

Compilateur Léa – Projet à rendre le 21 décembre 2012

Le but de ce projet est de livrer à terme un compilateur pour un langage élémentaire algorithmique (Léa) que nous illustrons en Fig.2. Ce langage est destiné à enseigner la programmation dans les cursus de sciences humaines. Nous nous inspirerons des langages Pascal, Python, Ruby etc. avec cette exigence : les connaissances en programmation avancées ne sont pas exigées pour savoir programmer avec Léa, et les habitudes récentes de programmation objet seront conservées.

Le projet Léa est surtout un prétexte pour proposer un travail de compilation complet. Il ne vous sera d'aucune façon reproché dans ce travail

1. De ne pas avoir livré un compilateur exploitable dans l'état
2. D'avoir été associé à une équipe faible

Il est demandé à chacun de traiter les sections écrites en noir en s'aidant des explications données en TD. Il n'est pas exigé de traiter l'ensemble du projet, on pourra en faire une partie seulement.

Les parties en **vert** et en **bleu** sont de difficulté croissante, vous pouvez tenter de les faire.

## 1 Modalités pratiques

Ce projet est à réaliser en groupe de 4 personnes, pour le xx décembre 2012 au plus tard. Chaque groupe est composé de :

1. Un(e) correspondant(e)  
C'est le porte-parole du projet, celui qui écrit sa documentation et assure sa cohérence.
2. Un(e) syntacticien(ne)  
Il écrit le code CUP/JFlex.
3. Un(e) sémanticien(ne)  
Il écrit le code intermédiaire
4. Un développeur  
Il écrit le code cible.

Chaque correspondant devra envoyer un courrier électronique au nom du groupe à (**frederique.carrere@labri.fr** ET **lionel.clement@labri.fr**)

- Dont l'objet est "DM Compilation"
- Signé de tous les prénoms et noms du groupe
- Contenant une archive **tar.bz2** avec les fichiers suivants :

1. **dm1.tar**, une archive tar, en excluant tout autre format d'archivage. Cette archive contiendra toutes les sources et toutes les données permettant la compilation et l'exécution du programme complet par simple appel à la commande **ant main**.
2. **Prénom-NomCorrespondant.pdf**, une explication de votre travail, sous la forme d'un texte en format .pdf en excluant tout autre format. Ce texte contiendra :
  - (a) Un cahier des charges précisé
  - (b) Un exposé des problèmes envisagés
  - (c) Une explication des choix retenus

FIG. 1 – Exemple de code Léa

```

1  /* *****
2  * Un exemple complet de code Léa
3  ***** */
4
5  // Définition des constantes
6  MAX = 500;
7
8  // Définition des types
9  fiche = struct {
10     nom: string;
11     age: int;
12 }
13
14 // Définition de la fonction main
15 function main(argn: int, argv: list of string): int
16 {
17     // Définition des variables
18     table: list of fiche;
19     echange: bool;
20
21     // Instructions
22     for i in [0 .. MAX-1] {
23         table[i].nom := "";
24         table[i].age := random(100);
25     }
26
27     repeat {
28         echange := false;
29         for i in [0 .. MAX-2] {
30             if (T[j].age > T[j+1].age){
31                 // Variable locale
32                 aux: fiche;
33                 aux := T[j];
34                 T[j] := T[j+1];
35                 T[j+1] := aux;
36                 echange := true;
37             }
38         }
39     } while (echange);
40     return 0;
41 }

```

## 2 Léa

Le compilateur doit être écrit en CUP/JFlex/Java et doit produire un programme Java équivalent. La gestion de la mémoire doit être entièrement dynamique et les types complexes sont implémentés en utilisant les bibliothèques standards de Java.

## 3 Généralités sur le langage

Le langage est algorithmique procédural. Il utilise un typage statique explicite. C'est-à-dire que toutes les variables sont déclarées avec un type immuable durant l'exécution du code. La syntaxe objet est utilisée pour les fonctions membres de l'objet désigné. Par exemple on peut accéder à la fonction membre `length()` d'une variable `x` de type `string` ainsi : `x.length()`

### 3.1 Structure générale d'un programme Léa

1. Déclaration des modules requis
2. Déclaration des types et constantes
3. Déclaration des signatures des fonctions et procédures
4. Implémentation des fonctions et procédures dont la fonction `main`

### 3.2 Déclaration des constantes, des types et des fonctions

Une constante et un type n'ont qu'une seule déclaration possible pour l'ensemble du programme.

Les procédures et fonctions peuvent avoir des déclarations multiples en cas de surcharge.

Un programme doit contenir une déclaration de la fonction `main(int, list of string): int` qui correspond à l'implémentation de la procédure principale. La signature `procedure main();` est également acceptée (le programme retournera 0 et les arguments seront inaccessibles).

1. Déclaration de constante :

*identifier* = *constant* ;

où *constant* est un littéral ou une constante complexe (expression de liste, d'ensemble, de tuple, etc constante)

2. Déclaration de type construit :

*identifier* = *type* ;

où *type* peut faire référence à des types déjà déclarés

3. Déclaration de procédure :

**procedure** *identifier* (*liste d'arguments*) *bloc*

4. Déclaration de fonction :

**function** *identifier* (*liste d'arguments*): *type bloc*

5. On peut déclarer la signature des fonctions et procédures avant leur implémentation comme on le fait en langage C :

**procedure** *identifier* (*liste d'arguments*) ;  
**function** *identifier* (*liste d'arguments*): *type* ;

### 3.3 Types

Les types suivants sont disponibles ( $T, T_1, T_2, \dots, T_k$  sont des types quelconques) :

- Types élémentaires `char`, `int`, `float`, `bool`, `string`
- Énumérés (vecteur de bits) `enum`  
On accède à l'un des bits par le nom d'un champs
- Listes (implémentation des suites mathématiques) `list of T`
- Tuples ( $T_1, T_2, \dots, T_k$ )
- Structures `struct {slot1 : T1; slot2 : T2; ... slotk : Tk}`  
On accède au champs `sloti` avec l'opérateur `.sloti`
- Ensembles `set of T`
- Dictionnaires `map of (T1, T2)`
- Variables de type. Il conviendra de ne pas confondre avec les noms de type (un identificateur utilisé pour désigner un type), ni les types nommés (un identificateur utilisé comme type distingué pour nommer les slots des structures et des classes). Ces variables de type sont utiles pour écrire les fonctions polymorphes comme dans cet exemple-ci :

```
1 function plus_petit_element(l: list of 'x): 'x {  
2     min: 'x := l.first ();  
3     for i in [0 .. l.len()] // [0 .. l.len()] est une  
4                             liste d'entiers  
5         if (min > l[i])  
6             min := l[i];  
7     return min;  
8 }
```

- Classes  
Les classes contiennent
  - Les variables, procédures et fonctions de classe (notées *static*)
  - Les variables d'instance
  - Les constructeurs d'instance
  - Les opérateurs d'instance

Des opérateurs et fonctions membres sont accessibles pour chaque type. Nous en dressons la liste plus loin.

L'allocation mémoire est implicitement dynamique. C'est-à-dire que les variables sont dynamiquement allouées lors de l'exécution sans utilisation de `new`, `malloc` ou autre instruction d'allocation, de réallocation.

Dans le cas d'une structure réentrante `S`, l'allocation est `null` par défaut pour chaque champs de type `S`. D'une manière générale, la valeur d'une variable non déclarée est `null`, mais cela n'est jamais visible de l'utilisateur qui reçoit un message d'erreur bien avant le calcul du type. Une variable dont le type est marqué comme réentrant est de valeur `null`, sinon en langage cible, elle aura la valeur de l'objet correspondant.

Voici un exemple :

```
1 arbre = struct {
2     left: arbre;
3     right: arbre;
4     etiq: int;
5 }
6
7 function parcours(root: arbre): string {
8     s: string := "";
9     if (root.left != null)
10         s += parcours(root.left)
11     if (root.right != null)
12         s += parcours(root.right)
13     s += root.etiq;
14     return s;
15 }
16
17 procedure main() {
18     a: arbre; // allocation implicite de a
19     a.etiq := 0;
20     a.left.etiq := 1; // allocation implicite de a.left
21
22     writeln (parcours(r));
23 }
```

### 3.4 Instructions et structures de contrôle

#### 3.4.1 Affectation

Pour dissuader l'apprenant de confondre l'égalité avec l'affectation, nous utilisons le signe `:=`

Sont affectables, les variables, les éléments de type complexe (liste, ensemble, etc). Dans le cas de l'affectation d'un élément d'une suite mathématique, l'allocation dynamique (ou la réallocation) doit prévoir  $n^2$  éléments disponibles où  $n$  est l'indice. En code cible java, on ne souciera pas de ce problème d'allocation, on utilisera par exemple immédiatement `ensureCapacity(int)` de `ArrayList` qui désaloue et réalloue dynamiquement.

*AffectableExpr := expr;*

On pourra associer l'affectation et la déclaration

*variable: type := expr;*

#### 3.4.2 Appel de procédure

Si `foo` est une procédure préalablement déclarée, `foo(...)` est accessible comme instruction.

Remarque. Contrairement à de nombreux langages, une expression n'est pas une instruction ni l'inverse. Il n'est donc pas possible d'écrire `fact(7)` comme instruction où `fact` est une fonction ni une affectation comme expression `if (x := foo(i)) ...`

*ProcedureName* (*expr*<sub>1</sub>, *expr*<sub>2</sub>, ..., *expr*<sub>*k*</sub>) ;

### 3.4.3 Entrées sorties

Les procédures `write(string)`, `writeln(string)` permettent un accès aux sorties standard. Ajoutons `read()` comme expression pour un accès à l'entrée standard (une chaîne de caractères terminée par `return`) Le type de ces fonctions et de leur argument est `string`.

### 3.4.4 Conditionnelles et choix multiples

Si (*I*, *I*<sub>1</sub>, *I*<sub>2</sub>, ... *I*<sub>*k*</sub> sont des instructions), les instructions conditionnelles suivantes sont admises :

```
if (expr) I  
if (expr) I1 else I2
```

*expr* doit être du type `bool` contrairement à de nombreux langages

```
case expr of  
    enumslot1 : I1  
    enumslot2 : I2  
    ...  
    enumslotk : Ik
```

*expr* doit être de type `enum`

### 3.4.5 Boucles

Si *I* est une instruction, les boucles suivantes sont admises :

1. boucle `while` et `repeat`

```
while (expr) I  
repeat I while (expr) ;
```

Comme pour les conditionnelles, *expr* doit être du type `bool`.

2. boucle `for`

```
for var in expr I
```

Expr est une expression de type

- Liste
- Ensemble
- Dictionnaire
- Intervalle
- Chaîne de caractères
- Énuméré

La variable *var* est déclarée localement et son type se rapporte au type de *expr* (`char` si *expr* est un `string`, l'un des énumérés pour `enum`, `int` pour `interval`, élément pour la liste l'ensemble ou le dictionnaire).

exemples

```
1 {
2 k: map of (int, string);
3 s: string;
4 t: enum (ROUGE, BLEU, VERT);
5 for i in [1 .. 100] print i;
6 for i in [1, 2, 3] print i;
7 for i in {1, 2, 3} print i;
8 for i in k
9     print(j.first().toString() + " => " + j.second());
10 for i in s
11     print(i);
12 for i in t
13     print(i.toString());
14 }
```

Techniquement, cette instruction `for` est équivalente à

```
1 {
2 uniqIterator: int := 0;
3 variable: typeofelement;
4 while (uniqIterator <= expr.size()) {
5     variable := expr[uniqIterator];
6     I
7     uniqIterator++;
8 }
9 }
```

où *exprIterator* est un itérateur de la suite *expr*

### 3.4.6 Blocs

Un bloc est un ensemble d'instructions précédé d'un ensemble de déclarations de variables. On ne mettra pas de déclaration de types ni de constantes dans les blocs.

## 3.5 Expressions

Les expressions manipulent les variables, les constantes élémentaires et complexes, les opérateurs, les fonctions [et les fonctions membres](#) .

On peut trouver des expressions qui ont des effets comme  $x++$  ou  $foo(y)$  où  $y$  est un argument passé en paramètre. Nous rappelons qu'une expression ne peut être assimilée à une instruction, donc nous écrirons  $x := x + 1$ ; ou  $x+ = 1$ ; au lieu de  $x++$ ;

Les opérateurs et fonctions membres seront repris du langage Java comme traditionnellement. On notera que le signe `=` est réservé pour l'égalité.

Nous dressons dans le tableau suivant les expressions du langage **Léa**

Expression	Type	Valeur
"..."	string	chaîne de caractères où sont échappés les caractères spéciaux <code>\t, \n, \r, \\", \"</code>
'a'	char	caractère où est échappé le caractère spécial <code>\'</code>
<i>nombre entier</i>	int	nombre entier décimal signé de Integer.MIN_VALUE à Integer.MAX_VALUE
<i>nombre flottant</i>	int	nombre flottant décimal signé de Float.MIN_VALUE à Float.MAX_VALUE
<i>True, False</i>	bool	
$E + E$	type des opérandes	addition, concaténation
$E - E$	type des opérandes	soustraction
$E * E$	type des opérandes	multiplication
$E / E$	type des opérandes	division
$E \% E$	type des opérandes	modulo
$- E$	type de l'opérande	$* - 1$
$E < E$	bool	inférieur
$E > E$	bool	supérieur
$E \leq E$	bool	inférieur ou égal
$E \geq E$	bool	supérieur ou égal
$E = E$	bool	égal
$E \neq E$	bool	différent
$E \&\& E$	bool	et
$E \ \  E$	bool	ou
$!E$	bool	non
$(E)$	type de $E$	identité
$(E_1, E_2, \dots, E_k)$	tuple $(type(E_1), type(E_2), \dots, type(E_k))$	construction d'un tuple
$[E_1, E_2, \dots, E_k]$	list of $type(E_i)$	construction d'une liste
$[E_1 \dots E_k]$	list of $type(E)$	construction d'une liste $[E_1, \dots, E_k]$ où $i < j \Rightarrow E_i < E_j$ L'opérateur $<$ est donné pour les types élémentaires, il doit être défini pour les classes
$\{E_1, E_2, \dots, E_k\}$	set of $type(E_i)$	construction d'un ensemble. Les opérateurs $<$ et $=$ sont donnés pour les types élémentaires, ils doivent être définis pour les classes
$\{(key_1, val_1), \dots, (key_k, val_k)\}$ ou $\{key_1 \Rightarrow val_1, \dots, key_k \Rightarrow val_k\}$	map of $(type(key_i), type(val_i))$	construction d'une relation. Les opérateurs $<$ et $=$ sont donnés pour les types de $key_i$ élémentaires, ils doivent être définis pour les classes
$foo(E_1, E_2, \dots, E_k)$	'y où $type(foo) = 'x \rightarrow 'y$	appel d'une fonction $foo$ où ' $x$ s'unifie avec $type(E_1) \times type(E_2) \dots type(E_k)$
$t[E]$	'x où $type(t) = liste\ of\ 'x$	$E^{ieme}$ élément de la liste $t$ ( $E$ doit être de type $int$ )
$t[E]$	'y où $type(t) = map\ of\ ('x, 'y)$	deuxième élément du $E^{ieme}$ élément de la relation $t$ ( $E$ doit être de type $('x, 'y)$ )
<i>E.identififier</i>	type du champs	champs d'une structure, fonction ou variable membre d'une instance de classe ou d'une classe



## 4 Travail des uns et des autres

### 1. Le(a) correspondant(e)

C'est le porte-parole du projet, celui qui écrit sa documentation et assure sa cohérence.

Il est aussi en charge de coordonner le projet et de résoudre les difficultés théoriques qui se trament au fil des ajouts. Il résout les problèmes d'interface entre les différents éléments du projet.

Il devra écrire le fichier `build.xml` et s'assurer de l'architecture complète du projet et de sa compilation. Il devra écrire des programmes tests en langage **Léa**. Il devra être initiateur d'un serveur `svn` pour ses camarades.

Le fichier `build.xml` produira automatiquement des éléments de débogage, des graphiques, des documentations, des tests, etc. Ce fichier permettra de produire un archive automatique et permettra de forcer le nettoyage des fichiers produits par la compilation (`ant clean`)

Il pourra proposer automatiquement la documentation sous forme d'un hypertexte à partir des commentaires de documentation en tête de chaque fonction et chaque classe. Les commentaires de documentation sont marqués avec `/**` et `**/`

### 2. Le(a) syntacticien(ne)

Il écrit le code CUP/JFlex. Il doit intervenir sur la gestion des erreurs de syntaxe et permettre au compilateur une poursuite de calcul en cas d'erreur. Il doit s'efforcer à produire une grammaire LALR.

On peut ajouter une syntaxe plus complexe pour les classes, l'héritage, l'encapsulation, etc.

Il serait intéressant de pouvoir traiter des modules et donc de permettre la compilation d'un ensemble de fichiers. Le lexer doit pouvoir traiter un ensemble de flux grâce à des `include`.

### 3. Le(a) sémanticien(ne)

Il écrit le code intermédiaire permettant au compilateur de gérer le typage des variables, la gestion du code intermédiaire et sa bonne formation.

Il serait intéressant de surcharger les fonctions, et donc de ne pas effacer la définition du type d'une fonction par une nouvelle définition. Une solution possible est d'ajouter une table des symboles avec surcharge (qui revient à un ensemble de signatures différentes pour le même nom) concurrente de la table des symboles qui masque une variable d'un bloc à un autre. Il est nécessaire de pouvoir trouver le code de la fonction surchargée par unification entre l'appel et la fonction.

Exemple :

```
1 function foo(i: int): int{
2     code1
3 }
4
5 function foo(i: int): bool{
6     code2
7 }
8
9 bool a := foo(6); // appel du code code2
```

On pourra ajouter le polymorphisme en utilisant les variables de type. La variable doit alors recevoir un type constant à la compilation (c'est le but du typage statique). Et chaque portion de code dépendant d'une variable de type doit être copié dans le code cible.

### 4. Un développeur

Il écrit le code cible en Java en utilisant au maximum les bibliothèques proposées pour ce langage.

Chaque variable du langage Léa correspond à un objet Java. Il convient d'allouer de la mémoire pour cet objet quand cela est pertinent.

Le type `enum` mérite un encodage avec des vecteurs de bits. Chaque élément d'un type énuméré correspond alors au poids d'un nombre binaire.

Les fonctions et opérateurs surchargés doivent être distingués par leur signature (le type de leurs paramètres et de la valeur rendue). Il convient de faire la gestion de ces noms pour le code cible.

Les paramètres des fonctions sont passés par valeur ou référence explicitement (comme en C++). En java, on a un passage par référence implicite si l'argument est un objet, sinon c'est un passage par valeur.

Pour écrire les fonctions, il est possible de produire le code correspondant en Java presque immédiatement. Mais dans le cas d'un passage par valeur pour une fonction récursive, cela pose problème. Il faut trouver une solution pour permettre une programmation récursive avec des arguments de type non élémentaires.

Les fonctions polymorphes sont recopiées pour chaque instance de type calculée en statique. Il serait intéressant de partager une portion constante de ces codes redondants (très difficile).

FIG. 2 – Exemple de code Léa

```

/* *****
* Un exemple complet de code Léa étendu
***** */

// Définition des types
fiche = class {
    static id: int;
    nom: string;
    age: int;

    // définition d'un constructeur
    fiche(){
        this.nom := "";
        this.age := 0;
    }

    // définition de l'opérateur >
    operator >(a: fiche, b: fiche){
        return (a.age > b.age);
    }
}

// Définition d'une fonction
function sort(T: list of 'x) {
    repeat {
        echange : bool := false;
        for i in [0 .. MAX-2] {
            if (T[j] > T[j+1]){
                // Variable locale
                aux: 'x;
                aux := T[j];
                T[j] := T[j+1];
                T[j+1] := aux;
                echange := true;
            }
        }
    } while (echange);
}

function main(argn:int, argv list of int): int{
    // Définition des variables
    table: list of fiche;

    // Instructions
    for i in [0 .. MAX-1] {
        table[i].nom := "";
        table[i].age := random(100);
    }

    sort(table);
    return 0;
}

```