

École Polytechnique  
INF564 : Compilation  
examen 2022 (X2019)

Jean-Christophe Filiâtre

14 mars 2022

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

L'épreuve dure 3 heures.

Les deux problèmes sont indépendants.

## 1 Opérateur &

Dans ce problème, on considère un tout petit fragment du langage C dont la syntaxe abstraite est donnée figure 1. Dans ce fragment, on a des entiers et pointeurs. Si  $p$  est un pointeur de type  $\tau^*$ , vers un emplacement mémoire contenant une valeur de type  $\tau$ , alors  $*p$  désigne la valeur de type  $\tau$  à cet emplacement. Si  $e$  est une valeur gauche de type  $\tau$  alors  $\&e$  désigne un pointeur vers cet emplacement, et ce pointeur a le type  $\tau^*$ . Dans ce fragment, une valeur gauche est soit une variable, soit une expression de la forme  $*e$ .

Les fonctions peuvent être récursives, mais pas mutuellement récursives : le corps d'une fonction  $f$  ne peut appeler que des fonctions définies préalablement ou la fonction  $f$  elle-même. Les fonctions ne renvoient pas de valeur. Toutes les variables locales sont déclarées en début de fonction. On rappelle que la construction `if` du langage C teste si son argument est non nul. Ici, on ne teste que des valeurs entières (il n'y a pas de pointeur nul et donc aucune raison de tester un pointeur). Voici un exemple de programme dans ce fragment :

```
void loop(int *a, int *b, int n) {
    if (n) { *b = *b - (0 - *a); *a = *b - *a; loop(a, b, n - 1); }
}
void fib(int *x, int n) {
    int b; *x = 0; b = 1; loop(x, &b, n);
}
```

La fonction `fib` calcule le  $n$ -ième nombre de Fibonacci et l'écrit dans `*x`. Pour cela, elle alloue une variable locale `b`, dont elle passe l'adresse à la fonction `loop`. On prendra le temps de bien lire et de bien comprendre cet exemple.

### 1.1 Sémantique

**Question 1** Donner le code d'une fonction `void mul(int *x, int y)` qui a pour effet de multiplier `*x` par `y`, c'est-à-dire le même effet que l'affectation `*x = *x * y` si on avait une multiplication à notre disposition. On suppose que la valeur de `y` est positive ou nulle.

---

**Correction :**

$e ::=$	<ul style="list-style-type: none"> <li>  <math>n</math></li> <li>  <math>x</math></li> <li>  <math>e - e</math></li> <li>  <math>*e</math></li> <li>  <math>\&amp;e</math></li> </ul>	<p><b>expression</b></p> <ul style="list-style-type: none"> <li>constante entière</li> <li>variable</li> <li>soustraction</li> <li>déréférencement de pointeur</li> <li>prise d'adresse</li> </ul>
$s ::=$	<ul style="list-style-type: none"> <li>  <math>e = e</math></li> <li>  <math>f(e, \dots, e)</math></li> <li>  <b>if</b> (<math>e</math>) <math>s</math> <b>else</b> <math>s</math></li> <li>  <math>\{ s \dots s \}</math></li> </ul>	<p><b>instruction</b></p> <ul style="list-style-type: none"> <li>affectation</li> <li>appel de fonction</li> <li>conditionnelle</li> <li>bloc</li> </ul>
$d ::=$	<ul style="list-style-type: none"> <li>  <b>void</b> <math>f(\tau x, \dots, \tau x)</math></li> <li>  <math>\{ \tau x; \dots \tau x; s \dots s \}</math></li> </ul>	<p><b>définition de fonction</b></p>
$\tau ::=$	<ul style="list-style-type: none"> <li>  <b>int</b></li> <li>  <math>\tau^*</math></li> </ul>	<p><b>type</b></p> <ul style="list-style-type: none"> <li>entier</li> <li>pointeur</li> </ul>
$p ::=$	<ul style="list-style-type: none"> <li>  <math>d \dots d</math></li> </ul>	<p><b>programme</b></p>

FIGURE 1 – Syntaxe abstraite.

---

```

void mul(int *x, int y) {
  int z;
  if (y) {
    z = *x;
    mul(x, y-1);
    *x = *x - (0 - z);
  } else
    *x = 0;
}

```

---

**Sémantique opérationnelle des expressions.** On souhaite munir notre langage d'une sémantique opérationnelle à grands pas sous la forme d'un premier jugement

$$E, M, e \rightarrow v$$

qui s'interprète comme « dans l'environnement  $E$  et la mémoire  $M$ , l'évaluation de l'expression  $e$  termine, avec la valeur  $v$  ». On suppose que toutes les variables sont allouées en mémoire et  $E$  est une fonction donnant l'emplacement de chaque variable dans la mémoire. Une valeur  $v$  peut être de

deux sortes :

$$v ::= \begin{array}{l} \mathbf{valeur} \\ | \quad n \text{ constante entière} \\ | \quad a \text{ adresse mémoire} \end{array}$$

Une valeur de la forme  $a$  désigne une adresse mémoire et la valeur située à cette adresse est donc  $M(a)$ .

**Question 2** Donner les règles d'inférence définissant la relation  $E, M, e \rightarrow v$ .

---

**Correction :**

$$\frac{}{E, M, n \rightarrow n} \quad \frac{x \in \text{dom}(E)}{E, M, x \rightarrow M(E(x))} \quad \frac{E, M, e \rightarrow a \quad a \in \text{dom}(M)}{E, M, *e \rightarrow M(a)}$$

$$\frac{E, M, e_1 \rightarrow n_1 \quad E, M, e_2 \rightarrow n_2}{E, M, e_1 - e_2 \rightarrow n_1 - n_2} \quad \frac{x \in \text{dom}(E)}{E, M, \&x \rightarrow E(x)} \quad \frac{E, e \rightarrow a}{E, M, \&*e \rightarrow a}$$


---

**Question 3** Donner un environnement  $E$ , une mémoire  $M$  et une expression  $e$  pour lesquels il n'existe pas de valeur  $v$  telle que  $E, M, e \rightarrow v$ .

---

**Correction :** Il y a plein de solutions, dont

- une variable qui n'est pas dans le domaine de  $E$ , c'est-à-dire  $\emptyset, \emptyset, x \not\rightarrow v$  ;
  - le déréférencement d'un pointeur qui n'est pas dans le domaine de  $M$ , c'est-à-dire  $\{x \mapsto a\}, \emptyset, *x \not\rightarrow v$  ;
  - la prise d'adresse d'une expression qui n'est pas une valeur gauche, comme  $\&(x-y)$  ;
  - la soustraction de deux expressions qui ne sont pas des entiers, comme  $\&x - y$ .
-

**Sémantique opérationnelle des instructions.** Pour les instructions, on se donne un second jugement

$$E, M, s \rightarrow M'$$

qui se lit comme « dans l'environnement  $E$  et la mémoire  $M$ , l'évaluation de l'instruction  $s$  termine, en aboutissant à une mémoire finale  $M'$  ». On donne la règle pour un appel de fonction :

$$\frac{\begin{array}{c} \text{void } f(\tau_1 x_1, \dots, \tau_n x_n) \{ \tau'_1 y_1; \dots; \tau'_m y_m; s_1; \dots; s_k \} \\ E, M, e_i \rightarrow v_i \\ E' \stackrel{\text{def}}{=} \{x_1 \mapsto a_1, \dots, x_n \mapsto a_n, y_1 \mapsto a'_1, \dots, y_m \mapsto a'_m\} \quad a_j, a'_j \text{ fraîches} \\ M' \stackrel{\text{def}}{=} M + \{a_j \mapsto v_j \mid 1 \leq j \leq n\} + \{a'_j \mapsto v'_j \mid 1 \leq j \leq m\} \\ E', M', \{s_1; \dots; s_k\} \rightarrow M'' \end{array}}{E, M, f(e_1, \dots, e_n) \rightarrow M''|_M}$$

Ici, les valeurs  $v'_j$  désignent des valeurs quelconques utilisées pour initialiser les variables locales  $y_1, \dots, y_m$  de la fonction. La notation  $M''|_M$  désigne la fonction  $M''$  restreinte au domaine de la fonction  $M$ . On exprime ainsi que les variables locales de la fonction  $f$  ne sont plus visibles après l'appel à  $f$ .

**Question 4** Donner les autres règles d'inférence définissant la relation  $E, M, s \rightarrow M'$ . (Qu'on se rassure, c'est *beaucoup* plus simple que pour la règle précédente.)

---

**Correction :**

$$\frac{x \in \text{dom}(E) \quad E, M, e \rightarrow v}{E, M, x = e \rightarrow M[E(a) \leftarrow v]} \quad \frac{E, M, e_1 \rightarrow a \quad E, M, e_2 \rightarrow v}{E, M, *e_1 = e_2 \rightarrow M[a \leftarrow v]}$$

$$\frac{E, M, e \rightarrow n \quad n \neq 0 \quad E, M, s_1 \rightarrow M'}{E, M, \text{if } (e) s_1 \text{ else } s_2 \rightarrow M'} \quad \frac{E, M, e \rightarrow 0 \quad E, M, s_2 \rightarrow M'}{E, M, \text{if } (e) s_1 \text{ else } s_2 \rightarrow M'}$$

$$\frac{}{E, M, \{ \} \rightarrow M} \quad \frac{E, M, s_1 \rightarrow M_1 \quad E, M_1, \{s_2; \dots; s_n\} \rightarrow M'}{E, M, \{s_1; s_2; \dots; s_n\} \rightarrow M'}$$


---

## 1.2 Analyse syntaxique

On souhaite réaliser l'analyse syntaxique de notre petit langage avec un outil de type YACC (par exemple CUP pour Java ou Menhir pour OCaml). Pour les expressions, on écrit notamment un morceau de grammaire de la forme suivante (la syntaxe précise peut légèrement varier selon l'outil), où **expr** est l'unique non terminal :

```

expr :
  | IDENT          { ... }
  | CONST          { ... }
  | expr MINUS expr { ... }
  | STAR expr     { ... }
  | AMP expr     { ... }
  | LPAR expr RPAR { ... }

```

Les actions sémantiques ne nous intéressent pas ici et sont omises ( $\{ \dots \}$ ).

**Question 5** Lorsque l'outil YACC est lancé sur le fichier ci-dessus, il déclare trois conflits de type lecture/réduction (**shift/reduce**). Les identifier, les expliquer, faire le choix de favoriser lecture ou réduction et expliquer comment indiquer ce choix à l'outil YACC. Les règles de grammaire ne doivent pas être modifiées.

---

**Correction :** Le premier conflit provient d'une expression telle que

$$e_1 - e_2 - e_3$$

où on a le choix entre réduire  $((e_1 - e_2) - e_3)$  ou lire  $(e_1 - (e_2 - e_3))$ . Le deuxième conflit provient d'une expression telle que

$$*e_1 - e_2$$

où on a le choix entre réduire  $((*e_1) - e_2)$  ou lire  $(*e_1 - e_2)$ . Enfin, le troisième conflit provient d'une expression telle que

$$\&e_1 - e_2$$

où on a le choix entre réduire  $((\&e_1) - e_2)$  ou lire  $(\&e_1 - e_2)$ .

Pour ce qui est du premier conflit, on choisit naturellement une associativité à gauche, c'est-à-dire qu'on choisit de privilégier la réduction. Pour ce qui est des deux autres conflits, il est naturel de donner à  $*$  et  $\&$  une priorité plus forte que celle de la soustraction. On déclare donc ceci :

```
%left MINUS
%nonassoc STAR AMP
```

On note qu'il n'y a pas de conflit entre  $*$  et  $\&$  et donc aucune raison ici de leur donner des priorités différentes. (Dans une grammaire plus réaliste où  $*$  serait utilisé également pour la multiplication, il serait alors nécessaire de donner des priorités différentes, mais ce n'est pas le cas ici.)

---

### 1.3 Typage

On se propose maintenant de réaliser le typage statique de notre langage, notamment sous la forme d'un jugement de typage

$$\Gamma \vdash e : \tau$$

qui signifie « dans l'environnement  $\Gamma$ , l'expression  $e$  est bien typée de type  $\tau$  ». L'environnement  $\Gamma$  est l'ensemble des variables connues, avec pour chacune son type.

**Question 6** Donner les règles d'inférence définissant le jugement  $\Gamma \vdash e : \tau$ . Attention à bien limiter l'utilisation de l'opérateur  $\&$  à des valeurs gauches.

---

**Correction :**

$$\frac{}{\Gamma \vdash n : \text{int}} \quad \frac{x : \tau \in \Gamma}{\Gamma \vdash x : \tau} \quad \frac{\Gamma \vdash e_1 : \tau*}{\Gamma \vdash *e_1 : \tau}$$

$$\frac{\Gamma \vdash e_1 : \text{int} \quad \Gamma \vdash e_2 : \text{int}}{\Gamma \vdash e_1 - e_2 : \text{int}} \quad \frac{\Gamma \vdash x : \tau}{\Gamma \vdash \&x : \tau*} \quad \frac{\Gamma \vdash e_1 : \tau*}{\Gamma \vdash \&*e_1 : \tau*}$$


---

**Question 7** Donner un exemple d'expression mal typée ne contenant aucune variable.

---

**Correction :** C'est facile : `*0`. (C'est effectivement une expression mal typée en C.)

---

**Question 8** A-t-on la propriété de sûreté de typage, c'est-à-dire qu'une expression bien typée s'évalue toujours en une valeur? Si oui, le démontrer. Sinon, donner un contre-exemple.

---

**Correction :** On n'a pas la propriété de sûreté du typage. En effet, si on a reçu un pointeur `x` en argument, par exemple de type `int*`, alors on peut le déréréferencer, c'est-à-dire écrire `*x` et cela est bien typé. Mais il pourrait s'agir d'un pointeur fantôme (on verra comment dans Q14) et on n'obtiendra alors pas de valeur (mais un plantage du programme).

---

**Typage des instructions.** Pour typer les instructions, on se donne un second jugement de typage, de la forme

$$\Delta, \Gamma \vdash s$$

qui signifie « dans les environnements  $\Delta$  et  $\Gamma$ , l'instruction  $s$  est bien typée ». L'environnement  $\Delta$  est la liste des fonctions connues, avec pour chacune son profil.

**Question 9** Donner les règles d'inférence définissant le jugement  $\Delta, \Gamma \vdash s$ .

---

**Correction :**

$$\frac{x : \tau \in \Gamma \quad \Gamma \vdash e : \tau}{\Delta, \Gamma \vdash x = e_2} \quad \frac{\Gamma \vdash e_1 : \tau^* \quad \Gamma \vdash e_2 : \tau}{\Delta, \Gamma \vdash *e_1 = e_2}$$

$$\frac{f(\tau_1 x_1, \dots, \tau_n x_n) \in \Delta \quad \forall i. \Gamma \vdash e_i : \tau_i}{\Delta, \Gamma \vdash f(e_1, \dots, e_n)} \quad \frac{\forall i. \Delta, \Gamma \vdash s_i}{\Delta, \Gamma \vdash \{s_1; \dots; s_n\}}$$

$$\frac{\Gamma \vdash e : \text{int} \quad \Delta, \Gamma \vdash s_1 \quad \Delta, \Gamma \vdash s_2}{\Delta, \Gamma \vdash \text{if } (e) s_1 \text{ else } s_2}$$


---

**Question 10** Donner enfin la règle permettant de vérifier qu'une définition de fonction est bien typée.

---

**Correction :** Pour qu'une définition de fonction

$$\text{void } f(\tau_1 x_1, \dots, \tau_n x_n) \{ \tau'_1 y_1; \dots; \tau'_m y_m; s_1; \dots; s_k \}$$

soit bien typée, il faut vérifier que le bloc  $\{s_1; \dots; s_k\}$  est bien typé dans l'ensemble  $\Gamma$  défini comme

$$\Gamma \stackrel{\text{def}}{=} \{ \tau_1 x_1, \dots, \tau_n x_n, \tau'_1 y_1, \dots, \tau'_m y_m \}$$

c'est-à-dire réunissant toutes les variables de la fonction, et en ajoutant la fonction  $f$  à  $\Delta$  pour permettre une définition récursive.

---

## 1.4 Compilation vers x86-64

On se propose maintenant de compiler notre langage vers l'assembleur x86-64. (Un aide-mémoire est donné en annexe.) On adopte le schéma de compilation suivant. Toutes les valeurs sont sur 64 bits (qu'il s'agisse d'entiers ou de pointeurs). Les variables sont allouées et passées dans des registres, lorsque c'est possible, et sur la pile sinon. On note qu'une variable  $x$  dont l'adresse est prise avec  $\&x$  doit nécessairement être allouée sur la pile. Le tableau d'activation prend donc la forme ci-contre. L'adresse de retour est celle déposée par l'instruction `call`. Les valeurs de  $n'$  et  $m'$  dépendent du nombre de variables qu'on a dû allouer sur la pile. On suppose qu'on est ici dans le contexte d'un compilateur optimisant suivant l'architecture vue en cours (sélection d'instruction, RTL, ERTL, LTL, linéarisation).

	argument $n'$
	⋮
	argument 1
	adresse de retour
$\%rbp \rightarrow$	sauvegarde $\%rbp$
	locale 1
	⋮
	locale $m'$

**Question 11** Donner un code x86-64 possible pour les fonctions `loop` et `fib` de la page 1, en optimisant l'appel terminal à `loop`.

---

**Correction :** Pour `fib`, les variables `x` et `n` sont dans  $\%rdi$  et  $\%rsi$  et la variable `b` est allouée sur la pile.

```
fib:   mov    %rsi, %rdx
       movq  $0, (%rdi)
       pushq $1
       mov  %rsp, %rsi      # b
       call loop
       addq $8, %rsp
       ret
```

Pour `loop`, les variables `a`, `b` et `n` sont dans  $\%rdi$ ,  $\%rsi$  et  $\%rdx$ .

```
loop:  testq %rdx, %rdx      # if(n)
       jz   2f
       movq (%rdi), %rcx
       addq %rcx, (%rsi)
       movq (%rsi), %r8
       subq %rcx, %r8
       movq %r8, (%rdi)
       dec  %rdx
       jmp  loop        # appel terminal
2:     ret
```

---

**Question 12** Indiquer à quelle(s) étape(s) du compilateur optimisant on va assurer qu'une variable  $x$  dont l'adresse est prise avec  $\&x$  est bien allouée sur la pile.

---

**Correction :** Si on repère  $\&x$  dans le code, alors on note que la variable `x` devra être allouée sur la pile. On peut le faire très tôt, par exemple pendant le `typage` ou juste après.

Ensuite, on alloue pour  $x$  un pseudo-registre, comme pour toute autre variable, et rien ne change jusqu'à l'allocation de registres. Au moment de la coloration, on force le vidage en mémoire des variables ainsi repérées. Pour autant, elles participent tout de même de l'allocation, notamment pour diminuer la valeur de  $m'$ .

---

**Question 13** Expliquer comment compiler la construction  $\&$ , en détaillant ce qui est fait aux différentes étapes du compilateur optimisant.

---

**Correction :** Pour la compilation de la construction  $\&$ , on distingue deux cas sur la valeur gauche qui constitue l'opérande :

— pour  $\&x$ , on peut introduire une nouvelle instruction RTL

```
addr #i r -> L
```

où  $\#i$  est le pseudo-registre associé à la variable  $x$ . Même chose en ERTL. Pendant la coloration, le vidage en mémoire de  $\#i$  est forcé (cf question précédente), et on obtient donc un emplacement de pile  $ofs$  à l'issue de la coloration. En LTL, on a donc

```
leaq ofs(%rbp), r
```

si  $r$  est un registre physique et

```
leaq ofs(%rbp), tmp1
```

```
mov tmp1, r
```

si  $r$  est alloué en pile.

— pour  $\&*e$ , on compile tout simplement l'expression  $e$ , à l'étape RTL.

---

**Question 14** Montrer que, bien que les fonctions ne renvoient pas de valeur, on a tout de même le problème des références fantômes (*dangling reference*) dans ce fragment de C.

---

**Correction :** Il est possible de récupérer une référence fantôme, en la stockant dans un emplacement mémoire passé sous la forme d'un argument de type `int**`, de la manière suivante :

```
void dangling(int **p, int y) {
    int z; ...
    *p = &z;
}
```

Maintenant, il suffit de récupérer la référence fantôme dans une variable de type `int*`, pour ensuite s'en servir pour faire planter le programme.

```
void ous() {
    int *p;
    dangling(&p);
    ... on plante le programme en écrivant dans *p ...
}
```

Il faut alors un peu d'imagination et de tâtonnements pour que l'écriture dans `*p` provoque un plantage. Une solution consiste à ce qu'une affectation comme `*p = 100` écrase justement la valeur de `p` avec 100, puis on tente d'accéder à l'adresse 100 avec `*p` et le programme plante.

---



#2 $f(\#1)$		
entrée $L_1$ sortie $L_{16}$		
$L_1$ : mov 0 #3	$\rightarrow L_2$	$L_8$ : binop sub #7 #6 $\rightarrow L_9$
$L_2$ : mov #1 #4	$\rightarrow L_3$	$L_9$ : binop add #6 #3 $\rightarrow L_{10}$
$L_3$ : mov 1 #5	$\rightarrow L_4$	$L_{10}$ : mov #1 #8 $\rightarrow L_{11}$
$L_4$ : binop sub #5 #4	$\rightarrow L_5$	$L_{11}$ : mov 1 #9 $\rightarrow L_{12}$
$L_5$ : ubranch >0 #4	$\rightarrow L_6, L_{15}$	$L_{12}$ : binop sub #9 #8 $\rightarrow L_{13}$
$L_6$ : mov #1 #6	$\rightarrow L_7$	$L_{13}$ : mov #8 #1 $\rightarrow L_{14}$
$L_7$ : mov 1 #7	$\rightarrow L_8$	$L_{14}$ : goto $\rightarrow L_2$
		$L_{15}$ : mov #3 #2 $\rightarrow L_{16}$

FIGURE 2 – Code RTL d’une fonction  $f$ .

## 2 Analyse de flot de données

Le contexte de ce problème est celui du langage RTL (*Register Transfer Language*). On va y réaliser une analyse statique, dans le but de réaliser des optimisations. On rappelle ici les différentes instructions du langage RTL :

$i ::=$	mov $n$ $r \rightarrow L$	chargement d’une constante
	mov $r$ $r \rightarrow L$	copie
	load $n(r)$ $r \rightarrow L$	lecture en mémoire
	store $r$ $n(r) \rightarrow L$	écriture en mémoire
	unop $op$ $r \rightarrow L$	opération unaire
	binop $op$ $r$ $r \rightarrow L$	opération binaire
	ubranch $br$ $r \rightarrow L, L$	branchement unaire
	bbranch $br$ $r$ $r \rightarrow L, L$	branchement binaire
	call $r \leftarrow f(r, \dots, r) \rightarrow L$	appel de fonction
	goto $\rightarrow L$	saut incondtionnel

Les conventions sont les suivantes :  $n$  désigne une constante entière,  $r$  un pseudo-registre et  $L$  une étiquette de code (à laquelle on trouve une instruction RTL). La figure 2 contient le code RTL d’une fonction  $f$  qui reçoit un argument dans le pseudo-registre #1 et renvoie un résultat dans le pseudo-registre #2. Le point d’entrée est l’étiquette  $L_1$  et le point de sortie l’étiquette  $L_{16}$ .

**Question 15** Donner un programme C possible pour le code RTL de la figure 2. (Note : On n’est plus dans le fragment C du problème 1.)

---

**Correction :**

```
int f(int x) {
    int s = 0;
    while (x-1 > 0) {
        s += x-1;
        x = x-1;
    }
    return s;
}
```

---

**Expressions disponibles.** Le but de notre analyse est de calculer, en chaque point du code RTL d'une fonction donnée, un ensemble d'*expressions disponibles*, c'est-à-dire de valeurs qui ont déjà été calculées et sont contenues dans des pseudo-registres. On représente de telles valeurs *symboliquement*, à l'aide de la syntaxe abstraite suivante :

$$\begin{array}{l|l}
 v ::= & \alpha_i \quad \text{valeur arbitraire, inconnue} \\
 & n \quad \text{constante entière} \\
 & op \ v \quad \text{opération unaire} \\
 & v \ op \ v \quad \text{opération binaire} \\
 & v[n] \quad \text{accès en mémoire à l'adresse } v + n
 \end{array}$$

En particulier, la construction  $\alpha_i$  nous permet de représenter une valeur arbitraire inconnue de l'analyse statique. C'est le cas notamment des arguments de la fonction. Ainsi, sur l'exemple de la figure 2, on donnera initialement au pseudo-registre #1 une valeur arbitraire  $\alpha_1$ .

Plus précisément, on va calculer pour chaque étiquette  $L$  du graphe de flot de contrôle un ensemble  $in(L)$  d'expressions disponibles avant l'exécution de l'instruction à l'étiquette  $L$  et un ensemble  $out(L)$  d'expressions disponibles après son exécution. Sur l'exemple de l'instruction  $L_4$  de la figure 2, on a les ensembles suivants

$$\begin{aligned}
 in(L_4) &= \{\alpha_2, \alpha_3, 1\} \\
 L_4 : \text{binop sub \#5 \#4} &\rightarrow L_5 \\
 out(L_4) &= \{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}
 \end{aligned}$$

où  $\alpha_2$  désigne la valeur contenue dans #1 et #4, et  $\alpha_3$  celle contenue dans #3.

**Question 16** Donner des ensembles  $in(L)$  et  $out(L)$  possibles pour toutes les instructions de la figure 2 sous la forme d'un tableau :

étiquette $L$	$in(L)$	$out(L)$
$L_1$	$\{\alpha_1\}$	$\{\alpha_1, 0\}$
etc.		

Comme le code contient une boucle, on s'autorisera des approximations, en introduisant des valeurs arbitraires  $\alpha_i$  aux endroits pertinents.

---

**Correction :**

étiquette $L$	$in(L)$	$out(L)$
$L_1$	$\{\alpha_1\}$	$\{\alpha_1, 0\}$
$L_2$	$\{\alpha_2, \alpha_3\}$	$\{\alpha_2, \alpha_3\}$
$L_3$	$\{\alpha_2, \alpha_3\}$	$\{\alpha_2, \alpha_3, 1\}$
$L_4$	$\{\alpha_2, \alpha_3, 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$
$L_5$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$
$L_6$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$
$L_7$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$
$L_8$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$
$L_9$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{10}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{11}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{12}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{13}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{14}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$	$\{\alpha_2, \alpha_3, 1, \alpha_2 - 1, \alpha_3 + \alpha_2 - 1\}$
$L_{15}$	$\{\alpha_4\}$	$\{\alpha_4\}$
$L_{16}$	$\{\alpha_4\}$	—

**Calcul des expressions disponibles.** Pour calculer les ensembles  $in(L)$  et  $out(L)$  de façon systématique, on va commencer par calculer les valeurs respectivement définies et invalidées par une instruction RTL. On notera respectivement  $gen(L)$  et  $kill(L)$  ces deux ensembles pour l'instruction située à l'étiquette  $L$ .

**Question 17** Définir les ensembles  $gen(L)$  et  $kill(L)$  pour chaque instruction du langage RTL, en complétant le tableau suivant :

Instruction RTL	$gen(L)$	$kill(L)$
<b>binop</b> $op\ r_1\ r_2$	$\{e(r_1)\ op\ e(r_2)\}$	toute expression associée à $r_2$
<i>etc.</i>		

où  $e(r)$  dénote la valeur de  $in(L)$  contenue dans  $r$ , le cas échéant, et une nouvelle valeur  $\alpha_i$  sinon.

**Correction :**

Instruction RTL	$gen(L)$	$kill(L)$
<b>mov</b> $n\ r_1$	$\{n\}$	toute expression associée à $r_1$
<b>mov</b> $r_2\ r_1$	$\{e(r_2)\}$	toute expression associée à $r_1$
<b>load</b> $n(r_2)\ r_1$	$\{e(r_2)[n]\}$	toute expression associée à $r_1$
<b>store</b> $r_2\ n(r_1)$	$\emptyset$	toute expression contenant un $n'(r_1)$
<b>unop</b> $op\ r_1$	$\{op\ e(r_1)\}$	toute expression associée à $r_1$
<b>binop</b> $op\ r_2\ r_1$	$\{e(r_1)\ op\ e(r_2)\}$	toute expression associée à $r_1$
<b>ubbranch</b> $br\ r_1$	$\emptyset$	$\emptyset$
<b>bbranch</b> $op\ r_1\ r_2$	$\emptyset$	$\emptyset$
<b>call</b> $r_0 \leftarrow f(r_1, \dots, r_n)$	$\{\alpha_i\}$ frais	toute expression associée à $r_0$
<b>goto</b>	$\emptyset$	$\emptyset$

**Question 18** Donner des équations définissant  $in(L)$  et  $out(L)$

$$\begin{cases} in(L) = \dots \\ out(L) = \dots \end{cases}$$

en fonction de  $gen(L)$ ,  $kill(L)$ ,  $in(L)$  et  $out(L)$ .

---

**Correction :**

$$\begin{cases} in(L) = \bigcap_{p \in pred(L)} out(L) \\ out(L) = gen(L) \cup (in(L) - kill(L)) \end{cases}$$

L'intersection traduit le fait qu'une expression est disponible en  $L$  si elle est disponible sur *tous* les chemins menant à  $L$ .

---

**Question 19** Donner les types et fonctions OCaml ou classes et méthodes Java que vous écririez pour réaliser une telle analyse, en supposant donnés des types/classes pour les instructions RTL, pour les pseudo-registres et pour les étiquettes de code (comme dans le projet). On ne demande pas d'écrire le code, mais seulement de décrire son architecture.

---

**Correction :** On le fait ici en OCaml.

On introduit bien sûr un type `expression` pour représenter les expressions disponibles :

```
type expression = Alpha of int | Const of int | ...
```

La difficulté, ensuite, vient du fait qu'on ne peut pas se contenter de calculer les expressions disponibles. Il faut aussi conserver l'information sur les pseudo-registres qui les contiennent à chaque instant. Là, plusieurs choix sont possibles, selon la finesse de l'analyse. On peut par exemple indiquer pour chaque expression disponible un *ensemble* de pseudo-registres la contenant.

```
type available_expressions = Register.Set.t Expression.Map.t
```

en supposant un module `Register.Set` pour des ensembles de pseudo-registres et un module `Expression.Map` pour des dictionnaires indexés par des expressions.

On peut alors des fonctions `gen` et `kill` avec le type

```
val gen, kill: available_expressions -> Rtl.t -> available_expressions
```

puis des fonctions ensemblistes pour réaliser les opérations utilisées dans la question précédente :

```
val inter, diff, union: available_expressions ->  
    available_expressions ->  
    available_expressions
```

Reste alors la fonction d'analyse proprement dite, qui prend un graphe de flot de contrôle en argument et renvoie le résultat de l'analyse :

```
type control_flow_graph = Rtl.t Label.Map.t
```

```
type result = { in: available_expressions Label.Map.t;  
    out: available_expressions Label.Map.t; }
```

```
val available_expressions: control_flow_graph -> result
```

C'est dans cette dernière fonction que se cache un calcul de point fixe (qui n'était pas demandé ici).

---

**Sous-expressions communes.** On souhaite maintenant exploiter le résultat de l'analyse des expressions disponibles pour effectuer l'optimisation dite des *sous-expressions communes*, qui consiste à éviter de faire plusieurs fois les mêmes calculs.

**Question 20** Identifier les calculs redondants sur l'exemple de la figure 2, en montrant notamment comment l'analyse des expressions disponibles a permis de les trouver. Donner un code RTL simplifié en conséquence pour la fonction  $f$ .

---

**Correction :** On calcule trois fois de suite  $\#1 - 1$  pour la même valeur de  $\#1$ . Les deux derniers calculs peuvent être supprimés. L'analyse des expressions disponibles permet de l'identifier, car les instructions  $L_8$  et  $L_{12}$  calculent  $\alpha_2 - 1$  avec la valeur  $\alpha_2 - 1$  disponible, en l'occurrence dans  $\#4$ . Il suffit donc de remplacer ces deux instructions par un `mov` à chaque fois. Du coup, les instructions  $L_7$ ,  $L_8$ ,  $L_{11}$  et  $L_{12}$  deviennent du code mort (ce que l'analyse de durée de vie, mise à jour, montre facilement) et peuvent être éliminées. Au final, on obtient

```
#2 f(#1)
   entrée L1 sortie L17
L1 : mov 0 #3          L2
L2 : mov #1 #4         L3
L3 : mov 1 #5          L4
L4 : binop sub #5 #4   L5
L5 : ubranch >0 #4    L9, L15
L9 : binop add #4 #3   L13
L13: mov #4 #1         L14
L14: goto              L2
L15: mov #3 #2         L16
```

---

**Question 21** Donner un exemple de code RTL où une même expression est disponible en un certain point du programme, mais contenue dans deux pseudo-registres différents selon le chemin du graphe de flot de contrôle qui a permis d'atteindre ce point de programme.

---

**Correction :** Dans un langage de haut niveau, ce pourrait être

```
if (...) y = x-1; else z = x-1;
```

Après cette conditionnelle, l'expression  $x-1$  est disponible, mais dans deux pseudo-registres différents selon la branche de la conditionnelle qui a été suivie. Traduit en RTL, cela pourrait être

```
L1: ubranch =0 #2    --> L2, L6
L2: mov #1 #3        --> L3
L3: mov 1 #4          --> L4
L4: binop sub #4 #3  --> L5
L5: goto              --> L9
```

```
L6: mov #1 #5      --> L7
L7: mov 1 #6       --> L8
L8: binop sub #6 #5 --> L9
L9:
```

et le point qui nous intéresse ici est  $L_9$ .

---

**Question 22** Décrire de manière générale comment l'analyse des expressions disponibles doit être utilisée pour réaliser l'optimisation des sous-expressions communes. On prendra soin d'expliquer comment le code mort résultant de cette optimisation est éliminé.

---

**Correction :** On a déjà en partie répondu dans la question 20 : si une instruction calcule une expression qui est disponible et que celle-ci se trouve dans un pseudo-registre, il suffit de remplacer le calcul par un `mov`. Mais la question précédente montre qu'une expression peut être disponible sans se trouver pour autant dans un registre donné. Il peut y avoir par exemple deux branches dans le graphe de flot de contrôle qui précède, l'une mettant la valeur de l'expression dans  $r_1$  et l'autre dans  $r_2$ . Dans ce cas, on choisit un nouveau pseudo-registre  $r$ , dans chaque branche on y copie la valeur de l'expression, puis enfin on réalise l'optimisation avec un `mov` prenant la valeur dans  $r$ .

Le code mort résultant de l'optimisation peut être éliminé facilement en mettant à jour l'analyse de durée de vie. En effet, on peut éliminer toute instruction qui affecte un pseudo-registre qui n'est plus vivant après cette instruction. Dans notre exemple, l'instruction  $L_6$  affecte le pseudo-registre `#6` qui n'est plus vivant après cette instruction suite à l'optimisation. De même pour les instructions  $L_7$ ,  $L_{10}$  et  $L_{11}$ .

---

## Annexe : aide-mémoire x86-64

Dans ce qui suit,  $L$  désigne une étiquette,  $r$  un registre,  $n$  une constante entière et  $o$  une opérande qui est soit un registre  $r$ , soit une opérande indirecte  $n(r)$ .

<code>mov <math>r, o</math></code>	copie le registre $r$ dans l'opérande $o$
<code>mov <math>o, r</math></code>	copie l'opérande $o$ dans le registre $r$
<code>mov <math>\\$n, r</math></code>	charge la constante $n$ dans le registre $r$
<code>add <math>o, r</math></code>	calcule $r + o$ et l'affecte à $r$
<code>add <math>r, o</math></code>	calcule $o + r$ et l'affecte à $o$
<code>add <math>\\$n, o</math></code>	calcule $o + n$ et l'affecte à $o$
<code>shl <math>\\$n, r</math></code>	décale les bits de $r$ vers la gauche, $n$ fois
<code>lea <math>n(r_1), r_2</math></code>	affecte à $r_2$ l'adresse représentée par $n(r_1)$ , c'est-à-dire $r_1 + n$
<code>push <math>r</math></code>	empile la valeur de $r$
<code>pop <math>r</math></code>	dépile une valeur dans le registre $r$
<code>test <math>r_2, r_1</math></code>	positionne les drapeaux en fonction de la valeur de $r_1$ ET $r_2$
<code>jz <math>L</math></code>	saute à l'adresse désignée par l'étiquette $L$ si les drapeaux signalent un résultat nul
<code>jmp <math>L</math></code>	saute à l'adresse désignée par l'étiquette $L$
<code>call <math>L</math></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 pourra utiliser librement toute autre instruction x86-64. Mais les instructions ci-dessus suffisent pour répondre aux questions de ce sujet.