

École Normale Supérieure
Langages de programmation et compilation
examen 2021–2022

Jean-Christophe Filliâtre
21 janvier 2022 — 8h30–11h30

Les notes de cours manuscrites ou reprographiées sont les seuls documents autorisés.

L'épreuve dure 3 heures.

Les figures sont regroupées en fin de sujet, page 7. Suggestion : détacher la dernière feuille.

Dans tout ce sujet, on considère un petit langage de programmation dont la syntaxe abstraite est donnée figure 1. Dans ce langage, les seules valeurs sont des listes, construites à partir de la liste vide `[]` et de l'opération `::`. Les fonctions sont notées f . Elles sont récursives et uniquement globales. Chaque fonction a un seul argument.

Ce langage est moins pauvre qu'il n'y paraît. On peut en effet encoder des données en utilisant uniquement des listes. Par exemple, on peut représenter un entier n comme une liste contenant n éléments tous égaux à `[]`. Ainsi, l'entier 2 est représenté par la liste `[] :: [] :: []`. On peut également représenter une paire (x, y) par la liste à deux éléments `x :: y :: []`, et s'affranchir ainsi du fait que les fonctions n'ont qu'un seul argument. Avec toutes ces idées, on peut proposer les définitions suivantes

```
rec fst x = match x with | [] -> [] | y :: _ -> y
rec snd x = match x with | [] -> [] | _ :: y -> fst y
rec add p = match fst p with
  | [] -> snd p
  | _ :: x -> [] :: add (x :: snd p :: [])
```

où la fonction `add` calcule donc la somme de deux entiers.

L'instruction `print` affiche la valeur d'une expression, suivie d'un retour chariot. Les éléments d'une liste sont affichés entre crochets, chaque élément étant suivi d'un point-virgule et d'un espace. Ainsi, avec les définitions ci-dessus, l'évaluation de l'instruction

```
print
  let three = [] :: [] :: [] :: [] in
  add (three :: three :: [])
```

va provoquer l'affichage suivant

```
[[]; []; []; []; []; []; ]
```

c'est-à-dire une liste qui représente l'entier 6.

Question 1 Donner la définition d'une fonction `fib` telle que, si n est une liste représentant l'entier n , alors `fib n` est une liste représentant l'entier F_n avec

$$\begin{cases} F_0 & = 0 \\ F_1 & = 1 \\ F_{n+2} & = F_n + F_{n+1} \quad \text{pour } n \geq 0. \end{cases}$$

Sémantique opérationnelle. On souhaite munir notre langage d'une sémantique opérationnelle à petits pas, sous la forme d'une relation $e \rightarrow e'$ définie entre deux expressions par

$$\frac{e_1 \xrightarrow{\epsilon} e_2}{E(e_1) \rightarrow E(e_2)}$$

où E désigne un contexte d'évaluation et $\xrightarrow{\epsilon}$ désigne une réduction en tête.

Question 2 Proposer une définition pour les contextes E et pour la réduction $\xrightarrow{\epsilon}$, qui traduisent un passage par valeur et une évaluation de la gauche vers la droite. Ici, les valeurs sont définies par

$$v ::= [] \mid v :: v$$

et on suppose qu'on a défini une opération de substitution $e[x \leftarrow v]$.

Indications : Il doit y avoir six cas dans la définition de E et quatre dans la définition de $\xrightarrow{\epsilon}$.

Analyse syntaxique. On souhaite réaliser l'analyse syntaxique de notre petit langage avec l'outil Menhir. La figure 2 contient un fichier d'entrée pour Menhir contenant la grammaire de notre langage. Les actions sémantiques ne nous intéressent pas ici et sont omises ($\{\dots\}$).

Question 3 Lorsque l'outil Menhir est lancé sur le fichier donné en figure 2, il déclare quatre conflits de type lecture/réduction (**shift/reduce**). Les identifier, les expliquer, faire le choix de favoriser lecture ou réduction et modifier en conséquence le fichier Menhir.

Compilation vers x86-64. On se propose maintenant de compiler notre petit langage vers l'assembleur x86-64. (Un aide-mémoire est donné en annexe.) On adopte le schéma de compilation suivant. Une valeur est soit l'entier 0 (sur 64 bits), représentant la liste vide, soit un pointeur vers un bloc de 16 octets alloué sur le tas, contenant deux valeurs v_1 et v_2 et représentant la liste $v_1 :: v_2$. Les valeurs ont donc toutes la même taille (64 bits). Pour toute définition **rec** $f\ x = e$ contenue dans le programme, on introduit une fonction assembleur f , recevant son paramètre x sur la pile. Le tableau d'activation de la fonction a la forme ci-contre, les adresses croissant vers le haut. Dans ce tableau, on trouve les variables locales à l'expression e , qu'elles soient introduites par **let** ou par **match**. Toutes les variables sont accessibles à partir de **%rbp**, soit avec +16 pour le paramètre, soit avec -8, -16, etc., pour les variables locales. Le résultat de la fonction est renvoyé dans le registre **%rax**.

paramètre
adresse retour
%rbp sauvegardé
variable locale
⋮
variable locale

Question 4 Donner un code assembleur pour la fonction **fst** suivante :

```
rec fst x = match x with | [] -> [] | y :: _ -> y
```

Allocation des variables. Pour déterminer la position de chaque variable sur la pile, on se donne une fonction $A(a, n, e)$ qui renvoie une expression identique à e , où chaque variable est remplacée par sa position sur la pile. Le paramètre a est une fonction donnant les positions des variables déjà déterminées. Le paramètre n est la prochaine position libre sur la pile. Ainsi, pour une définition de fonction **rec** $f\ x = e$, on va réaliser l'allocation avec $A(\{x \mapsto 16\}, -8, e)$, c'est-à-dire que l'on indique que le paramètre x est à la position +16 et que la prochaine position libre est -8.

Question 5 Proposer une définition de la fonction A pour chacune des six constructions du langage. Sans chercher à faire quelque chose de trop sophistiqué, on essaiera de proposer un schéma qui permet d'allouer plusieurs variables sur un même emplacement.

Question 6 Donner une expression e contenant plusieurs variables locales, dont

- au moins deux peuvent être allouées au même emplacement ;
- au moins deux ne peuvent pas être allouées au même emplacement.

Compilation. La compilation d'une expression e est un code assembleur noté $C(e)$ dont l'exécution a pour effet de stocker la valeur de l'expression e dans le registre `%rax`.

Question 7 Donner la définition de $C(e)$ pour chacune des six constructions du langage.

Question 8 Proposer un code assembleur pour une fonction `print` qui imprime une liste au format décrit au début de l'énoncé. On suppose que l'argument de `print` est passé dans le registre `%rdi`. On se servira de la fonction `putchar` pour imprimer un caractère (voir l'annexe). On cherchera à optimiser tous les appels terminaux. *Attention : un seul retour chariot doit être imprimé, tout à la fin.*

Optimisation. Dans la fonction `add` donnée en exemple au tout début de l'énoncé, à savoir

```
rec add p = match fst p with
| []      -> snd p
| _ :: x -> [] :: add (x :: snd p :: [])
```

l'appel récursif à `add` n'est pas un appel terminal, car il faut ensuite ajouter un élément au début de la liste renvoyée par cet appel récursif. Il existe cependant un moyen de le transformer en un appel terminal. Pour cela, on alloue le bloc correspondant à `[] :: ...` avant de faire l'appel récursif et on le passe à la fonction pour qu'elle y stocke son résultat, avec une affectation, plutôt que de le renvoyer. On obtient alors une fonction de la forme suivante :

```
rec add_dps d p = match fst p with
| []      -> d.next <- snd p
| _ :: x -> let d' = [] :: [] in
            let _ = d.next <- d' in
            let d = d' in
            add_dps d (x :: snd p :: [])
```

L'argument `d` est la *destination* du calcul, c'est-à-dire l'emplacement où il faut mettre le résultat. L'opération `d.next <- ...` modifie cet emplacement pour y stocker une valeur (dans le second champ du bloc mémoire, ici intitulé `next` puisqu'il s'agit de listes). Cette version de la fonction `add` ne renvoie plus rien, mais elle stocke son résultat dans `d`. On constate que son appel récursif, sur la dernière ligne, est maintenant un appel terminal. En particulier, il peut être optimisé pour conduire à une exécution qui se fera en espace de pile constant.

Pour retrouver une fonction `add` qui calcule la même chose que la fonction d'origine, il suffit alors de la définir ainsi :

```

rec add p = match fst p with
| []      -> snd p
| _ :: x -> let d = [] :: [] in
            let _ = add_dps d (x :: snd p :: []) in
            d

```

On note que cette fonction `add` n'est pas récursive : elle se charge uniquement du premier appel à la fonction récursive `add_dps`, lorsque la liste `fst p` n'est pas vide.

Ce style de programmation s'appelle *destination passing style*. On se propose de réaliser cette optimisation sur notre petit langage. Pour cela, on étend la syntaxe des expressions avec de nouvelles constructions :

$$\begin{aligned}
 e ::= & \dots \\
 & | \quad d.\text{next} \leftarrow e \\
 & | \quad f^{dps} \ x \ e
 \end{aligned}$$

L'expression $d.\text{next} \leftarrow e$ affecte la valeur de e au second champ de la liste désignée par la variable d . Cette construction suppose que la variable d désigne une liste non vide. Cette construction ne renvoie rien (mais pour la faire rentrer dans le formalisme précédent, on pourrait supposer qu'elle renvoie une valeur non significative, comme par exemple la liste vide).

Par ailleurs, une définition de fonction `rec f x = e` dans le langage d'origine devient maintenant une paire de fonctions mutuellement récursives

$$\begin{aligned}
 \text{rec } f \ x &= T(\text{false}, e) \\
 \text{and } f^{dps} \ d \ x &= T(\text{true}, e)
 \end{aligned}$$

La fonction f se comporte comme la fonction originale. La fonction f^{dps} ne renvoie rien et se comporte comme $d.\text{next} \leftarrow f \ e$. Le premier paramètre formel de f^{dps} s'appelle toujours d . La fonction T réalise l'optimisation. Elle prend en argument un booléen indiquant si on se trouve dans la fonction f ou dans la fonction f^{dps} et une expression à traduire.

Question 9 Donner un exemple de fonction justifiant le caractère mutuellement récursif des fonctions f et f^{dps} .

Question 10 Proposer une définition de la fonction T pour chacune des six constructions originales du langage.

Auto-interprète. Pour conclure en beauté, on se propose d'explorer dans les questions suivantes l'idée d'écrire un interprète de notre langage *dans le langage lui-même*, c'est-à-dire écrire une fonction

```

rec eval p = ...

```

qui reçoit en argument le codage sous forme de listes d'un programme et renvoie la liste des (codes des) caractères imprimés par l'évaluation de ce programme. On considère ici uniquement le langage original, tel que défini dans la figure 1, c'est-à-dire sans les extensions introduites dans les deux questions précédentes.

- Pour écrire cet auto-interprète, on reprend notamment deux idées introduites au début du sujet :
- un entier n est représenté par une liste contenant n éléments tous égaux à `[]` ;
 - une paire (x, y) est représentée par la liste à deux éléments $x :: y :: []$.

Question 11 Donner la définition d'une fonction `equal` qui reçoit en argument (le codage d')une paire (x, y) et renvoie `[]` si les listes x et y sont identiques et `[] :: []` sinon. La comparaison des listes x et y doit être faite en profondeur.

Liste d'association. Une liste d'association est une liste de paires (k, v) où k est une clé et v une valeur. Ainsi, la liste

`([] :: ([] :: []) :: []) :: (([] :: [] :: []) :: ([] :: [] :: [] :: []) :: []) :: []`

associe à la clé `[]` la valeur `[] :: []` et à la clé `[] :: [] :: []` la valeur `[] :: [] :: [] :: []`.

Question 12 Donner la définition d'une fonction `assoc` qui reçoit en argument (le codage d')une paire (k, a) , où a est une liste d'association, et renvoie la première valeur v associée à k dans a , le cas échéant, et `[]` si la clé k n'apparaît pas dans a .

Question 13 Proposer un codage pour la syntaxe abstraite du langage, telle qu'elle est définie dans la figure 1. On distinguera un codage pour les expressions (e) , un autre pour les déclarations (d) et un troisième pour les programmes (p) .

Question 14 Donner les grandes lignes de la fonction `eval`, c'est-à-dire sa décomposition en sous-fonctions et la spécification de chacune de ces sous-fonctions. On ne demande pas d'écrire le code de ces fonctions.

$e ::= x$	<i>variable</i>
$f e$	<i>application</i>
$\text{let } x = e \text{ in } e$	<i>variable locale</i>
$[]$	<i>liste vide</i>
$e :: e$	<i>cons</i>
$\text{match } e \text{ with } [] \rightarrow e \mid x :: x \rightarrow e$	<i>filtrage</i>
$d ::= \text{rec } f x = e$	<i>fonction réursive</i>
$\text{print } e$	<i>affichage</i>
$p ::= d \dots d$	<i>programme</i>

FIGURE 1 – Syntaxe abstraite.

```

%token <string> IDENT
%token LET IN REC PRINT MATCH WITH
%token EQ BAR NIL CONS ARROW LEFTPAR RIGHTPAR EOF
%start program
%type <...> program
%%
program:
| list(decl) EOF           {...}
decl:
| REC IDENT IDENT EQ expr {...}
| PRINT expr               {...}
expr:
| IDENT                    {...}
| IDENT expr               {...}
| LET IDENT EQ expr IN expr {...}
| NIL                      {...}
| expr CONS expr           {...}
| MATCH expr WITH
  BAR? NIL ARROW expr BAR IDENT CONS IDENT ARROW expr {...}
| LEFTPAR expr RIGHTPAR   {...}

```

FIGURE 2 – Fichier Menhir.

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 la valeur à l'adresse L dans le registre r_1
<code>mov $\\$L, r_1$</code>	charge l'adresse de l'étiquette L dans le registre r_1
<code>add r_2, r_1</code>	calcule la somme de r_1 et r_2 dans r_1 (on a de même <code>sub</code> et <code>imul</code>)
<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>cmp r_2, r_1</code>	positionne les drapeaux en fonction de la valeur de $r_1 - r_2$
<code>test r_2, r_1</code>	positionne les drapeaux en fonction de la valeur de $r_1 \& r_2$
<code>je L</code>	saute à l'adresse désignée par l'étiquette L en cas d'égalité (on a de même <code>jne</code> , <code>kg</code> , <code>jge</code> , <code>jl</code> et <code>jle</code>)
<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>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`.

On imprime un caractère avec un appel à `putchar`, qui attend le code du caractère dans `%rdi`. Dans le code assembleur, on peut écrire (le code d')un caractère c symboliquement, avec la syntaxe `$$c'`. Ainsi, `movq $$'a', %rdi` met le code du caractère `a` dans le registre `%rdi`.