



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

PC3, TP7-8 – Let's practice compilation

Langages et logique – ELU 610

Objectives

At the end of the activity, you should be capable of:

- to define the syntax and implement the corresponding lexer and parser for a simple language,
- give the formal semantics of a simple language,
- implements a simple interpreter of it,
- implements a simple translation from one language to another one.

Part I

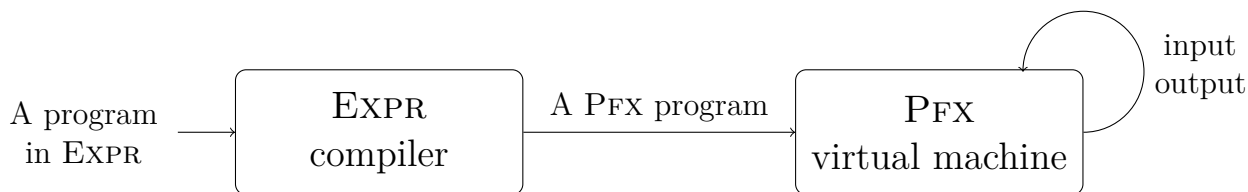
Context

1 The big picture

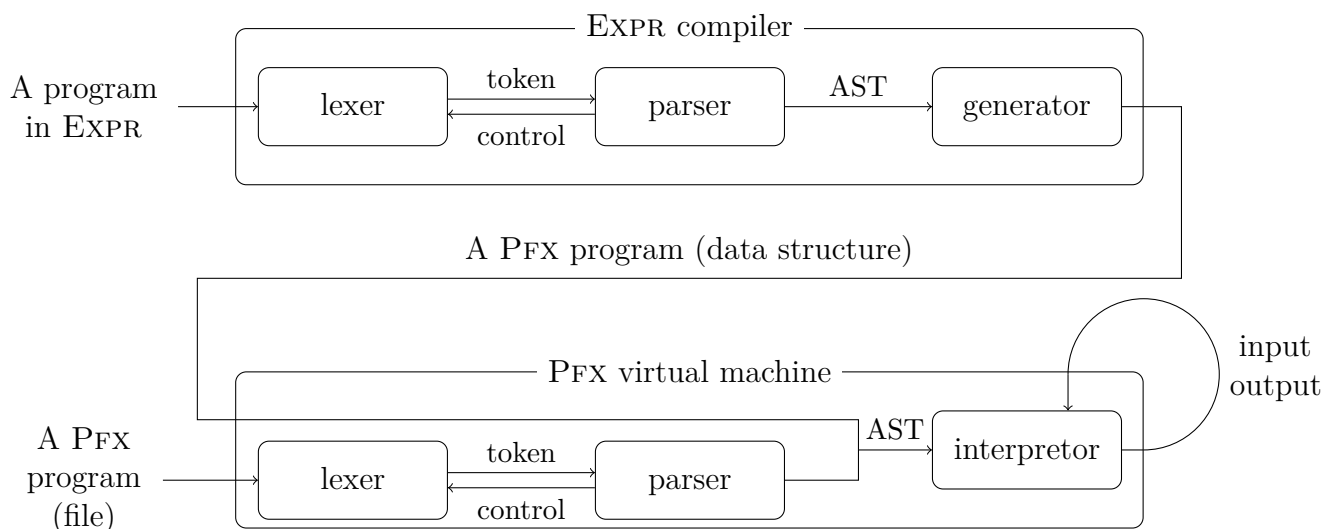
During our practice of compilation, we are going to write a compiler for a simple arithmetic expression language. This simple language is named `EXPR` and follow the exercise 6 directly from the *Discovering OCaml* document.

To make it simpler for you, instead of generating X86 code to execute it on a standard machine, we are going to use a *virtual machine*. Such a virtual machine is a program that takes a low-level program and executes it. We talk about interpretation. For pedagogical reason, we stick to a simple and minimal low-level language. This low-level language is called `PFX` and is a so-called stack language (this will become clear when we will define it).

To summarise, we are going to work on two components:



These two components, the EXPR compiler and the PFX virtual machine follow a similar architecture: first they parse their input language to produce an internal data structure (the AST), and then treat this AST. However, they differ in their treatment. The EXPR compiler produces a corresponding PFX program while the PFX virtual machine *executes* its input.



For pedagogical reason and testing purpose, the PFX virtual machine is able to receive a (PFX) AST directly or to parse a file containing a PFX program.

2 Work to be done

The overall objective is to produce a working EXPR compiler and a working PFX virtual machine. The work will happen both at a mathematical (conceptual) level and at a programming level, the first being the specification of the second. You will be required to answer three kinds of questions:

- **expl:** you are supposed to produce a textual explanation
- **math:** you are supposed to give a set of mathematical definition, often in the form of a set of inference rules
- **code:** you are supposed to produce code; as usual, this code should be of quality (well indented, commented, with a good choice of identifiers, tested, ...)

The work is organised in a set of questions you are supposed to work on linearly. Moreover, at every step, you will have both a working EXPR compiler and a working PFX virtual machine. All along the subject, we enrich both the user language EXPR and the low-level language PFX.

During your work, you will need:

- the document containing lecture notes about *Compilation with OCaml* (`notes.pdf`). To succeed in writing your compiler, **the reading of this document is necessary**: it gives hints related to the compilation process, to the tools you will use (`ocamllex` and `menhir`). It also gives you the syntax and links to the official documentation.
- this document presenting the questions and the work to be done.

We also provide you an additional document to help you to remember the tools and commands: `cheatsheet.pdf`.

If you are proficient in OCaml and in the compilation field, just read the current document (please note that we also extend the language in order to handle simple functions and closures), then shut up and hack!

Up to the exercise you reach, all questions have to be answered. Code-related ones are answered by writing OCaml code and a quick note somewhere (comments within the code, quick summary in another file such as the README, etc.). Other questions have to be answered in a separated file, using an interoperable format (plain text or PDF for instance). It is obviously possible to answer mathematical questions by learning and using \LaTeX , however this is not the point of this course (and could be time-consuming...). An efficient way to answer such questions is to scan a *readable* handwritten answer.

3 How and what to deliver?

To help you with your deliverable, we provide a project skeleton (`project_skeleton.tar.gz`) which should be renamed, adapted and extended following your needs. Please also read the README example.

We expect you to deliver a compressed archive (`.tar.gz` file¹) containing:

- a PFX virtual machine written in OCaml, using `ocamllex` and `menhir`, wrt the specifications that are expressed throughout the questions,
- an `EXPR` compiler written in OCaml, using `ocamllex` and `menhir`, wrt the specifications that are expressed throughout the questions,
- a README (plain text file) explaining your work.

To be sure to be evaluated and to hope to obtain a strictly positive mark, you should also pay close attention to the form of your deliverable:

- deliver only one compressed archive (`tar.gz` format) of your top level directory (meaning that decompressing it should result in a single directory),
- please respect the following naming convention: `NAME1-NAME2.tar.gz`,
- file encoding: UTF-8 (especially if you write non-ASCII characters...),
- clean your directory: remove all useless files and files that can be generated (for example the `_build` directory),

¹`tar.gz` is not `zip`, `rar`, `7z`, `tar.xz`, `tar.bz2`, etc. and should be built using the `tar` tool.

- choose meaningful names (file names and contents) understandable by other people,
- comment your code,
- indent your code consistently,
- organise your files (use subdirectories if needed),
- join a README file and answers to the non-code-related questions.

Part II

Questions

4 A simple stack language

During our practice of compilation, we are going to use a stack language: PFX. As expressed by its name, it is a language in the tradition of Postscript² and Forth³ relying on a stack to store value instead of variables. In such a language, all operations act on the elements of the stack.

Exercise 1 (*expl*)

- ▷ **What is a stack? What are the operations that you usually execute on a stack?**

PFX is inspired by the Postfix machine of the book "Design Concepts in Programming Languages"⁴. It is intentionally kept simple. Whenever, you are not satisfied with it, feel free to extend and modify it... But not during the lab sessions.

4.1 Informal description

A PFX program begins with an integer specifying the number of arguments it will need to be run, then it is composed of a sequence of basic instructions. Basic instructions are `push`, `pop`, `swap` (it exchanges the first two elements of the stack) and the five arithmetic operations `add`, `sub`, `mul`, `div` and `rem`. Only `push` takes an argument which is an integer. All the arithmetic operations behave similarly: they use the first two elements of the stack as arguments, remove them and push back the result. This first version only manipulates integers so all values are integers.

For example, the program `0 push 12 push 7 sub` returns `-5` while `0 push 12 push 7 swap sub` returns `5`.

Program arguments are pushed onto the stack from last to first, before the execution of the program commands.

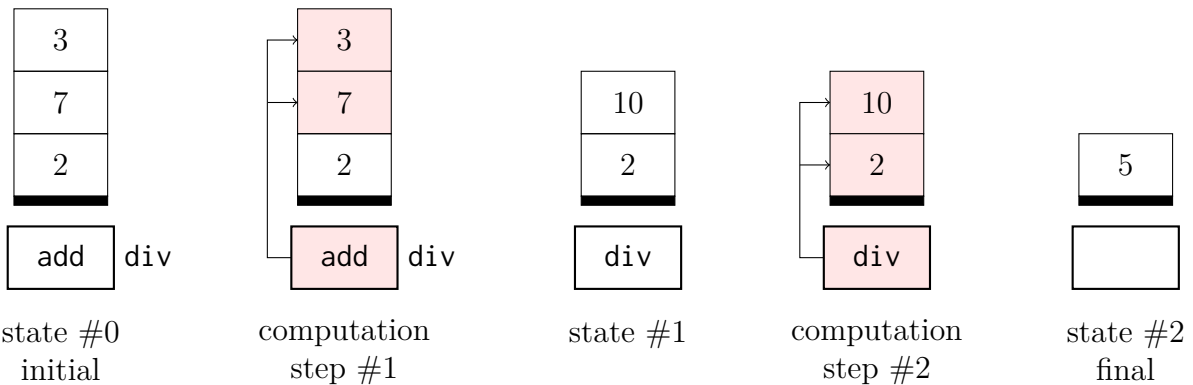
Here follow two examples of a program execution in detail.

²<https://en.wikipedia.org/wiki/PostScript>

³[https://en.wikipedia.org/wiki/Forth_\(programming_language\)](https://en.wikipedia.org/wiki/Forth_(programming_language))

⁴<https://mitpress.mit.edu/books/design-concepts-programming-languages>

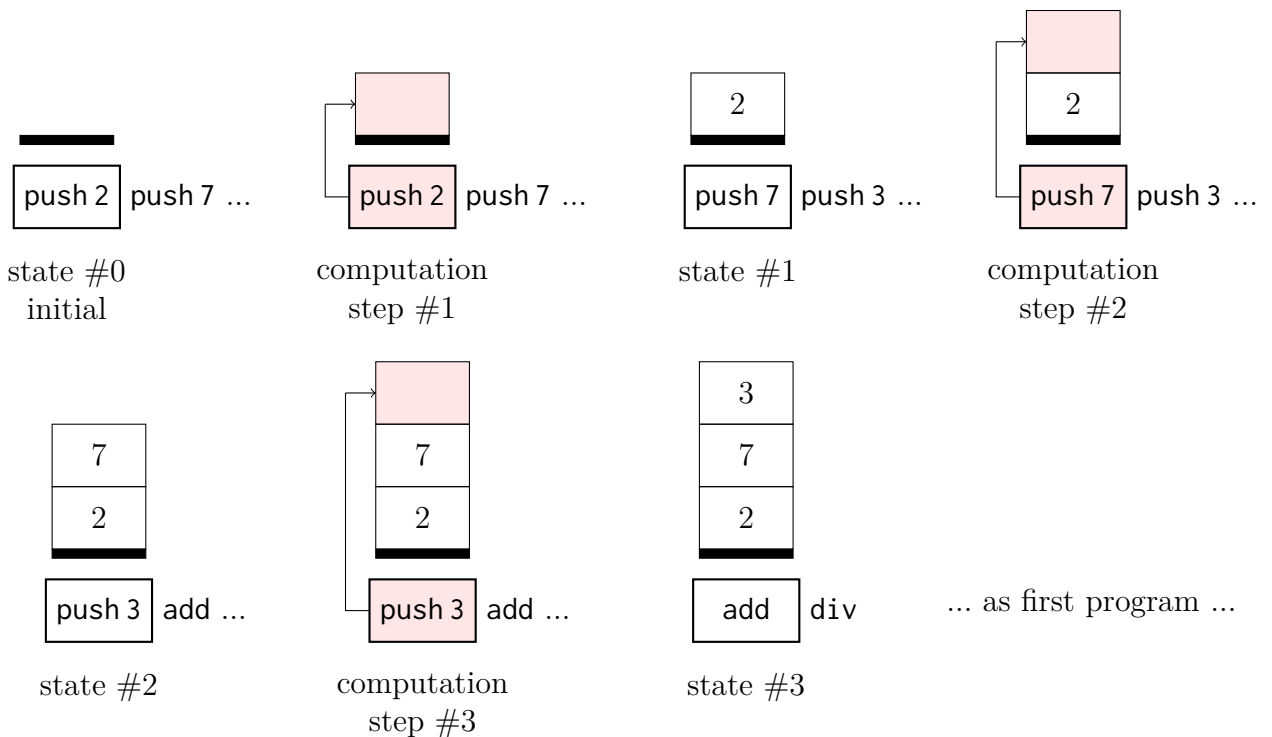
Program: 3 add div
Parameters: 3 7 2



⇒ result = 5

The first example is equivalent to the following one:

Program: 0 push 2 push 7 push 3 add div
Parameters: ∅



⇒ result = 5

Exercise 2 (expl)

▷ Detail in the same way the execution of 0 push 12 push 7 swap sub.

4.2 Formal semantics

To formalise the semantics of PFX, we will use the following notations:

- a program is a pair i, Q where i is the number of awaited arguments and Q the sequence of instructions composing the program, the empty instruction sequence is written \emptyset ,
- an instruction sequence is built by the constructor '.', the sequence composed of I then Q is $I.Q$,
- a stack is built by the constructor '::', adding an element n to the stack S is $n::S$, the empty stack is also written \emptyset , the length of the stack S can be obtained by $\#S$,

The semantics of PFX uses two sets of rules, one for programs that specify the execution of the complete program and one to describe each possible computational step on a current instruction sequence and the current stack. The rules for programs use the rule for computational steps.

Exercise 3

For programs, we have the following semantics:

$$(1) \frac{i \neq n}{v_1, \dots, v_n \vdash i, Q \Rightarrow \text{ERR}} \qquad (2) \frac{Q, v_1 :: \dots :: v_n :: \emptyset \rightarrow^* \text{ERR}}{v_1, \dots, v_n \vdash n, Q \Rightarrow \text{ERR}}$$

$$(3) \frac{Q, v_1 :: \dots :: v_n :: \emptyset \rightarrow^* \emptyset, v :: S}{v_1, \dots, v_n \vdash n, Q \Rightarrow v}$$

The reduction rule \rightarrow specifies the small step semantics of instructions and \rightarrow^* its transitive closure⁵. $Q, S \rightarrow Q', S'$ means that in one step the execution of instruction sequence Q with stack S leads to instruction sequence Q' and stack S' .

▷ **Question 3.1 (expl):**

Explain using plain words the semantics of programs.

▷ **Question 3.2 (math):**

A case is still missing, spot it out and give the corresponding rule.

▷ **Question 3.3 (math):**

Give the rules describing the small step semantics for instruction sequences. Beware to cover all cases of runtime errors.

4.3 Implementation

Exercise 4

▷ **Question 4.1 (code):**

Propose the OCaml code for a type **command** describing the Pfx instructions. It should be in the file **pfxAst.ml**.

⁵The rule is potentially applied several times in sequence. Mathematically, $A \rightarrow^* B$ if and only if $A = B$ or there exists C such that $A \rightarrow C$ and $C \rightarrow^* B$.

▷ **Question 4.2 (code):**

Write an OCaml function `step` that implements the small step reduction you defined above the Pfx instructions. It should be in the file `pfxEval.ml`.

⚠ You should test your code by running some PFX programs defined as OCaml values (not yet parsed from files, this will be done later).

5 A simple arithmetic expression language

Following on the exercise 6 from the *Discovering OCaml* document, we will define a simple arithmetic expression language, named `EXPR` in order to compile it to PFX.

This first version of `EXPR` AST is implemented in file `exprAst.ml` as follows:

```
type expression =
  | Const of int
  | Var of string
  | Binop of BinOp.t * expression * expression
  | Uminus of expression
```

where the type `BinOp.t` is implemented in file `binOp.ml` as follows:

```
type t =
  | Badd | Bsub | Bmul | Bdiv | Bmod
```

You can explore the code provided within the project skeleton.

Exercise 5

▷ **Question 5.1 (math):**

Propose a compilation schema of `Expr` in Pfx. Give its formal description. Notice that with the current definition of Pfx, we cannot implement variables. We defer their implementation to a later exercise.

▷ **Question 5.2 (code):**

Define a function `generate` implementing the semantics you defined in previous question. It should be in the file `exprToPfx.ml`.

⚠ At this point, you have a first working compiler of `EXPR` and you should be able to execute simple programs.

6 Parsing

Read the document containing the lecture notes entitled *Compilation with OCaml (notes.pdf)* in order to understand `ocamllex` and its syntax. You can explore the code for the lexer and the parser provided within the project skeleton.

Exercise 6 (*A first Pfx lexer*)

▷ Question 6.1 (*code*):

Write a lexer for the Pfx stack machine language. Complete the provided `pfxLexer.mll`. To test it without the parser, have a look at the file `exprLexer_standalone.mll` on Moodle.

The following OCaml code⁶ provides a way to read the string to parse from a file. The name of the file is given as an argument on the command line and is automatically passed to the function `compile`:

```
(* Entry point of the program, should contain your main function: here it is
   named parse_eval, it is the function provided after question 6.1 *)
(* The arguments, initially empty *)
let args = ref []
(* The main function *)
let parse_eval file =
  print_string ("File "^file^" is being treated!\n");
```

▷ Question 6.2 (*code*):

Reuse this code to be able to parse a file containing a Pfx program and prints all the tokens encountered in the process.

⚠ You should test your lexer and use test files.

Exercise 7 (*Locating errors, code*)

Generally, a compiler should be able to return an error message containing the location of the error to its user. OCaml module `Lexing` defines a type `position` for this purpose.

```
type position = {
  pos_fname : string; (* name of the file *)
  pos_lnum : int;     (* number of the line *)
  pos_bol : int;      (* nb of chars between the beginning of the file and the one of current line *)
  pos_cnum : int;      (* nb of characters since the beginning of the file *)
}
```

By default, the generated lexer only updates the last element (`pos_cnum`). The actions must take care of the others. The functions `lexeme_start_p` and `lexeme_end_p` of the module `Lexing` allows one to get respectively the location of the beginning and the end of the current token. To help you, we provide the module `Location`⁷ which defines helpful elements. For the moment, you should only use:

- the exception `Location.Error` carrying both a message and the location of the error;
- the type `Location.t` of a location composed of a starting position and an ending position;
- the function `Location.init` setting the file name of the given buffer;

⁶Provided in the project skeleton, on Moodle.

⁷As always, available on Moodle.

- the function `Location.incr_line` increasing a line in the given buffer;
- the function `Location.curr` return the current position of the given buffer;
- the function `Location.print` printing the given location.

▷ **Modify your code from the previous exercise to be able to return the location of errors.**

⚠ Now and for the rest of the UV, you should use `ocamlbuild` to compile and directly produce an executable file. You need to add the following line to bind the defined main function (here `compile`)⁸.

```
let _ = Arg.parse [] compile ""
```

Exercise 8 (*A first Pfx parser*)

▷ **Question 8.1 (*code*):**

Write a parser for the Pfx stack machine language.

▷ **Question 8.2 (*code*):**

Test it in combination with your Lexer. To do it, you will have to write a function that prints the AST of Pfx. You should now use the provided file `pfxVM.ml` as the main file for the Pfx virtual machine. It is the file that should be given to `ocamlbuild` as a target.

Notice that it requires that you modify slightly your lexer to remove the main functions and replace the token type definition by an open of the parser module.

⚠ You should test with more than one test your PFX parser!

7 Simple functions

Let's suppose we would like to add notions of function and application to `EXPR`. One way of doing it is to add a definition of function expression and an application expression. As in the λ -calculus, we limit ourselves to function with a unique argument. This leads to the following AST.

```
type expression =
  | Const of int
  | Var of string
  | Binop of BinOp.t * expression * expression
  | Uminus of expression
  (* For function support *)
  | App of expression * expression
  | Fun of string * expression
```

⁸For more details, do not hesitate to consult the manual section on the module `Arg`.

△ You can find the new AST and the modified lexer and parser in the directory `expr_fun` of the skeleton. You can replace the corresponding files in `expr` by their new version from `expr_fun`.

Currently, PFX cannot be used to generate code including functions. We first have to modify the language PFX. Three new instructions must be added:

- executable sequence (Q), where Q is an usual instruction sequence, when it is encountered the executable sequence is pushed on the top of the stack;
- `exec` which is an instruction that pops the top of the stack and executes it by appending it in front of the executing sequence, notice that the top of the stack must be an executable sequence;
- `get` pops the integer i on top of the stack, and copies on top of the stack the i -th element of the stack, it raises an error if there is not enough element on the stack.

Notice that these new constructions impose that the stack now contains two kinds of objects: integer or executable sequence.

Be sure to understand that PFX does not have functions, it only has executable sequence.

Exercise 9

▷ Question 9.1 (*expl*):

Do we need to change the rules for the already defined constructs?

▷ Question 9.2 (*math*):

Give the formal semantics of these new constructions.

▷ Question 9.3 (*code*):

If needed, extend the lexer and parser of Pfx to include these changes.

The translation of `EXPR` to `PFX` must be revised. The idea is that a λ -abstraction is translated to an executable sequence and an application is translated to an `exec` plus some code to *clean up* the stack. This executable sequence supposes that its parameter is on the top of the stack when it starts to execute. In its body, whenever it wants to use its parameter it `gets` it from where it is using an environment \mathcal{P} associating a variable to its position in the stack (its depth). Consequently, during the translation, we have to keep track of the depth at which each parameter is. When application terminates, it pops out the parameter from the stack. This behaviour is described in the figure 1:

- (a) the stack is in some state
- (b) a computation pushes the argument of the future call
- (c) another computation pushes the function to call, it must be an executable sequence
- (d) the call is made, probably using the stack and its argument
- (e) the call ends and pushes its results on the stack

Beware that during step (d), the stack may grow and therefore each time something is pushed onto the stack we need to update the position of the various reachable arguments.

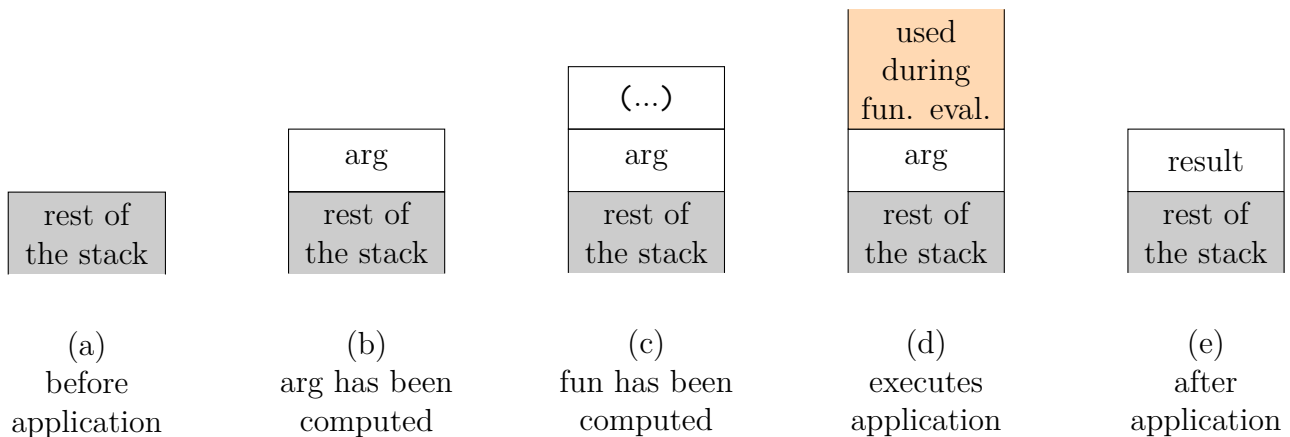


Figure 1: Usage of the stack during application

Exercise 10

- ▷ **Question 10.1 (*expl*):**
Give the compiled version of the expression $(\lambda x.x + 1) 2$. Then describe step by step the evaluation of its Pfx translation.
- ▷ **Question 10.2 (*math*):**
Give the formal rule for transformation.
- ▷ **Question 10.3 (*code*):**
Provide a new version of `generate`.
- ▷ **Question 10.4 (*expl*):**
Give the compiled version of the expression $((\lambda x.\lambda y.(x - y)) 12) 8$. Then describe step by step the evaluation of its Pfx translation. What do you think of the result? What is happening?

A straight forward extension of the syntax consists in adding a so-called *syntactic sugar* to support `let` definition (as of OCaml).

Exercise 11 (*Syntactic sugar*)

- ▷ **Question 11.1 (*expl*):**
What is the translation in the syntax of Expr already defined of a new `let $x = e_1$ in e_2` ?
- ▷ **Question 11.2 (*code*):**
What part of the code must be modified to get support for `let`?

8 Closure

To solve the problem illustrated in question 10.4, we need to change the function value. When defining a function value, we need to capture and store all its free variables (the variables defined

by its context). Indeed, the function may be executed in a different scope than the one where it was defined, so it must be able to retrieve the values of variables of its definition location when executed. This new kind of value is called a *closure* and it is a pair composed of an environment (of the free variables of the function) and the function. So the closure of a function $\lambda x.e$ in the context of an environment \mathcal{E} is denoted $\{\mathcal{E}, x, e\}$. Notice that the domain of \mathcal{E} is limited to set of free variables of $\lambda x.e$. When this domain is empty, we fall down on the function value of the previous section. This explains why and when it was working!

The formal semantics of **EXPR** is presented in figure 2. The values may be integer from set \mathbb{I} , closures (set $\mathbb{C} = X \rightarrow (\mathbb{I} \cup \mathbb{C}) \times X \times Expr^9$) or error (value **ERR**). All operations are extended to produce an error whenever they are applied to an error (*e.g.* $-\mathbf{ERR} = \mathbf{ERR}$ and $\mathbf{ERR} + v = \mathbf{ERR}$).

$$\begin{array}{c}
\text{(CONST)} \quad \mathcal{E} \vdash v \Rightarrow v \qquad \text{(VAR)} \quad \frac{x \in \text{dom}(\mathcal{E})}{\mathcal{E} \vdash x \Rightarrow \mathcal{E}(x)} \qquad \text{(VARERR)} \quad \frac{x \notin \text{dom}(\mathcal{E})}{\mathcal{E} \vdash x \Rightarrow \mathbf{ERR}} \\
\\
\text{(UMINUS)} \quad \frac{\mathcal{E} \vdash e \Rightarrow v}{\mathcal{E} \vdash -e \Rightarrow -v} \qquad \text{(BINOP)} \quad \frac{op \in \{+, -, *\} \quad \mathcal{E} \vdash e_1 \Rightarrow v_1 \quad \mathcal{E} \vdash e_2 \Rightarrow v_2}{\mathcal{E} \vdash e_1 \text{ op } e_2 \Rightarrow v_1 \text{ op } v_2} \\
\\
\text{(DIV)} \quad \frac{op \in \{/, \%\} \quad \mathcal{E} \vdash e_1 \Rightarrow v_1 \quad \mathcal{E} \vdash e_2 \Rightarrow v_2 \quad v_2 \neq 0}{\mathcal{E} \vdash e_1 \text{ op } e_2 \Rightarrow v_1 \text{ op } v_2} \\
\\
\text{(DIVERR)} \quad \frac{op \in \{/, \%\} \quad \mathcal{E} \vdash e_2 \Rightarrow 0}{\mathcal{E} \vdash e_1 \text{ op } e_2 \Rightarrow \mathbf{ERR}} \\
\\
\text{(APP)} \quad \frac{\mathcal{E} \vdash e_1 \Rightarrow \{\mathcal{E}_c, x, e\} \quad \mathcal{E} \vdash e_2 \Rightarrow v_2 \quad \mathcal{E}_c, x \mapsto v_2 \vdash e \Rightarrow v}{\mathcal{E} \vdash e_1 e_2 \Rightarrow v} \\
\\
\text{(APPERR1)} \quad \frac{\mathcal{E} \vdash e_1 \Rightarrow v_1 \quad v_1 \notin \mathbb{C}}{\mathcal{E} \vdash e_1 e_2 \Rightarrow v_1 \text{ op } v_2} \qquad \text{(APPERR2)} \quad \frac{\mathcal{E} \vdash e_1 \Rightarrow \{\mathcal{E}_c, x, e\} \quad \mathcal{E} \vdash e_2 \Rightarrow \mathbf{ERR}}{\mathcal{E} \vdash e_1 e_2 \Rightarrow \mathbf{ERR}} \\
\\
\text{(FUN)} \quad \frac{\text{dom}(\mathcal{E}') = \text{dom}(\mathcal{E}) \quad \forall x \in \text{dom}(\mathcal{E}), \mathcal{E}'(x) = \mathcal{E}(x)}{\mathcal{E} \vdash \lambda x.e \Rightarrow \{\mathcal{E}', x, e\}}
\end{array}$$

Figure 2: Big step operational semantics of **EXPR**

Exercise 12 (*math*)

- ▷ Give the proof derivation computing the value of the term of question 10.4 ($((\lambda x.\lambda y.(x - y)) 12) 8$).

⁹ X is the set of variables and $Expr$ the set of terms of **EXPR**.

Exercise 13

▷ **Question 13.1 (*expl*):**

Is it possible to translate Expr to Pfx? If yes, can you give the idea of the translation. If no, what would be necessary to add to Pfx ?

The chosen implementation of closure we are going to explore is the following.

A closure is an executable sequence that begins with instructions pushing onto the stack the values of the free variables. Notice that this part of the executable sequence must be added at runtime (the only moment when the values are known).

To enable the modification of an executable sequence at runtime we add a new construct to PFX: **append**. This construct expects the stack to contain a value on the top and an executable sequence below. It appends the command storing the "value" in the stack at the beginning of the executable sequence. When the value is an integer, this command is a push of it. The command is the value when it is an executable sequence. It also adds the operation to remove this value from the stack when the function ends.

▷ **Question 13.2 (*math*):**

Give the formal semantics of **append**.

▷ **Question 13.3 (*code*):**

If needed, extends the lexer and parser of Pfx to include these changes.

▷ **Question 13.4 (*math*):**

Give the formal rules of translation from Expr to Pfx to support closure.

▷ **Question 13.5 (*code*):**

Provide a new version of **generate**.

▷ **Question 13.6 (*expl*):**

Give the compiled version of the expression $((\lambda x.\lambda y.(x - y)) 12) 8$. Then describe step by step the evaluation of its Pfx translation. Is it better?

9 Extensions

To explore compilation in more depth, complete the following extensions by giving the formal version and then by implementing them. For each extension, a piece of information on the difficulty is provided.

1. (*regular*) Extend the stack machine with a **sel** instruction that requires a stack of at least three elements. Its result is the third element of the stack if the first is 0 and the second in the other case. Use this extension to add booleans, boolean operators, comparison operators and a condition operator $e_1?e_2 : e_3$ similar to the one of Java. Notice that the game is to try to add as little new instructions to Pfx as possible¹⁰.

¹⁰My solution just adds one.

2. (*medium*) Extend `EXPR` with a `let rec` construct enabling the definition of recursive functions. *Hint*¹¹. Propose a mechanism to compile it to `PFX`. *Hint*¹².
3. (*challenging*) Extend `PFX` with mechanisms to manipulate the stack. In this extension, a new value can be on the stack, a stack. Add the following constructs:
 - `pack` pushes a stack value containing current stack content after clearing it;
 - `unpack` pops the top of the stack which must be a stack value and replaces current stack by this value;
 - `switch` pops the top of the stack which must be a stack value, packs the rest of the current stack, replaces current stack by the stack value and pushes rest stack value.

These constructs provide so-called continuations. Use them to add to `EXPR`:

- (a) an exception mechanism,
- (b) basic threading

¹¹**Do not return before trying!**

Define a recursive function f by $\lambda f:\lambda v:e$ and a recursive closure $\{\mathcal{E}, f, x, e\}$. Evaluating such a recursive closure add it to the environment under the name f .

¹²**Do not return before trying!**

The easiest implementation is to define a notion of recursive instruction sequence that is kept on the stack when executed and only popped when terminating.