

Machine universelle

I. Introduction.

L'objectif de ce TP est de coder, en OCAML, la machine universelle \mathcal{U} permettant de calculer la fonction `interprète` définie par :

$$\text{interprète} (\langle \mathcal{M} \rangle, w) = \begin{cases} w' & \text{si } w \xrightarrow{\mathcal{M}} w' \\ \text{non défini} & \text{sinon} \end{cases}$$

pour une certaine machine \mathcal{M} . De même que dans le TD Théorie n°1, on utilise la notation $\langle \cdot \rangle$ pour dénoter la sérialisation. La fonction `interprète` aura la signature OCAML ci-dessous :

```
interprete : string -> string
```

L'entrée w a été supprimée : seul le code $\langle \mathcal{M} \rangle$ de la machine \mathcal{M} est donné à la fonction `interprete`. L'entrée de la machine est « ajoutée » au début du code, sous la forme d'une expression de la forme

```
let input = [on placera l'entrée de la machine ici] in
[on placera le code de la machine ici]
```

Ainsi définie, la fonction `interprete` aura le même fonctionnement que `utop` : il évalue l'expression, sans prendre en compte une certaine entrée.

Pour les deux premières sections, on utilisera le dossier `tp2-1`. Dans la troisième section (vérification de types), on utilisera le dossier `tp2-2`. On copiera les fonctions écrites dans le dossier `tp2-1` dans le dossier `tp2-2`.

Ce TP est long, on n'attendra pas qu'il soit fini en 2 heures. Mais, si vous voulez continuer ce TP par la suite, vous pouvez. Des informations sur l'environnement de développement sont disponibles en fin de sujet (Section V).

La chaîne de production d'un programme est divisée en plusieurs parties. Premièrement, l'*analyse lexicale* (ou *lexer* en anglais) décompose le code source en une liste de symboles (*tokens*). Puis, l'*analyse syntaxique* (ou *parser* en anglais) permet de construire un *arbre syntaxique* (*abstract syntax tree* en anglais, souvent abrégé AST). Et, l'*analyse sémantique* vérifie que l'arbre syntaxique a un sens : dans notre cas, ce sera la vérification des types (cette étape est optionnelle pour ce TP, mais elle est importante dans le cadre d'un véritable langage de programmation). Enfin, l'*interprétation* de l'arbre syntaxique exécute le code entré.

Dans ce TP, nous considérons un « sous »-langage de OCAML, nommé OCAML-- (« OCAML moins moins »). Il contient les structures fondamentales de OCAML :

- la définition de variables (par exemple `let x: int = 3 in`)
- les opérations `+`, `-`, `*`, `/` sur des entiers,
- les opérations `>`, `>=`, `<`, `<=`, `=` et `<>` sur des entiers (et aussi des booléens pour les deux derniers).
- les types `int` et `bool`,
- les fonctions anonymes `fun x -> 2 * x`.

II. Quelques notes sur l'analyse syntaxique et lexicale.

Dans la figure suivante (Fig. 1, page 2), on représente l'analyse lexicale et sémantique sur un exemple.

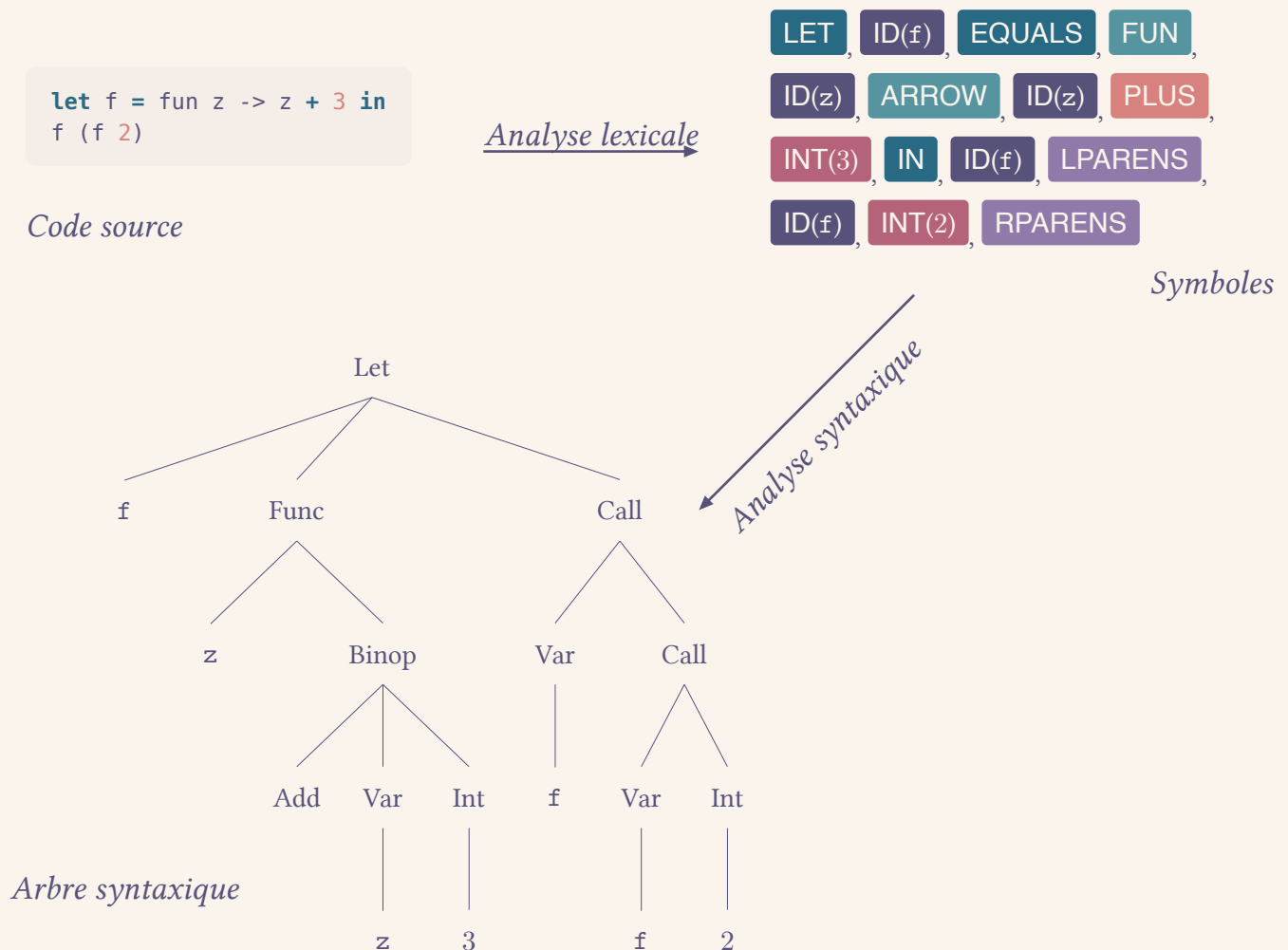


Fig. 1. Représentation de l'analyse lexicale et sémantique sur un exemple

L'analyse lexicale est réalisé à partir des automates : il suffit de définir des motifs à reconnaître, et comment les catégoriser (*i.e.* définir les symboles). L'analyse syntaxique est réalisée via des grammaires (notion vue plus tard dans l'année) : il suffit de définir les assemblages de symboles possibles, et quel impact cela aura sur l'arbre syntaxique. Cette partie n'est pas à réaliser, vous trouverez les fonctions nécessaires dans les documents associés au TP.

Le type `expr` représente l'arbre syntaxique, comme représenté sur la figure. Une fonction `parse` permet de convertir une chaîne de caractères en un arbre syntaxique. En cas d'erreur de syntaxe, une exception est levée ; on n'aura pas besoin de gérer ces erreurs de syntaxe.

```

type expr =
| Var   of string          (* Une variable *)
| Int   of int             (* Un entier *)
| Bool  of bool           (* Un booléen *)
| Binop of bop * expr * expr (* Une opération binaire *)
| Let   of string * expr * expr (* Une définition "let x = ... in ..." *)
| If    of expr * expr * expr (* Une condition "if, then, else" *)

```

```
| Func  of string * expr      (* Une fonction "fun z -> ..." *)
| Call  of expr * expr        (* Un appel d'expression *)
```

Le type `bop` contient les différents opérateurs binaires reconnus. Ainsi, il contient l'addition (`Add`), la soustraction (`Sub`), la multiplication (`Mult`), et la division (`Div`) de nombres entiers. Il contient également les opérateurs de comparaison (sur des entiers uniquement, retournant un booléen) `<` (`Lt`, *less than*), `>` (`Gt`, *greater than*), `≤` (`Leq`, *less than or equal to*), et `≥` (`Geq`, *greater than or equal to*). D'autres opérateurs pourront être définis, afin que OCAML-- se rapproche plus de OCAML.

Dans la suite du TP, le *parser* est la combinaison de l'analyseur sémantique et l'analyseur syntaxique.

- Q1.** Afin de vérifier le fonctionnement du *parser*, tester la fonction `parse` sur le code source de la Fig. 1. On comparera l'arbre obtenu à celui de la figure. On devrait obtenir le résultat ci-dessous.

```
Let ("f", Func ("z", Binop (Add, Var "z", Int 3)),
    Call (Var "f", Call (Var "f", Int 2)))
```

III. Interpréteur OCAML--.

Réaliser un interpréteur peut se réaliser de plusieurs manières. Dans ce TP, nous prendrons l'approche de la *substitution*. Dans ce modèle, l'interpréteur s'exécute sur l'arbre de syntaxe du programme. Les variables sont stockées dans l'arbre de syntaxe, et non dans une autre mémoire.

Généralement, ce modèle n'est pas utilisé car il modifie l'arbre de syntaxe du programme. Mais, l'implémentation de cette méthode est assez simple.

III.1. Formalisation, notations.

Le type OCAML `expr` est représenté mathématiquement par un ensemble \mathcal{E} défini par *induction nommée* à partir des règles **Var** représentant le variant `Var`, **Int** représentant le variant `Int`, **Bool** représentant le variant `Bool`, **Binop** représentant le variant `Binop`, **Let** représentant le variant `Let`, **If** représentant le variant `If`, **Func** représentant le variant `Func`, et **Call** représentant le variant `Call`. On identifiera, dans ces règles, le type `int` et l'ensemble \mathbb{Z} ; le type `bool` et l'ensemble $\mathbb{B} = \{\mathbf{V}, \mathbf{F}\}$; le type `string` et un ensemble fini de variables $\mathcal{V} = \{x, y, \dots\}$; les opérateurs du type `bop` et leurs équivalents mathématiques, ainsi que la notation `?` quand l'opérateur est inconnu.

Les variables et les expressions seront différenciées dans leur notation. Les variables seront notées en police machine à écrire droite, par exemple `x`. Les expressions seront notées en police roman italique, par exemple *x*.

III.2. Substitutions.

Étant données deux expressions e et $f \in \mathcal{E}$, et une variable `x`, on définit la *substitution* de `x` par f dans e , notée $e[f \mapsto x]$, par induction sur e :

- $(\mathbf{Var} \ x)[x \mapsto f] = f$,
- $(\mathbf{Var} \ y)[x \mapsto f] = \mathbf{Var} \ y$ avec $y \neq x$,
- $(\mathbf{Int} \ n)[x \mapsto f] = \mathbf{Int} \ n$,
- $(\mathbf{Bool} \ b)[x \mapsto f] = \mathbf{Bool} \ b$,
- $(\mathbf{Binop} \ ? \ e_1 \ e_2)[x \mapsto f] = \mathbf{Binop} \ ? \ e_1[x \mapsto f] \ e_2[x \mapsto f]$ pour une opération `?`,
- $(\mathbf{If} \ a \ b \ c)[x \mapsto f] = \mathbf{If} \ a[x \mapsto f] \ b[x \mapsto f] \ c[x \mapsto f]$,

- $(\mathbf{Call} e_1 e_2)[x \mapsto f] = \mathbf{Call} e_1[x \mapsto f] e_2[x \mapsto f]$,
- $(\mathbf{Func} x e)[x \mapsto f] = \mathbf{Func} x e$,
- $(\mathbf{Func} y e)[x \mapsto f] = \mathbf{Func} y e[x \mapsto f]$ si $x \neq y$ et que y n'est pas une variable libre de f ,
- $(\mathbf{Func} y e)[x \mapsto f] = \mathbf{Func} y e[x \mapsto f']$ si $x \neq y$ et que y est une variable libre de f , et que l'on a renommé y dans f en f' ,^[1]
- $(\mathbf{Let} x a b)[x \mapsto f] = \mathbf{Let} x a[x \mapsto f] b$,
- $(\mathbf{Let} x a b)[x \mapsto f] = \mathbf{Let} x a[x \mapsto f] b[x \mapsto f]$ si $x \neq y$.

On définit l'ensemble des *variables libres* $FV(e) \subseteq \mathcal{V}$ d'une expression $e \in \mathcal{E}$ par

- $FV(\mathbf{Var} x) = \{x\}$,
- $FV(\mathbf{Call} e f) = FV(e) \cup FV(f)$,
- $FV(\mathbf{Binop} ? e f) = FV(e) \cup FV(f)$,
- $FV(\mathbf{Int} n) = \emptyset$,
- $FV(\mathbf{Bool} b) = \emptyset$,
- $FV(\mathbf{Func} x e) = FV(e) \setminus \{x\}$,
- $FV(\mathbf{Let} x e_1 e_2) = FV(e_1) \cup (FV(e_2) \setminus \{x\})$,
- $FV(\mathbf{If} a b c) = FV(a) \cup FV(b) \cup FV(c)$.

On représentera en OCAML un ensemble par une liste. On s'assurera que chaque élément n'apparaît qu'une fois dans la liste représentant un ensemble.

Q2. Coder, en OCAML, une fonction `union : 'a list -> 'a list -> 'a list` qui réalise l'union de deux ensembles, représentés par des listes. (On pourra utiliser la fonction `List.mem` qui teste l'appartenance d'un élément dans une liste.)

On pourra définir un opérateur binaire `@@` tel que `e @@ f` réalise l'union de `e` et de `f`. Pour cela, il suffit d'ajouter le code ci-dessous.

```
let @@ = union
```

Pour plus d'informations sur les opérateurs binaires, regardez la documentation de OCAML, et la fiche « topo » de Mr Journault.

Q3. Définir une fonction `fv : expr -> string list` telle que `fv e` corresponde à l'ensemble des variables libres $FV(e)$, où e est la représentation mathématique de `e`.

Q4. Définir une fonction `substitution : expr -> expr -> string -> expr` qui réalise la substitution $e[x \mapsto f]$ lors de l'appel `substitution e f x`. On pourra supposer qu'aucune variable ne serait libre, quitte à modifier le code de la fonction `substitution`.

III.3. L'interprétation.

On définit la fonction partielle $\triangleright : \mathcal{E} \rightarrow \mathcal{E}$ par induction :

- $\triangleright(\mathbf{Var} x) = \square$,
- $\triangleright(\mathbf{Bop} ? a b) = \triangleright(a) ? \triangleright(b)$,
- $\triangleright(\mathbf{Int} n) = \mathbf{Int} n$,
- $\triangleright(\mathbf{Bool} b) = \mathbf{Bool} b$,
- $\triangleright(\mathbf{Let} x e f) = \triangleright(f[x \mapsto \triangleright(e)])$,
- $\triangleright(\mathbf{If} a b c) = \begin{cases} \triangleright(b) & \text{si } \triangleright(a) = \mathbf{Bool} V, \\ \triangleright(c) & \text{si } \triangleright(a) = \mathbf{Bool} F, \\ \square & \text{sinon,} \end{cases}$

^[1]Dans un premier temps, on n'implémentera pas ce cas. On supposera que toutes les variables ne seront pas libres.

- $\triangleright(\mathbf{Func}\ x\ e) = \mathbf{Func}\ x\ e,$
- $\triangleright(\mathbf{Call}\ e\ f) = \begin{cases} \triangleright(g[x \mapsto \triangleright(f)]) & \text{si } \triangleright(e) = \mathbf{Func}\ x\ g, \\ \square & \text{sinon.} \end{cases}$

Le « \square » est un élément inconnu, dans le cas où la fonction \triangleright n'est pas définie. Informatiquement, c'est une exception, une erreur.

Exemple. Avec l'arbre de la Fig. 1 (page 2), on a

$$\begin{aligned} & \triangleright(\mathbf{Let}\ x\ (\mathbf{Func}\ z\ (\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3)))(\mathbf{Call}\ (\mathbf{Var}\ f)(\mathbf{Call}\ (\mathbf{Var}\ f)(\mathbf{Int}\ 2)))) \\ = & \triangleright((\mathbf{Call}\ (\mathbf{Var}\ f)(\mathbf{Call}\ (\mathbf{Var}\ f)(\mathbf{Int}\ 2)))[f \mapsto \triangleright(\mathbf{Func}\ z\ (\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))]) \\ = & \triangleright((\mathbf{Call}\ (\mathbf{Func}\ z\ (\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3)))(\mathbf{Call}\ (\mathbf{Func}\ z\ (\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3)))(\mathbf{Int}\ 2))) \\ = & \triangleright((\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))[z \mapsto \triangleright(\mathbf{Call}\ (\mathbf{Func}\ z\ (\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3)))(\mathbf{Int}\ 2))] \\ = & \triangleright((\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))[z \mapsto \triangleright(\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))[z \mapsto (\mathbf{Int}\ 2)]] \\ = & \triangleright((\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))[z \mapsto \triangleright(\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Int}\ 2)(\mathbf{Int}\ 3))] \\ = & \triangleright((\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Var}\ z)(\mathbf{Int}\ 3))[z \mapsto \mathbf{Int}\ 5]) \\ = & \triangleright(\mathbf{Binop}\ \mathbf{Add}\ (\mathbf{Int}\ 5)(\mathbf{Int}\ 3)) \\ = & \mathbf{Int}\ 8 \end{aligned}$$

Q5. Définir des exceptions pour les différents cas non définis de la fonction \triangleright .

Q6. Coder la fonction \triangleright sous forme d'une fonction de signature `interp : expr -> expr`. On créera des fonctions mutuellement récursives pour gérer les différents cas de la définition de la fonction \triangleright : `interp_bop`, `interp_let`, `interp_if`, `interp_call`. Cette question est bien plus longue que les autres.

III.4. Interface utilisateur.

Pour le moment, nous avons défini des fonctions OCAML, mais nous n'avons pas créé un outil prêt à être utilisé. C'est ce que nous allons réaliser dans cette sous-section.

Q7. Coder une fonction `string_of_value : expr -> string` qui associe une chaîne de caractère à une *valeur* (donc uniquement booléen, et entier).

Q8. Coder la fonction `interprete : string -> string`.

Passer ce point, nous avons codé la machine universelle \mathcal{U} . Mais, nous allons continuer quand même. Le programme `utop` est un *REPL* (*Read-Eval-Print Loop*) : il lit la commande de l'utilisateur, l'exécute, puis affiche le résultat, et recommence.

Q9. Définir une fonction `get_input : unit -> string` qui récupère l'entrée de la console. Attention, le programme peut être sur plusieurs lignes : on ne s'arrêtera qu'après `;;`.

Q10. Définir une fonction `run_interpreteur : unit -> unit` qui lance l'interpréteur.

IV. Vérification de types.

Dans la suite de ce TP, nous allons vérifier les types dans le programme à exécuter.

IV.1. Formalisation mathématique, notations.

Nous définissons l'ensemble des *types* \mathcal{T} en OCAML-- comme un ensemble défini par induction nommée à partir des règles **IntType** et **BoolType** d'arité 0, et du type **FuncType** d'arité 2. Ainsi, on définit en OCAML le type associé `typ` aux éléments de \mathcal{T} .

```
type typ
| BoolType
| IntType
| FuncType of typ * typ
```

La « vérification de types » consiste à vérifier que le type de toute expression incluse dans le programme a un type correct, en supposant les sous-expressions sont correctement typées.

Un *environnement de typage* est un ensemble de couples variable type, donc un élément de $\wp(\mathcal{V} \times \mathcal{T})$ tel que, pour chaque variable $v \in \mathcal{V}$, il n'existe qu'un unique $t \in \mathcal{T}$, ou aucun tel que (v, t) soit un élément de l'environnement. On le notera $\text{env} = \{x \mapsto t_x, y \mapsto t_y, \dots\}$. On notera `bool` le type associé à la règle **BoolType**, `int` le type associé à la règle **IntType**, $e \rightarrow f$ le type associé à la règle **Func** $e f$. Comme env est un ensemble, on pourra noter $\text{env} \cup \{x \mapsto t\}$ pour ajouter un type. Si x avait déjà un type associé auparavant, il est redéfini par cet union.

Pour $e \in \mathcal{E}$, $t \in \mathcal{T}$ et un environnement env , on notera $\text{env} \vdash e : t$ dès lors que l'expression e a pour type t . On définit inductivement cette notation.

- On a toujours $\text{env} \vdash \mathbf{Int} \ n : \text{int}$ et $\text{env} \vdash \mathbf{Bool} \ b : \text{bool}$.
- On a $\text{env} \vdash \mathbf{Var} \ x : t$ où $(x, t) \in \text{env}$.
- Si $\text{env} \vdash e_1 : t_1$ et $\text{env} \cup \{x \mapsto t_1\} \vdash e_2 : t_2$, alors $\text{env} \vdash (\mathbf{Let} \ x \ e_1 \ e_2) : t_2$.
- Si $\text{env} \vdash e_1 : \text{int}$ et $\text{env} \vdash e_2 : \text{int}$, alors $\text{env} \vdash (\mathbf{Binop} \ ? \ e_1 \ e_2) : \text{int}$ pour $? \in \{+, -, \times, /\}$.
- Si $\text{env} \vdash e_1 : \text{int}$ et $\text{env} \vdash e_2 : \text{int}$, alors $\text{env} \vdash (\mathbf{Binop} \ ? \ e_1 \ e_2) : \text{bool}$ pour $? \in \{\leq, \geq, <, >\}$.
- Si $\text{env} \vdash a : \text{bool}$, $\text{env} \vdash b : t$ et $\text{env} \vdash c : t$, alors $\text{env} \vdash (\mathbf{If} \ a \ b \ c) : t$.
- Si $\text{env} \cup \{x \mapsto t_1\} \vdash f : t_2$, alors $\text{env} \vdash (\mathbf{Func} \ z \ f) : t_1 \rightarrow t_2$.
- Si $\text{env} \vdash f : t_1$ et $\text{env} \vdash e : t_1 \rightarrow t_2$, alors $\text{env} \vdash (\mathbf{Call} \ e \ f) : t_2$.

Pour stocker un environnement, on utilisera une liste de couples variables type. On définit le type `environment` comme ci-dessous.

```
type environment = (string * typ) list
```

Q11. Définir une fonction `find_type : environment -> string -> typ option` qui, lors de l'appel `find_type v env`, renvoie `Some t` si $(v, t) \in \text{env}$, et `None` sinon.

On modifie le `parser` pour que les variants `Let` et `Func` contiennent également le type (de type `typ`) de la variable définie.

Q12. Ajuster le code défini précédemment avec cette modification.

Q13. Définir une fonction `type_of : environment -> expr -> typ` qui, à une expression e et un environnement env , associe un type t tel que $\text{env} \vdash e : t$. On réalisera, comme pour la question , des fonctions mutuellement récursives nommées `type_of_bop`, `type_of_if`, ...

Q14. Définir une fonction `type_check : expr -> expr`. Cette fonction sera l'identité^[2] pour toute expression correctement typée. Pour les expressions non typée, on renverra une erreur.

^[2]On rappelle $\text{id}_A : x \mapsto x$ est la fonction identité de A .

V. Environnement de développement

Ce TP utilise `dune`, un environnement de développement pour OCAML. On peut lancer `utop` via la commande `dune utop`. Une fois ouvert, il suffira d'exécuter `open Interpreter.Main;;` pour importer le fichier `main.ml`. Si la commande `dune` n'est pas reconnue, on pourra l'installer via `opam` (le gestionnaire de paquets d'OCAML) :

```
opam install dune
```

En cas de difficultés d'installation, consultez le site <https://dune.build/>. Si ces difficultés persistent, prévenez et un *Repl.it* sera mis en ligne.