`
 - `let id_char = (letter | digit | '_')`
- ▶ `letter id_char* as id` ⇒ identifiers, `id` contains the result
- ▶ integers ⇒ `digit+` as `nb`
- ▶ floating point numbers

OCamllex regexp syntax by example

- ▶ `[' '\014' '\t' '\012']+` ⇒ at least one space
- ▶ `(['\n' '\r'] | "\r\n")` ⇒ newline
- ▶ `[^ '\n' '\r']` ⇒ any character except newline
- ▶ `"//[^\n\r"]*` ⇒ C like line comment
- ▶ Suppose
 - `let digit = ['0'-'9']`
 - `let letter = ['a'-'z' 'A'-'Z']`
 - `let id_char = (letter | digit | '_')`
- ▶ `letter id_char* as id` ⇒ identifiers, `id` contains the result
- ▶ integers ⇒ `digit+ as nb`
- ▶ floating point numbers ⇒
`digit* '.' digit* (['e' 'E'] ['+' '-']? digit+)? as nb`

Generated OCaml code

- ▶ Each rule `<name> a1 ... an` gives a function name taking
 - ▶ `a1, ... an`, the user arguments
 - ▶ a buffer for the stream of characters of type `Lexing.lexbuf`
- ▶ This function matches the characters in the buffer and execute the corresponding accepting action
 - ▶ it selects the regexp giving the longest part matched
 - ▶ in case of equality it selects the first defined
- ▶ The standard library module `Lexing` also provides
 - ▶ two constructors for buffers: `from_channel` and `from_string`
 - ▶ `lexeme buf` returning the matched string of `buf`
- ▶ Only one automaton is generated
 - ▶ the automaton is determinized and minimized
 - ▶ its code is finally put between the prelude and postlude

An example and the flow

- ▶ File `formulaLexer.mll`

```
{
  type token = EOF | AND | OR | TRUE | FALSE
}
let space = [' ' '\t' '\n']
rule token = parse
  | space+ { token lexbuf }
  | eof    { EOF }
  | "and"  { AND }
  | "or"   { OR }
  | "true" { TRUE }
  | "false" { FALSE }
```

- ▶ `ocamllex formulaLexer.mll` produces `formulaLexer.ml`
- ▶ It then can be compiled using `ocamlc`
 - ⚠ it can contain errors if the `mll` file contained wrong OCaml code

Progress

1 The structure of a compiler

2 Lexing

3 Parsing

4 Core

Which kind of grammar?

- ▶ For most reasonable syntax, regexp are not sufficient
 - ▶ We must use more powerful grammars but keep efficiency of parsing
- ⇒ We use context-free grammars (CFG)
- ▶ defined only by production $A \rightarrow m$ where $A \in V$ and $m \in (X \cup V)^*$
 - ▶ In this course, we focus on LR(1) parsing by using Menhir
 - ▶ Menhir `http://gallium.inria.fr/~fpottier/menhir`
 - ▶ offers a DSL for defining grammars in `.mly` files
 - ▶ has a tool compiling a grammar spec. to OCaml code (`menhir`)
 - ▶ Menhir follows a flow similar to OCamllex

Some remarks

- ▶ The token type is now generated within the parser
- ⇒ The lexer does not define it anymore but imports it
- ▶ Each entry point (`%start<type>`) gives a parsing function of type `(Lexing.lexbuf -> token) -> Lexing.lexbuf -> type`
- ⇒ The lexer must be given to the parser

```
let compile file =
  try
    let input_file = open_in file in
    let result = formula token (Lexing.from_channel input_file) in
    close_in (input_file);
    printf "read %s\n" result
  with Sys_error s ->
    printf "Can't find file '%s'" file
let () = Arg.parse [] compile ""
```

The result of parsing

- ▶ In general, it is a data structure representing programs called an **Abstract Syntax Tree** (AST)
- ▶ Defined using recursive sum types

```
type t = Var of string
      | Bool of bool
      | And of t * t
      | Or of t * t
```

- ▶ Manipulated by recursive functions

```
let rec string_of = function
  | Var s -> s
  | Bool true -> "true"
  | Bool false -> "false"
  | And(f1, f2) -> "(" ^ (string_of f1) ^ " and " ^ (string_of f2) ^ ")"
  | Or(f1, f2) -> "(" ^ (string_of f1) ^ " or " ^ (string_of f2) ^ ")"
```

The example revisited

```
%{  
  open FormulaAst  
}%
```

Prelude

```
%token AND OR EOF TRUE FALSE  
%token <string> IDENT  
%start< FormulaAst.t > formula  
%%
```

```
formula: c=conj EOF
```

```
{ c }
```

```
conj:
```

```
| d=disj AND c=conj
```

```
{ And(d,c) }
```

```
| d=disj
```

```
{ d }
```

```
disj:
```

```
| s=ident_or_const OR d=disj
```

```
{ Or(s,d) }
```

```
| s=ident_or_const
```

```
{ s }
```

```
ident_or_const:
```

```
| id=IDENT
```

```
{ Var id }
```

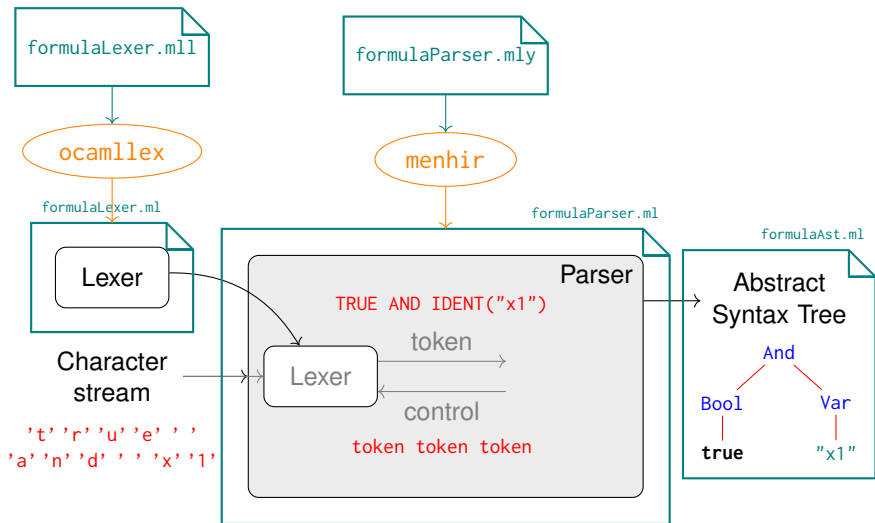
```
| TRUE
```

```
{ Bool true }
```

```
| FALSE
```

```
{ Bool false }
```

The new complete flow



Progress

1 The structure of a compiler

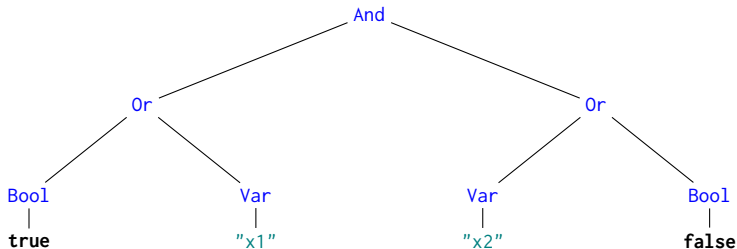
2 Lexing

3 Parsing

4 Core

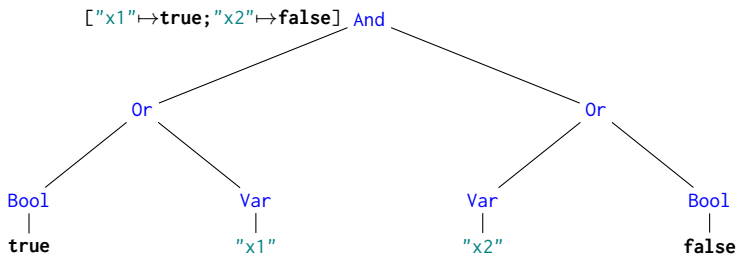
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



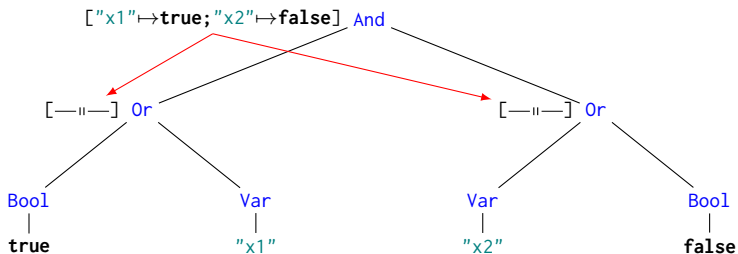
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



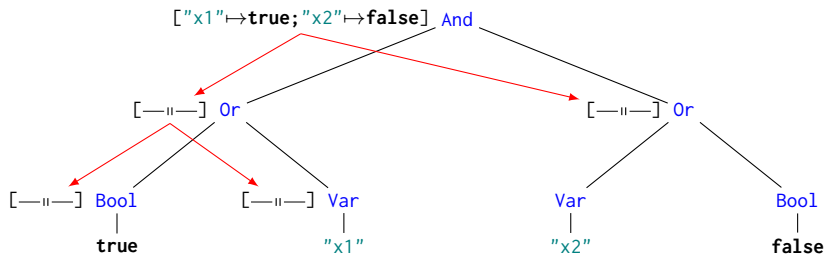
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



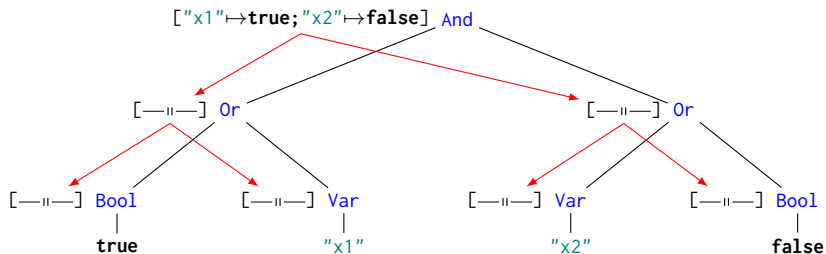
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



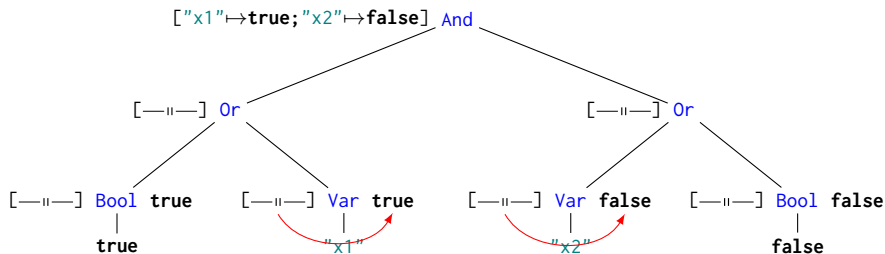
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



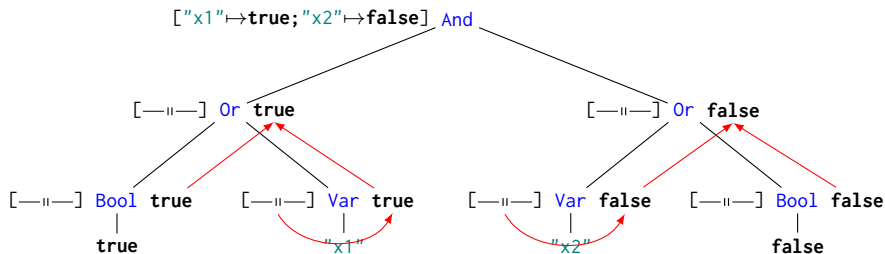
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



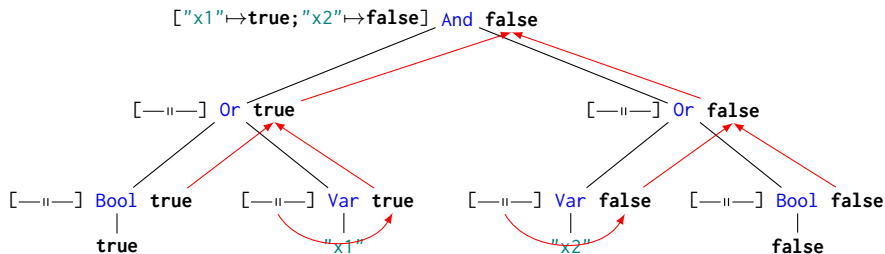
Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



Evaluation

- ▶ One way to execute a program is to use an **interpreter**
 - ▶ often called a Read Eval Print Loop
- ▶ Consists in producing a **value** from an AST
- ▶ It uses a recursive visit of the AST to synthesize the value
- ▶ While visiting the AST, naming information is collected



Formalization

- ▶ Such visit can be formalized using big step Structural Operational Semantics
 - ▶ values are from the boolean algebra $\mathbb{B} = \{\mathbf{T}, \mathbf{F}\}$ with \wedge and \vee
 - ▶ $\mathbf{T} \wedge \mathbf{T} = \mathbf{T}, \mathbf{F} \wedge b = b \wedge \mathbf{F} = \mathbf{F}, \mathbf{T} \vee b = b \vee \mathbf{T} = \mathbf{T}, \mathbf{F} \vee \mathbf{F} = \mathbf{F}$
 - ▶ an environment function \mathcal{E} mapping variable names to values
 - ▶ *dom* gives its domain, $\mathcal{E}(x)$ gives the value associated to x in \mathcal{E}
 - ▶ judgements of the form $\mathcal{E} \vdash \text{AST term} \Rightarrow \text{value}$

$$(1) \mathcal{E} \vdash \text{Bool true} \Rightarrow \mathbf{T} \quad (2) \mathcal{E} \vdash \text{Bool false} \Rightarrow \mathbf{F}$$

$$(3) \frac{x \in \text{dom}(\mathcal{E})}{\mathcal{E} \vdash \text{Var } x \Rightarrow \mathcal{E}(x)} \quad (4) \frac{\mathcal{E} \vdash F_1 \Rightarrow b_1 \quad \mathcal{E} \vdash F_2 \Rightarrow b_2}{\mathcal{E} \vdash \text{And}(F_1, F_2) \Rightarrow b_1 \wedge b_2}$$

$$(5) \frac{\mathcal{E} \vdash F_1 \Rightarrow b_1 \quad \mathcal{E} \vdash F_2 \Rightarrow b_2}{\mathcal{E} \vdash \text{Or}(F_1, F_2) \Rightarrow b_1 \vee b_2}$$

Implementation

```
open FormulaAst
let eval env formula =
  let rec eval_rec = function
    | Var s    -> List.assoc s env
    | Bool true -> true
    | Bool false -> false
    | And(f1, f2) -> (eval_rec f1) && (eval_rec f2)
    | Or(f1, f2) -> (eval_rec f1) || (eval_rec f2)
  in
  eval_rec formula
```

- ▶ Here, as the environment is constant during the visit, it is made global to the visiting function (`eval_rec`)
- ▶ Implementation should take care of errors ($x \notin \text{dom}(\mathcal{E})$)

Type

- ▶ Evaluation can lead to runtime errors
- ▶ Typing "approximates" evaluation to detect errors in advance
 - ▶ the "value" set is simplified and called type set
 - ▶ a new "value" is created to represent errors \perp
 - ▶ operations are defined on this simplified set
- ▶ For our example
 - ▶ all booleans values are approximated by the type *bool*
 - ▶ \wedge and \vee are both transformed in \otimes
 - ▶ $bool \otimes bool = bool$ and $\perp \otimes x = x \otimes \perp = \perp$
 - ▶ the environment is approximated by a type environment Γ

$$(1) \Gamma \vdash \text{Bool } b : bool$$

$$(2) \Gamma \vdash \text{Var } x : \Gamma(x)$$

$$(3) \frac{\Gamma \vdash F_1 : b_1 \quad \Gamma \vdash F_2 : b_2}{\Gamma \vdash \text{And}(F_1, F_2) : b_1 \otimes b_2}$$

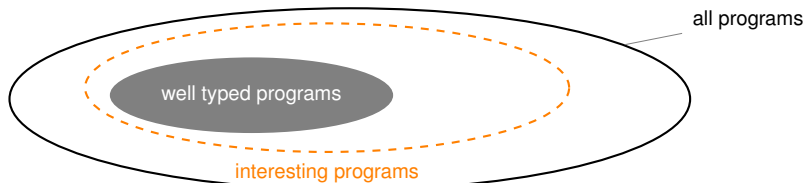
$$(4) \frac{\Gamma \vdash F_1 : b_1 \quad \Gamma \vdash F_2 : b_2}{\Gamma \vdash \text{Or}(F_1, F_2) : b_1 \otimes b_2}$$

Usefulness of types

Subject reduction theorem (safety)

$$\emptyset \vdash P : \tau \wedge \tau \neq \perp \implies (\emptyset \vdash P \Rightarrow v \wedge \emptyset \vdash v : \tau) \vee \emptyset \vdash P \implies \infty$$

- ▶ *Well-typed programs cannot "go wrong"* (produce errors)
- ▶ Computing types is much cheaper than evaluating
- ⚠ In general, Halting and Error discovery are **undecidable**
- ▶ Any typing must reject correct (too complicated) programs

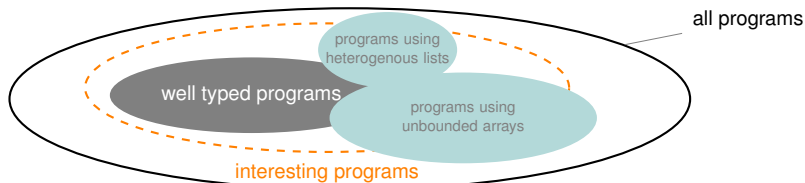


Usefulness of types

Subject reduction theorem (safety)

$$\emptyset \vdash P : \tau \wedge \tau \neq \perp \implies (\emptyset \vdash P \Rightarrow v \wedge \emptyset \vdash v : \tau) \vee \emptyset \vdash P \implies \infty$$

- ▶ *Well-typed programs cannot "go wrong"* (produce errors)
- ▶ Computing types is much cheaper than evaluating
- ⚠ In general, Halting and Error discovery are **undecidable**
- ▶ Any typing must reject correct (too complicated) programs
 - ⇒ needs a compromise between flexibility and safety
 - ⇒ to achieve safety, runtime checks are often needed



More on types

- ▶ Type algebra are often much more complex
 - ▶ more values
 - ▶ new types during typing...
 - ▶ specific relation between types (*e.g.* subtyping)
- ▶ Often the developer must add type annotations
 - ▶ **type checking**: given Γ , P and τ , is $\Gamma \vdash P : \tau$ true?
 - ▶ the simply typed λ -calculus
 - ▶ syntax: $M ::= \dots \mid \lambda x : \tau. M$ and types $\tau ::= \tau \rightarrow \tau$
 - ▶ typing

$$\frac{x : \tau \in \Gamma}{\Gamma \vdash x : \tau} \quad \frac{\Gamma, x : \tau \vdash M : \tau'}{\Gamma \vdash \lambda x : \tau. M : \tau \rightarrow \tau'} \quad \frac{\Gamma \vdash M_1 : \tau \rightarrow \tau' \quad \Gamma \vdash M_2 : \tau}{\Gamma \vdash M_1 M_2 : \tau'}$$

- ▶ With no annotation, it is **Type Inference** (*i.e.* OCaml)
 - ▶ **typeability**: given P finds Γ and τ such that $\Gamma \vdash P : \tau$ is true
 - ▶ much harder

Optimization

- ▶ The step using the most difficult algorithms and heuristics
- ▶ It aims at transforming AST to reduce some consumption
 - ▶ time, memory, energy, ...
- ▶ Can be generic or specific to a target machine
- ▶ For example
 - ▶ for all b , $\mathbf{F} \wedge b = \mathbf{F}$
 - ▶ so `And(Bool false, F)` can be simplified to `Bool false`
- ▶ In real life much more complex!
 - ▶ loop unrolling
 - ▶ propagating constants, inlining small functions
 - ▶ removing dead code
 - ▶ transforming variables into `StaticSingleAssignment`
 - ▶ tail-call, closure elimination
 - ▶ ...

Compilation

- ▶ Transform each element of the AST to machine operation
- ▶ See the practical session

Conclusion

- ▶ Just a very fast introduction to compilation
- ▶ Practice will help concretize!
- ▶ Formalization is important and often forgotten by engineers, that's an error!
- ▶ Vocabulary
 - ▶ abstract machine, token, sentence, typing, translation rule, pattern, action, abstract syntax tree, interpreter, value, undecidable, type checking, type inference, typeability,
- ▶ Acronym
 - ▶ LR1, AST, REPL, SOS, SSA