

École Normale Supérieure
Langages de programmation et compilation
examen 2017–2018

Jean-Christophe Filliâtre
26 janvier 2018

Les notes de cours manuscrites ou reprographiées sont les seuls documents autorisés.
L'épreuve dure 3 heures.

Dans tout ce sujet, on considère une variante de mini-ML dont la syntaxe abstraite est la suivante :

$e ::= n$	constante entière $n \in \mathbb{Z}$
x	variable
$e - e$	soustraction
fun $x \rightarrow e$	fonction anonyme
$e e$	application
print e	affichage
if $e \leq 0$ then e else e	conditionnelle
fix e	point fixe

Dans la suite, on pourra écrire **let** $x = e_1$ **in** e_2 comme un sucre syntaxique pour $(\mathbf{fun} \ x \rightarrow e_2) \ e_1$. Dans ce langage, un programme est réduit à une expression.

Les parties 1 et 2 ne sont pas indépendantes, la partie 2 utilisant des définitions et des résultats de la partie 1. Les parties 3 et 4 sont indépendantes entre elles et des parties 1 et 2.

1 Sémantique

On munit ce langage d'une sémantique opérationnelle à petits pas de la forme

$$e \xrightarrow{a} e'$$

où e et e' sont deux expressions et a l'affichage réalisé par le pas de réduction. Cet affichage vaut soit \emptyset , lorsque rien n'est affiché, soit un entier n , lorsque le pas d'exécution affiche l'entier n . La sémantique opérationnelle à petits pas est donnée figure 1. On prendra le temps de bien la lire et notamment de relever les petites différences avec celle vue en cours. Une *évaluation* d'une expression e est une séquence de réductions qui conduit à une valeur, c'est-à-dire

$$e \xrightarrow{a_1} e_1 \xrightarrow{a_2} e_2 \cdots \xrightarrow{a_n} v.$$

Question 1 Donner *une* évaluation de l'expression

$$(\mathbf{print} \ 2) - ((\mathbf{print} \ 60) - (\mathbf{print} \ 100))$$

Combien y a-t-il d'évaluations possibles pour cette expression ? Justifier.

valeurs

$v ::= n$ constante
| $\text{fun } x \rightarrow e$ fonction

affichage

$a ::= \emptyset$ on n'affiche rien
| n on affiche n

réductions de tête

$n_1 - n_2 \xrightarrow{\emptyset} n$ avec n l'entier $n_1 - n_2$
 $(\text{fun } x \rightarrow e) v \xrightarrow{\emptyset} e[x \leftarrow v]$
 $\text{print } n \xrightarrow{n} n$
 $\text{if } n \leq 0 \text{ then } e_1 \text{ else } e_2 \xrightarrow{\emptyset} e_1$ si $n \leq 0$
 $\text{if } n \leq 0 \text{ then } e_1 \text{ else } e_2 \xrightarrow{\emptyset} e_2$ si $n > 0$
 $\text{fix } (\text{fun } x \rightarrow e) \xrightarrow{\emptyset} e[x \leftarrow \text{fix } (\text{fun } x \rightarrow e)]$

contextes de réduction

$E ::= \square$
| $E - e$
| $e - E$
| $E e$
| $e E$
| $\text{print } E$
| $\text{if } E \leq 0 \text{ then } e \text{ else } e$
| $\text{fix } E$

réduction

$$\frac{e_1 \xrightarrow{a} e_2}{E(e_1) \xrightarrow{a} E(e_2)} \quad \text{pour } E \neq \square$$

FIGURE 1 – Sémantique opérationnelle à petits pas.

Question 2 Donner une expression e telle que, pour tout entier n , toute évaluation de e n affiche exactement un entier, égal à la factorielle de n .

Confluence du calcul. On va montrer maintenant que le résultat d'un calcul ne dépend pas de l'ordre d'évaluation, tant qu'on ignore les affichages. Dans les deux questions suivantes, on s'autorise à écrire $e \rightarrow e'$ pour une réduction dont on ne précise pas l'affichage.

Question 3 Montrer que si $e \rightarrow e_l$ et $e \rightarrow e_r$ alors soit $e_l = e_r$, soit il existe e' telle que $e_l \rightarrow e'$ et $e_r \rightarrow e'$.

Question 4 En déduire que si $e \rightarrow e_1 \rightarrow e_2 \cdots \rightarrow e_n$ et $e \rightarrow e'_1 \rightarrow e'_2 \cdots \rightarrow e'_{n'}$ alors il existe e_m telle que $e_n \rightarrow e_{n+1} \rightarrow e_{n+2} \cdots \rightarrow e_m$ ($m - n$ étapes) et $e'_{n'} \rightarrow e'_{n'+1} \rightarrow e'_{n'+2} \cdots \rightarrow e_m$ ($m - n'$ étapes) avec $m \leq n + n'$.

2 Typage avec effets

Dans cette partie, on va typer notre langage avec des types enrichis d'une notion d'*effet*. Un effet, noté ϕ , prend trois valeurs possibles :

- $\phi ::= \perp$ le calcul ne produit aucun affichage
- | \mathbf{P} le calcul peut produire un affichage, qui ne dépend pas de l'ordre d'évaluation
- | \mathbf{T} le calcul peut produire un affichage, qui peut dépendre de l'ordre d'évaluation

Les effets sont naturellement ordonnés de la manière suivante

$$\perp \leq \mathbf{P} \leq \mathbf{T}$$

et on s'autorisera à écrire $\max(\dots)$ pour calculer le maximum de plusieurs effets. On définit par ailleurs une opération binaire \sqcup sur les effets, de la manière suivante :

$$\begin{aligned} \perp \sqcup \phi &\stackrel{\text{def}}{=} \phi \\ \phi \sqcup \perp &\stackrel{\text{def}}{=} \phi \\ \phi_1 \sqcup \phi_2 &\stackrel{\text{def}}{=} \mathbf{T} \quad \text{sinon} \end{aligned}$$

Un environnement de typage, noté Γ , associe des types à des variables. Le jugement de typage prend la forme

$$\Gamma \vdash e : \tau \ \& \ \phi$$

et se lit comme « dans l'environnement Γ , l'expression e a le type τ et l'effet ϕ ». Les types sont ici des types simples, de la forme

$$\begin{aligned} \tau ::= \mathbf{int} &\quad \text{type des entiers} \\ | \tau \xrightarrow{\phi} \tau &\quad \text{type des fonctions} \end{aligned}$$

Le type d'une fonction contient un effet au-dessus de la flèche, appelé *effet latent*. C'est l'effet de l'évaluation d'une application de cette fonction. Les règles de typage sont données figure 2.

$$\begin{array}{c}
\frac{}{\Gamma \vdash n : \text{int} \ \& \ \perp} \quad \frac{x \in \text{dom}(\Gamma)}{\Gamma \vdash x : \Gamma(x) \ \& \ \perp} \quad \frac{\Gamma \vdash e_1 : \text{int} \ \& \ \phi_1 \quad \Gamma \vdash e_2 : \text{int} \ \& \ \phi_2}{\Gamma \vdash e_1 - e_2 : \text{int} \ \& \ \phi_1 \sqcup \phi_2} \\
\\
\frac{\Gamma \vdash x : \tau_1 \vdash e : \tau_2 \ \& \ \phi}{\Gamma \vdash \text{fun } x \rightarrow e : \tau_1 \xrightarrow{\phi} \tau_2 \ \& \ \perp} \quad \frac{\Gamma \vdash e_1 : \tau_2 \xrightarrow{\phi} \tau_1 \ \& \ \phi_1 \quad \Gamma \vdash e_2 : \tau_2 \ \& \ \phi_2}{\Gamma \vdash e_1 \ e_2 : \tau_1 \ \& \ \max(\phi_1 \sqcup \phi_2, \phi)} \\
\\
\frac{\Gamma \vdash e : \text{int} \ \& \ \phi}{\Gamma \vdash \text{print } e : \text{int} \ \& \ \max(\phi, \text{P})} \quad \frac{\Gamma \vdash e_1 : \text{int} \ \& \ \phi_1 \quad \Gamma \vdash e_2 : \tau \ \& \ \phi_2 \quad \Gamma \vdash e_3 : \tau \ \& \ \phi_3}{\Gamma \vdash \text{if } e_1 \leq 0 \ \text{then } e_2 \ \text{else } e_3 : \tau \ \& \ \max(\phi_1, \phi_2, \phi_3)} \\
\\
\frac{\Gamma \vdash e : (\tau_1 \xrightarrow{\phi} \tau_2) \xrightarrow{\phi_3} (\tau_1 \xrightarrow{\phi_1} \tau_2) \ \& \ \phi_3 \quad \phi \stackrel{\text{def}}{=} \max(\phi_1, \phi_2)}{\Gamma \vdash \text{fix } e : \tau_1 \xrightarrow{\phi} \tau_2 \ \& \ \phi_3}
\end{array}$$

FIGURE 2 – Typage avec effets.

Question 5 Donner trois expressions e_1 , e_2 et e_3 ayant respectivement les types suivants :

$$\begin{array}{l}
\vdash e_1 : \text{int} \ \& \ \top \\
\vdash e_2 : \text{int} \xrightarrow{\text{P}} \text{int} \ \& \ \text{P} \\
\vdash e_3 : (\text{int} \xrightarrow{\top} \text{int}) \xrightarrow{\text{P}} \text{int} \ \& \ \perp
\end{array}$$

Question 6 L'expression suivante est-elle typable ?

$$\text{fun } f \rightarrow f (\text{fun } x \rightarrow 0) - f (\text{fun } x \rightarrow \text{print } 0)$$

Si oui, lui donner un type. Sinon, justifier qu'elle n'est pas typable.

Sûreté du typage. Comme vu en cours, la sûreté du typage se déduit des résultats de progrès et de préservation, qui font l'objet des deux questions suivantes.

Question 7 Montrer la propriété de progrès : si $\vdash e : \tau \ \& \ \phi$, alors soit e est une valeur, soit il existe e' et a tels que $e \xrightarrow{a} e'$.

Question 8 Montrer la propriété de préservation du typage : si $\Gamma \vdash e : \tau \ \& \ \phi$ et $e \xrightarrow{a} e'$, alors $\Gamma \vdash e' : \tau \ \& \ \phi'$ avec $\phi' \leq \phi$.

On admettra la propriété de substitution : si $\Gamma \vdash x : \tau' \vdash e : \tau \ \& \ \phi$ et $\Gamma \vdash e' : \tau' \ \& \ \perp$ alors $\Gamma \vdash e[x \leftarrow e'] : \tau \ \& \ \phi$.

Correction de l'effet \perp . On va maintenant chercher à montrer que si l'effet d'une expression est \perp , alors son évaluation ne produit aucun affichage.

Question 9 Montrer que si $\Gamma \vdash e : \tau \ \& \ \phi$ et $e \xrightarrow{n} e'$, alors $\phi \neq \perp$.

Question 10 Dédire du résultat précédent que, si $\Gamma \vdash e : \tau \ \& \ \phi$ et $e \xrightarrow{a_1} e_1 \xrightarrow{a_2} e_2 \cdots \xrightarrow{a_n} e_n$ et s'il existe i tel que $a_i \neq \emptyset$, alors $\phi \neq \perp$.

Correction de l'effet P. Enfin, on va chercher à montrer que si l'effet d'une expression est P, alors son évaluation produit un affichage qui ne dépend pas de l'ordre d'évaluation. Pour cela, on se donne une opération binaire \cdot sur les affichages, associative et d'élément neutre \emptyset .

Question 11 Montrer que si $\Gamma \vdash e : \tau \ \& \ \phi$, avec $\phi \leq P$, et si $e \xrightarrow{a_l} e_l$ et $e \xrightarrow{a_r} e_r$, avec $e_l \neq e_r$, alors il existe e' telle que $e_l \xrightarrow{a'_l} e'$, $e_r \xrightarrow{a'_r} e'$ et $a_l \cdot a'_l = a_r \cdot a'_r$.

Question 12 Montrer que si $\Gamma \vdash e : \tau \ \& \ \phi$, avec $\phi \leq P$, et si $e \xrightarrow{a_1} e_1 \xrightarrow{a_2} e_2 \cdots \xrightarrow{a_n} e_n$ et $e \xrightarrow{a'_1} e'_1 \xrightarrow{a'_2} e'_2 \cdots \xrightarrow{a'_{n'}} e'_{n'}$, avec $e_n = e'_{n'} = v$ une valeur, alors $a_1 \cdot a_2 \cdots a_{n-1} \cdot a_n = a'_1 \cdot a'_2 \cdots a'_{n'-1} \cdot a'_{n'}$ et $n = n'$.

On pourra admettre que si $e \xrightarrow{a} e'$ et $e \xrightarrow{a'} e'$ alors $a = a'$.

3 Production de code

On se propose maintenant de compiler notre langage vers l'assembleur x86-64. Comme expliqué en cours, on procède en deux temps, en commençant par une étape d'*explicitation des fermetures*. Plus précisément, un programme e est traduit vers un ensemble de définitions de fonctions globales d_1, \dots, d_n d'une part et une expression c d'autre part, dans la syntaxe abstraite suivante :

$$\begin{array}{l}
 c ::= n \\
 \quad | x \\
 \quad | c - c \\
 \quad | \text{clos } f [E] \\
 \quad | c c \\
 \quad | \text{print } c \\
 \quad | \text{if } c \leq 0 \text{ then } c \text{ else } c \\
 \quad | \text{fix } c \\
 d ::= \text{letfun } f [E] x = c \\
 E ::= x, \dots, x
 \end{array}$$

Dans ce langage, la construction $\text{fun } x \rightarrow e$ a été remplacée par une opération explicite de construction de fermeture $\text{clos } f [E]$ où f est maintenant une fonction globale et $E = x_1, \dots, x_m$ est l'*environnement*, c'est-à-dire la liste des variables libres de $\text{fun } x \rightarrow e$. Cette liste est ordonnée de façon arbitraire. On note E_i le i -ième élément de la liste E , les éléments étant numérotés à partir de 1.

Question 13 Donner le résultat de l'explicitation des fermetures sur le programme

```
print ((fix (fun fib → fun n → if n - 1 <= 0 then n else (fib (n - 1)) - (0 - (fib (n - 2)))))) 10)
```

Assembleur. L'étape suivante consiste à traduire ce langage intermédiaire vers l'assembleur x86-64. On adopte le schéma de compilation suivant :

- une valeur est soit un entier signé 64 bits, soit un pointeur vers une fermeture sur le tas (toutes les valeurs ont donc la même taille) ;
- une fermeture est un bloc alloué sur le tas contenant $n+1$ mots de 64 bits, le premier contenant une adresse de code et les suivants les valeurs de l'environnement ;
- la valeur d'une expression est calculée dans `%rax` ;
- dans le code d'une fonction :

```

compile(E, n) = movq $n, %rax
compile(E, x) = movq %rdi, %rax      si x ∉ E
compile(E, x) = movq 8i(%rsi), %rax  si x = Ei
compile(E, c1 - c2) = ...
compile(E, clos f [x1, ..., xn]) = ...
compile(E, c1 c2) = ...
compile(E, print c) = ...
compile(E, if c1 <= 0 then c2 else c3) = ...
compile(E, fix c) = compile(E, c)
                    pushq %rax
                    pushq %rdi
                    pushq %rsi
                    movq $16, %rdi
                    call malloc
                    popq %rsi
                    popq %rdi
                    movq $_fix, (%rax)
                    popq %rcx
                    movq %rcx, 8(%rax)

compile(letfun f E x = c) = f :
                           compile(E, c)
                           ret

```

FIGURE 3 – Compilation vers x86-64.

-
- la valeur de l'argument est contenue dans `%rdi`,
 - la valeur de la fermeture est contenue dans `%rsi` ;
 - `%rcx` est utilisé comme temporaire.

La figure 3 contient une partie du compilateur. Une expression c est compilée par un appel à $compile(E, c)$ où E est l'environnement de la fermeture, lorsqu'on compile le corps d'une fonction, et une liste vide sinon, lorsqu'on compile le programme principal. On notera qu'une variable qui n'est pas dans E est nécessairement l'argument de la fonction. Un aide-mémoire x86-64 est donné en annexe.

Question 14 Donner le code du compilateur pour

- la soustraction $c_1 - c_2$;
- la conditionnelle `if c1 <= 0 then c2 else c3` ;
- l'application $c_1 c_2$.

Point fixe. Pour compiler la construction `fix c`, on choisit de représenter son résultat comme une fermeture de 16 octets, le premier mot contenant l'adresse du code d'une fonction assembleur `_fix` et le second la valeur de c , qui se trouve être une fermeture (c'est garanti par le typage).

Question 15 Justifier l'utilisation des instructions `pushq` et `popq` sur les registres `%rdi` et `%rsi` dans la compilation de la construction `fix` donnée figure 3.

Question 16 Donner le code assembleur de la fonction `_fix`.

Question 17 Quel est l'intérêt de représenter le résultat de `fix c` par une fermeture ?

4 Analyse syntaxique

On cherche maintenant à définir une syntaxe concrète pour notre langage, qui puisse être reconnue par un analyseur LR(1). Pour lever un certain nombre d'ambiguïtés telles que

```
print 1 - 2
```

on choisit de limiter les expressions en position d'argument d'une application, de `print` et de `fix` à des expressions réduites à des constantes entières, des identificateurs et des expressions parenthésées. On écrit donc notre grammaire sous la forme suivante :

```

E → S
   | E - E
   | E S
   | let ident = E in E
   | fun ident -> E
   | if E <= 0 then E else E
   | fix S
   | print S
S → constante
   | ident
   | ( E )
```

Les non terminaux sont E et S . Tous les autres symboles sont terminaux. Les mots `let`, `in`, `fun`, `if`, `then`, `else`, `fix` et `print` sont des mots-clés.

Question 18 Une telle grammaire est cependant toujours refusée par l'outil `menhir` comme n'étant pas LR(1). Expliquer pourquoi.

Question 19 Proposer des règles de priorités et d'associativités qui permettent à `menhir` d'accepter cette grammaire (*i.e.* sans signaler de conflit).

Annexe : aide-mémoire x86-64

On donne ici un fragment du jeu d'instructions x86-64. Vous êtes libre d'utiliser tout autre élément de l'assembleur x86-64. Dans ce qui suit, r_i désigne un registre, n une constante entière et L une étiquette.

<code>mov r_2, r_1</code>	copie le registre r_2 dans le registre r_1
<code>mov $\\$n, r_1$</code>	charge la constante n dans le registre r_1
<code>mov $\\$L, r_1$</code>	charge l'adresse de l'étiquette L dans le registre r_1
<code>sub r_2, r_1</code>	calcule $r_1 - r_2$ et l'affecte à r_1
<code>mov $n(r_2), r_1$</code>	charge dans r_1 la valeur contenue en mémoire à l'adresse $r_2 + n$
<code>mov $r_1, n(r_2)$</code>	écrit en mémoire à l'adresse $r_2 + n$ la valeur contenue dans r_1
<code>push r_1</code>	empile la valeur contenue dans r_1
<code>pop r_1</code>	dépile une valeur dans le registre r_1
<code>test r_2, r_1</code>	positionne les drapeaux en fonction de la valeur de r_1 ET r_2
<code>jg L</code>	saute à l'adresse désignée par l'étiquette L si les drapeaux signalent un résultat signé > 0
<code>jmp L</code>	saute à l'adresse désignée par l'étiquette L
<code>call L</code>	saute à l'adresse désignée par l'étiquette L , après avoir empilé l'adresse de retour
<code>jmp $*o$</code>	saute à l'adresse désignée par l'opérande o
<code>call $*o$</code>	saute à l'adresse désignée par l'opérande o , après avoir empilé l'adresse de retour
<code>ret</code>	dépile une adresse et y effectue un saut

On alloue de la mémoire sur le tas avec un appel à `malloc`, qui attend un nombre d'octets dans `%rdi` et renvoie l'adresse du bloc alloué dans `%rax`.