

École Normale Supérieure
Langages de programmation et compilation
examen 2018–2019

Jean-Christophe Filliâtre

25 janvier 2019

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

Dans ce tout ce sujet, on étudie FORTH, un langage à piles datant de 1970. Un programme FORTH est une suite de *mots* séparés par des caractères blancs (espace, tabulation, retour-chariot). Un mot peut être composé de n'importe quelle séquence de caractères non blancs. Ainsi, le programme

```
42 dup emit emit 10 emit
```

contient exactement six mots. L'exécution d'un programme FORTH utilise une pile contenant des valeurs entières. Initialement la pile est vide. On exécute alors chaque mot en séquence. Si le mot correspond à une opération prédéfinie, comme `dup` ou `emit`, on effectue cette opération. Par exemple, `dup` duplique la valeur se trouvant au sommet de la pile et `emit` dépile un entier et affiche le caractère dont cet entier est le code ASCII. Si en revanche le mot à exécuter ne correspond pas à une opération prédéfinie, alors ce mot doit être une constante entière et sa valeur est mise au sommet de la pile. Ainsi, le programme ci-dessus affiche deux fois le caractère `*` (de code ASCII 42) suivi d'un retour-chariot (de code ASCII 10).

Les deux mots `if` et `endif` sont interprétés de façon particulière. Chaque occurrence d'un mot `if` doit nécessairement être associée à un mot `endif` se trouvant plus loin dans le texte du programme. Lorsque le mot `if` est rencontré, on dépile une valeur. Si elle est nulle, l'exécution se poursuit après la première occurrence du mot `endif` qui suit. Sinon, l'exécution se poursuit immédiatement après le mot `if`. Lorsque l'exécution arrive plus tard au mot `endif`, celui-ci est ignoré et l'exécution se poursuit au-delà. Ainsi, le programme

```
dup if dup emit endif drop
```

a pour effet de dépiler un entier et d'afficher le caractère correspondant uniquement s'il est non nul. L'opération prédéfinie `drop` dépile une valeur. L'ensemble des opérations prédéfinies est donné figure 1.

Question 1 Donner un programme qui attend deux entiers n_1 et n_2 au sommet de la pile et affiche soit le caractère `!` (code ASCII 33) s'ils sont différents, soit le caractère `=` (code ASCII 61) s'ils sont égaux. Les deux entiers doivent être dépilés au final.

Grammaire. Pour effectuer l'analyse syntaxique du langage FORTH, on se donne la grammaire suivante :

$$\begin{array}{l} S \rightarrow L \# \\ C \rightarrow \text{word} \\ \quad | \text{if } L \text{ endif} \\ L \rightarrow \epsilon \\ \quad | C L \end{array}$$

L'ensemble des terminaux est $\{\text{word}, \text{if}, \text{endif}, \#\}$ et l'ensemble des non terminaux est $\{S, C, L\}$. Le symbole de départ est S .

primitive	sémantique
dup	duplique le sommet de la pile
emit	dépile un entier n et affiche le caractère de code ASCII n
drop	dépile un entier
-	dépile un entier n_2 puis un entier n_1 et empile l'entier $n_1 - n_2$
0=	dépile un entier n et empile soit -1 si $n = 0$, soit 0 si $n \neq 0$
@	dépile un entier a et empile la valeur se trouvant en mémoire à l'adresse a
!	dépile un entier a puis un entier n et écrit l'entier n en mémoire à l'adresse a

FIGURE 1 – Opérations primitives.

Question 2 Calculer NULL, FIRST et FOLLOW pour chacun des non terminaux de cette grammaire.

Question 3 Cette grammaire est-elle LL(1) ? LR(0) ? SLR(1) ? LR(1) ? ambiguë ? Justifier à chaque fois.

Variables et fonctions. Le langage FORTH contient également des déclarations de variables et de fonctions. Une variable est déclarée avec la syntaxe

variable *ident*

où *ident* est le nom de la variable. Après une telle déclaration, le nom *ident* peut être utilisé et son effet est de déposer sur la pile l'adresse où la variable est stockée en mémoire. Grâce aux opérations primitives ! et @ données figure 1, on peut donc accéder au contenu d'une variable et le modifier. Ainsi, le programme

```
variable x
42 x !
x @ emit
```

déclare une variable **x**, y stocke l'entier 42, puis affiche le caractère dont le code est stocké dans la variable. Une fonction est déclarée avec la syntaxe

: *ident* ...code... ;

La définition commence avec le mot « : » et se termine avec le mot « ; ». Le mot immédiatement après « : », noté ici *ident*, est le nom de la fonction. Le corps de la fonction, noté ici « ...code... », est une liste de mots obéissant à la grammaire L donnée plus haut. En particulier, le corps de la fonction ne peut pas contenir récursivement de définition de fonction, ni une définition de variable. On appelle une fonction en donnant son nom et sa définition est alors exécutée. Ainsi, le programme

```
: print
  dup if dup emit endif drop ;
: print3
  dup print dup print print 10 emit ;
```

définit deux fonctions, dont une fonction **print3** qui dépile un entier et affiche trois fois le caractère correspondant, s'il est non nul, puis un retour-chariot. Le corps d'une fonction ne peut faire référence qu'à des fonctions et des variables définies *précédemment* dans le programme. En particulier, une fonction ne peut pas être récursive.

```

type code =
  | Wrd of string          (* un mot          *)
  | If  of code list      (* if ... endif  *)
type decl =
  | Var of string         (* variable x     *)
  | Def of string * code list (* : f ... ;     *)
  | Run of code list      (* du code à exécuter *)
type file =
  decl list

```

FIGURE 2 – Types OCaml pour la syntaxe abstraite (dans un fichier `ast.mli`).

```

%{ open Ast %}

%token <string> WORD
%token COLON SEMICOLON VARIABLE IF ENDIF
%token EOF

%start file
%type <Ast.file> file

%%

file:
| EOF          { []      }
| d=decl f=file { d :: f }
;

decl:
| VARIABLE x=WORD          { Var x          }
| COLON f=WORD c=codes SEMICOLON { Def (f, c)  }
| c=code l=codes          { Run (c :: l)  }
;

codes:
|          { []      }
| c=code l=codes { c :: l }
;

code:
| IF l=codes ENDIF { If l }
| w=WORD          { Wrd w }
;

```

FIGURE 3 – Analyse syntaxique avec Menhir.

Question 4 Donner le code d'une fonction `swap` qui échange les deux valeurs au sommet de la pile. Est-il possible de le faire en utilisant une seule variable ?

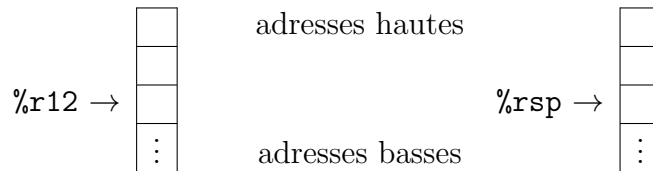
Analyse syntaxique avec Menhir. La figure 3 contient la syntaxe complète du langage FORTH, au format de l'outil Menhir. Volontairement, on n'utilise pas ici les fonctionnalités de Menhir pour reconnaître des listes (`+`, `*`, `nonempty_list`, etc.). On note qu'il y a cinq mot-clés, à savoir

```
if   endif  variable  :   ;
```

Cet analyseur syntaxique renvoie des arbres de syntaxe abstraite dont les types OCaml sont donnés figure 2.

Question 5 Lorsque l'outil Menhir est lancé sur ce fichier, il déclare un conflit de type *shift/reduce*. Identifier et expliquer ce conflit. Expliquer ensuite comment le résoudre en ajoutant dans ce fichier des indications de priorité et/ou d'associativité. Les règles de grammaire ne doivent pas être modifiées.

Compilation vers x86-64. On se préoccupe maintenant de la compilation du langage FORTH vers l'assembleur x86-64. (Un petit aide-mémoire x86-64 est donné à la fin du sujet.) Toutes les valeurs manipulées occupent 64 bits. Elles sont interprétées soit comme des adresses, soit comme des entiers, selon le contexte. Il y a *deux piles* matérialisées dans la mémoire. La pile de FORTH, dont nous avons parlé jusqu'à présent, est matérialisée quelque part sur le tas et son sommet est repéré grâce au registre `%r12`. Elle croît en direction des adresses basses. L'autre pile est la pile usuelle de x86-64, dont le sommet est repéré par le registre `%rsp` et qui croît également en direction des adresses basses. La première pile est appelée la *pile de données* et la seconde est appelée la *pile de retour*.



On pourra supposer que la compilation du programme commence par quelque chose comme

```
movq $0x4000, %rdi
call sbrk
leaq 0x4000(%rax), %r12
```

qui alloue de l'espace sur le tas pour la pile de données (ici 16 ko). La pile de retour s'utilise de la manière usuelle. En particulier, l'instruction assembleur `call` dépose l'adresse de retour au sommet de cette pile et l'instruction `ret` la dépile et y revient.

La compilation d'un programme FORTH est très simple. Toutes les variables sont allouées statiquement sur le segment de données. Chaque fonction FORTH devient une fonction assembleur, dont le code est obtenu en compilant successivement toutes les instructions qui la composent. La compilation de chaque instruction est donnée figure 4. Enfin, la compilation d'un programme contient également le code assembleur de chaque opération primitive. Ainsi, le code assembleur de la primitive `emit` est le suivant :

```
emit:
    movq (%r12), %rdi
    addq $8, %r12
    call putchar
    ret
```

<code>compile(x)</code>	=	<code>addq \$-8, %r12</code>	
		<code>movq \$x, (%r12)</code>	si x est une variable ou un entier
<code>compile(x)</code>	=	<code>call x</code>	si x est une fonction
<code>compile(if code... endif)</code>	=	<code>movq (%r12), %rcx</code>	
		<code>addq \$8, %r12</code>	
		<code>testq %rcx, %rcx</code>	
		<code>jz L</code>	
		<code>compile(code)...</code>	
		<code>L:</code>	avec L une étiquette fraîche
<code>compile(: x code...;)</code>	=	<code>x: compile(code)... ret</code>	compilation d'une fonction
<code>compile(variable x)</code>	=	<code>.data</code>	compilation d'une variable
		<code>x: .quad 0</code>	

FIGURE 4 – Compilation de FORTH vers x86-64.

Question 6 Donner le code assembleur des primitives `dup`, `drop`, `-`, `0=`, `@` et `!`.

Des boucles. Pour l'instant, les programmes que nous pouvons écrire sont extrêmement limités, car il n'y a aucune notion de boucle. (On rappelle que les fonctions ne sont pas récursives). On va y remédier en ajoutant deux nouvelles opérations primitives, `begin` et `until`, de telle sorte que

`begin code1... until code2...`

exécute `code1` tant que la valeur trouvée en sommet de pile (et dépilée) par `until` est nulle. Ensuite, l'exécution se poursuit avec `code2`. Ainsi, la fonction `many-stars` définie par

```
: many-stars
  begin 42 emit 1 - dup 0= until drop 10 emit ;
```

dépile un entier n , supposé strictement supérieur à zéro, et affiche n fois le caractère `*`, suivi d'un retour-chariot.

Question 7 Donner le code FORTH d'une fonction `fibonacci` qui trouve au sommet de la pile un entier $n \geq 1$ et affiche n lignes de texte contenant respectivement F_1, F_2, \dots, F_n caractères `*`, dans cet ordre¹. Ainsi, l'exécution de `6 fibonacci` doit donner la sortie

```
*
*
**
***
****
*****
*****
```

On pourra se resservir de la fonction `many-stars` donnée plus haut. L'entier n doit être dépilé au final.

1. S'il faut le rappeler, on a $F_0 = 0$, $F_1 = 1$ et $F_n = F_{n-2} + F_{n-1}$ pour $n \geq 2$.

Compilation de begin et until. Nous allons ajouter les deux primitives `begin` et `until` *sans changer* ni la syntaxe, ni la compilation de notre langage, mais uniquement en fournissant du code assembleur pour ces deux nouvelles primitives. L'idée est d'utiliser la pile de retour pour y stocker l'adresse du code situé immédiatement après `begin` et pouvoir ainsi y revenir autant de fois que nécessaire lorsque l'on exécute `until`.

Question 8 Donner le code assembleur des primitives `begin` et `until`.

Compilation d'un langage structuré vers FORTH. On se donne un petit langage de programmation impératif et structuré, avec des variables globales et des valeurs entières. Les expressions sont formées de constantes, variables et soustractions, avec la syntaxe abstraite suivante :

$$\begin{array}{ll}
 e ::= n & \text{constante entière} \\
 & | x \quad \text{variable} \\
 & | e - e \quad \text{soustraction}
 \end{array}$$

Les instructions comprennent des affectations, des conditionnelles, des boucles `while` et des séquences, avec la syntaxe abstraite suivante :

$$\begin{array}{ll}
 s ::= x \leftarrow e & \text{affectation} \\
 & | \text{if } e = 0 \text{ then } s \text{ else } s \quad \text{conditionnelle} \\
 & | \text{while } e = 0 \text{ do } s \quad \text{boucle} \\
 & | s ; s \quad \text{séquence} \\
 & | \text{putchar}(e) \quad \text{afficher un caractère} \\
 & | \text{skip} \quad \text{ne rien faire}
 \end{array}$$

On se propose de compiler ce petit langage vers le langage FORTH. Plus précisément, on cherche à définir

- pour toute expression e , un programme FORTH $compile(e)$ dont l'exécution a pour effet de déposer la valeur de e au sommet de la pile ;
- pour toute instruction s , un programme FORTH $compile(s)$ dont l'exécution produit les mêmes affichages que le programme s et laisse la pile inchangée (au final).

On suppose que toutes les variables apparaissant dans le programme initial ont été déclarées comme autant de variables au début du programme FORTH que l'on construit.

Question 9 Donner la définition de $compile(e)$ pour chacune des trois constructions et la définition de $compile(s)$ pour chacune des six constructions.

Estimation de la taille de pile. Dans cette dernière partie, on va chercher à estimer l'utilisation de la pile par un programme FORTH. Plus précisément, pour chaque morceau de code c , on va calculer cette estimation sous la forme d'un intervalle $[x, y]$ avec la signification « l'exécution du code c augmente la taille de pile d'une valeur comprise entre x et y ». La borne inférieure x est soit un entier, soit $-\infty$; la borne supérieure y est soit un entier, soit $+\infty$. Lorsque ce sont deux entiers, on a nécessairement $x \leq y$. Le problème n'étant pas décidable, on va se contenter d'une approximation correcte.

Question 10 Donner un code FORTH pour lequel la meilleure estimation est $[3, +\infty]$.

Question 11 Donner l'intervalle pour chacune des primitives `dup`, `emit`, `drop`, `-`, `0=`, `@` et `!`. (On ne demande pas de traiter ici les primitives `begin` et `until`.)

Question 12 On suppose que les primitives `begin/until` sont bien parenthésées et bien imbriquées avec les `if/endif`. Dit autrement, on peut se donner la syntaxe abstraite OCaml suivante pour le code FORTH :

```
type code =
  | Wrd of string      (* un mot      *)
  | If  of code list  (* if ... endif  *)
  | Beg of code list  (* begin ... until *)
```

Donner alors le code d'une fonction OCaml `eval: code -> bounds` qui renvoie l'intervalle d'estimation pour un code donné. L'intervalle est représenté par le type OCaml

```
type bounds = int option * int option
```

où la valeur `None` représente une valeur infinie. On pourra supposer l'existence de deux tables de hachage globales contenant respectivement les variables et les fonctions connues. La table des fonctions donne, pour chaque fonction, son intervalle d'estimation (supposé déjà calculé).

* *
* *

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>jg</code> , <code>jge</code> , <code>j1</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

Énoncé en français.