

Projet de compilation

PREMIÈRE PARTIE

28 mars 2019

1 Introduction

Le langage dont il est question dans ce projet est inspiré d'un langage classique de programmation impérative (Algol, Pascal, C, etc) . Nous l'appellerons « *Léa* » (Langage Élémentaire Algorithmique).

Le projet consiste à analyser ce langage impératif et à le traduire en code intermédiaire.

Dans une seconde partie, qui sera donnée dans environ dix jours, nous dirons comment générer le code pour une machine Y86.

2 Travail à réaliser

1. Un analyseur lexical du langage Léa. Outre que cet analyseur doit reconnaître l'ensemble des tokens utilisés pour l'analyse syntaxique, il devra aussi reconnaître et encoder correctement :
 - Les commentaires de type C (`/* */`)
 - Les commentaires de type C++ (`//`)
 - Les littéraux **integer** codés en décimal et en hexadécimal (*i.e.* `0x1CFE`)
 - Les chaînes de caractères écrites entre `"`.

Cet analyseur lexical conservera le numéro de ligne, le numéro de colonne pour d'éventuels messages d'erreur.

2. Compléter l'analyseur syntaxique pour le langage *Léa*.

Cette fonction teste que le texte est bien engendré par la grammaire de *Léa* et envoie un message d'erreur précis dans le cas contraire.

3. Un analyseur sémantique pour *Léa*

Chaque expression doit avoir un type bien formé : l'analyseur sémantique envoie un message explicite en cas de mauvaise formation d'un type, et construit le typage des expressions dans le cas contraire.

L'analyseur construit une table de symboles pour

- Les définitions de types
- Les variables globales
- Les fonctions et les procédures
- Les variables locales

4. Un compilateur de *Léa* vers le code intermédiaire.

3 Éléments fournis

- Le fichier `build.xml`
- La classe principale `ubordeaux.deptinfo.compilation.project.main.Main`
- La grammaire incomplète `parser/Parser.grammar`
- Le package `ubordeaux.deptinfo.compilation.project.main.abstractTree.*` qui contient toutes les classes permettant de construire les arbres de syntaxe abstraite.
- Le package `ubordeaux.deptinfo.compilation.project.main.type.*` qui contient toutes les classes permettant de construire les types.
- Les fichiers de test `question-*/test-*.lea` qui doivent pouvoir être analysés en fonction des questions posées.

4 Modalités de réalisation

Groupes Le projet doit être réalisé par des groupes de 4 étudiants¹ et très exceptionnellement moins sur autorisation des enseignants.

Chaque groupe est représenté par l'un des trois étudiants qui est le *correspondant*. Le *correspondant* n'a aucune autorité sur le reste du groupe, son rôle est :

- De correspondre avec l'enseignant
- De correspondre avec les trois autres étudiants
- D'administrer le dépôt *GIT* ou *SVN*

Livable Le livable doit être déposé par le correspondant du groupe sur un dépôt du Cremi accessible par l'enseignant selon les modalités exigées lors des séances de TD.

La dernière version devra être déposée avant le 26 avril à 23h59 et ne plus être modifiée après. Le livable contient les fichiers suivants :

- `AUTHORS`

Qui contient le noms des 4 étudiants dans l'ordre alphabétique avec les réalisations effectuées (plusieurs réalisations par étudiant et plusieurs étudiants par réalisation sont possibles).

- `documentation.pdf`

Un rapport d'environ 4 pages qui explique le plan de votre programme, les difficultés rencontrées, les solutions apportées, les tests effectués. Ce document doit pouvoir être lu indépendamment des fichiers sources.

Si un ou plusieurs étudiants du groupe a des difficultés à rédiger correctement en langue française parce qu'il ne serait pas francophone, il lui est autorisé de rédiger en langue anglaise une section ou un chapitre entier. Dans tous les cas, la rédaction doit être soignée, tant dans l'expression que dans la forme (erreurs typographiques, mise en page, etc.).

- Les sources complètes

Soutenances des projets Les 29 ou 30 avril selon le calendrier fixé par Lionel Clément, votre groupe fera une démonstration commentée de 20 minutes, questions comprises.

1. dans ce qui suit *étudiant*, *correspondant*, etc. valent pour *étudiant.e*, *correspondant.e* etc.

5 Langage source : définition de la syntaxe

Le langage *Léa* est défini par la grammaire suivante :

Programme

```
1 program=
2   type_declaration_part
3   variable_declaration_part
4   procedure_definition_part
5   'begin'
6   statement_list.stm
7   'end'
```

Déclaration des types nommés

```
1 type_declaration_part:
2   /* empty */
3   | 'type' type_declaration_list
4
5 type_declaration_list:
6   type_declaration_list type_declaration
7   | type_declaration
8
9 type_declaration:
10  IDENTIFIER '=' type ';' ;
```

Définition des types

```
1 type:
2   simple_type
3   | named_type
4   | index_type
5   | array_type
6   | pointer_type
7   | structure_type
8
9 simple_type:
10  'string'
11  | 'integer'
12  | 'boolean'
13
14 named_type:
15  IDENTIFIER
16
17 index_type:
18  enumerated_type
19  | subrange_type
```

```

20
21 enumerated_type:
22     '(' identifier_list ')'
23
24 subrange_type:
25     INTEGER '..' INTEGER
26     | IDENTIFIER '..' IDENTIFIER
27
28 array_type:
29     'array' '[' range_type ']' 'of' type
30
31 range_type:
32     enumerated_type
33     | subrange_type
34     | named_type
35
36 pointer_type:
37     '^' type
38
39 structure_type:
40     'struct' '{' feature_list_type '}'
41
42 feature_list_type:
43     feature_list_type feature_type
44     | feature_type
45
46 feature_type:
47     IDENTIFIER ':' type ';'

```

Déclaration des variables typées

```

1 variable_declaration_part:
2     /* empty */
3     | 'var' variable_declaration_list
4
5 variable_declaration_list:
6     variable_declaration_list variable_declaration
7     | variable_declaration
8
9 variable_declaration:
10    identifier_list ':' type ';'
11
12 identifier_list:
13    identifier_list ',' IDENTIFIER
14    | IDENTIFIER

```

Déclaration et définition des procédures et des fonctions

```

1 procedure_definition_part=

```

```

2 | procedure_definition_list
3
4 procedure_definition_list=
5   procedure_definition_list
6   procedure_definition
7 | procedure_definition
8
9 procedure_definition=
10  procedure_definition_head block
11 | procedure_definition_head TOKEN_SEMIC
12
13 procedure_definition_head=
14  'procedure' IDENTIFIER '(' argt_part ')'
15 | 'function' IDENTIFIER '(' argt_part.args ')' ':' type
16
17 argt_part:
18   /* empty */
19 | argt_list
20
21 argt_list:
22   argt_list ',' argt
23 | argt
24
25 argt:
26   IDENTIFIER ':' type

```

Blocs

```

1 block:
2   variable_declaration_part
3   'begin'
4   statement_list
5   'end'

```

Instructions

```

1 statement_list:
2   statement_list statement
3 | statement
4
5 statement:
6   simple_statement
7 | structured_statement
8
9 simple_statement:
10  assignment_statement
11 | procedure_statement
12 | new_statement
13 | dispose_statement

```

```

14 | println_statement
15 | readln_statement
16 | return_statement
17
18 assignment_statement:
19     variable_access '=' expression ';'
20
21 procedure_statement:
22     procedure_expression ';'
23
24 procedure_expression:
25     IDENTIFIER '(' expression_part ')'
26
27 expression_part:
28     /* empty */
29     | expression_list
30
31 expression_list:
32     expression_list ',' expression
33     | expression
34
35 new_statement:
36     'new' variable_access ';'
37
38 dispose_statement:
39     'dispose' variable_access ';'
40
41 println_statement:
42     'println' expression ';'
43
44 readln_statement:
45     'readln' expression ';'
46
47 return_statement:
48     'return' expression ';'
49
50 structured_statement:
51     block
52     | if_statement
53     | while_statement
54     | switch_statement
55
56 if_statement:
57     'if' expression 'then' statement
58     | 'if' expression 'then' statement 'else' statement
59
60 while_statement:
61     'while' expression 'do' statement
62
63 while_statement=
64     'while' expression 'do' statement

```

```

65
66 switch_statement=
67     'switch' expression 'begin' case_statement_list 'end'
68
69 case_statement_list=
70     case_statement_list case_statement
71     | case_statement
72
73 case_statement=
74     'case' IDENTIFIER ':' statement
75     | 'default' ':' statement

```

Expressions

```

1 variable_access:
2     IDENTIFIER
3     | variable_access '[' expression ']'
4     | expression '^'
5
6 expression:
7     expression '+' expression
8     | expression '-' expression
9     | '-' expression
10    | expression '*' expression
11    | expression '/' expression
12    | expression '||' expression
13    | expression '&&' expression
14    | '!' expression
15    | expression '<' expression
16    | expression '<=' expression
17    | expression '>' expression
18    | expression '>=' expression
19    | expression '==' expression
20    | expression '!=' expression
21    | '(' expression ')'
22    | procedure_expression
23    | variable_access
24    | literal

```

Expressions littérales

```

1 literal:
2     INTEGER
3     | STRING
4     | 'true'
5     | 'false'
6     | 'null'

```

Où les symboles IDENTIFIER, INTEGER,STRING et IDENTIFIER représentent les unités lexicales suivantes :

IDENTIFIEUR Identificateur au format Java.

INTEGER Une constante numérique entière comprise entre 0 et 0xFFFFFFFF.

STRING Une chaîne d'au plus 64 caractères entre deux ".

Toutes les autres unités lexicales sont représentées dans la grammaire entre apostrophes.

Les opérateurs ont les mêmes propriétés algébriques que ce qui est défini en langage Java.

Remarque à retenir au moment d'écrire l'analyseur syntaxique : Le programme **Beaver** ne permet pas l'usage habituel de **%nonassoc**. Par conséquent, il vous est conseillé d'utiliser **%left** pour les opérateurs de comparaison.

6 Langage source : types

Définition : Expression de type est et n'est rien d'autre que :

- Types simples
 - « **string** », « **integer** », « **boolean** » sont des expressions de type
- Énumérés
 - une valeur d'un type énuméré est un type
 - si $\tau_1, \tau_2, \dots, \tau_k$ sont des valeurs de type énuméré, alors « $enum(\tau_1, \tau_2, \dots, \tau_k)$ » est un type énuméré
- Tableau
 - si min et max sont deux entiers avec $max > min$, « $[min..max]$ » est un type intervalle
 - si τ_1 est un type, et τ_2 un type intervalle ou un énuméré, alors « $array(\tau_1, \tau_2)$ » est un type tableau
- Pointeur
 - si τ est un type, alors « $pointer(\tau)$ » est un type pointeur
- Fonctions et procédures
 - « *void* » est un type
 - si $\tau_1, \tau_2 \dots \tau_k$ sont des types, alors « $\tau_1 \times \tau_2 \times \dots \tau_k$ » est un type produit
 - si τ_1 est un type produit et τ_2 un type, alors « $\tau_1 \rightarrow \tau_2$ » est un type fonction
 - si τ_1 est un type fonction et τ_2 un type produit, alors « $\tau_1(\tau_2)$ » est un type application
- Structures
 - si $name$ est un nom de champ et τ est un type, alors « $feature(name, \tau)$ » est un type feature
 - si $\tau_1, \tau_2 \dots \tau_k$ sont des types feature, alors « $struct(\tau_1, \tau_2 \times, \tau_k)$ » est un type structure

Bonne formation d'un code écrit en langage Léa

Un programme écrit en *Léa* est un texte engendré par la grammaire décrite plus haut qui vérifie les contraintes suivantes :

1. Tous les identificateurs (types définis, valeurs d'énumérés, variables, fonctions et procédures) sont déclarés avant d'être utilisés. Les fonctions et procédures peuvent être déclarées puis ensuite définies. Si elles sont seulement définies, leur déclaration est implicite.
2. Toute expression est bien typée, c'est-à-dire qu'il lui correspond une expression de type telle que définie plus haut et définie selon les tableaux suivants :

Opérateur binaire

	Type de opérande 1	Type de opérande 2	Type de la valeur
<code>+, *, -, /</code>	<code>integer</code>	<code>integer</code>	<code>integer</code>
<code><, >, <=, >=</code>	<code>integer</code>	<code>integer</code>	<code>boolean</code>
<code>==, !=</code>	<code>integer</code> <code>string</code> <code>pointer(τ)</code>	<code>integer</code> <code>string</code> <code>pointer(τ)</code>	<code>boolean</code>
<code>&&, </code>	<code>boolean</code>	<code>boolean</code>	<code>boolean</code>
<code>[]</code>	<code>array(τ_1, τ_2)</code>	<code>integer</code> ou valeur d'énuméré compris dans τ_2	τ_1
<code>()</code>	$\tau_1^1 \times \tau_1^2 \dots \tau_1^k \rightarrow \tau_2$	$\tau_3^1 \times \tau_3^2 \dots \tau_3^k$	τ_2
<code>.</code>	<code>struct($\tau_1 \times \tau_2 \dots \tau_k$)</code>	<code>feature(name, τ)</code>	τ

Opérateur unaire

	Type de l'Opérande	Type de la Valeur
<code>-</code>	<code>integer</code>	<code>integer</code>
<code>!</code>	<code>boolean</code>	<code>boolean</code>
<code>^</code>	<code>pointer(τ)</code>	τ

Expression

	Valeur
<code>true</code>	<code>boolean</code>
<code>false</code>	<code>boolean</code>
<code>null</code>	<code>pointer(τ)</code>
<code>INTEGER</code>	<code>integer</code>
<code>STRING</code>	<code>string</code>
<code>IDENTIFIER</code>	type de la variable dans son environnement courant

3. Toute instruction est bien typée :
 - L'affectation doit avoir des membres gauche et droit de même type.
 - Le test et la boucle ont comme argument une expression de type booléen.
 - L'instruction `return` a comme argument une expression de même type que le type de retour de la fonction qui la contient.
 - Les instruction `new` et `dispose` ont comme argument un pointeur.

7 Langage source : définition de la sémantique

Valeurs des expressions typées

Aux expressions de type correspondent un ensemble de valeurs dans $\mathbb{Z}/n\mathbb{Z}$.

Le tableau suivant donne pour chaque type, le nombre d'octets alloués en mémoire et l'encodage attendu :

Type	Image	Taille (en octets)	Encodage
integer	$\mathbb{Z}/2^{32}\mathbb{Z}$	4	encodés comme le sont les <i>signed int</i> en C (1 bit pour le signe, 31 pour le nombre). C'est-à-dire un relatif compris dans l'intervalle $[-2^{32-1}, 2^{32-1} - 1]$.
string	$\mathbb{Z}/2^{512}\mathbb{Z}$	64	une chaîne UTF-8 sur 64 octets
boolean	$\mathbb{Z}/2^8\mathbb{Z}$	1	0 pour false, $\neq 0$ pour true
Valeur d'énuméré	$\mathbb{Z}/32\mathbb{Z}$	4	encodé 2^x où x est l'indice de la valeur comprise entre 0 et 31. Par exemple l'énuméré (ROUGE, NOIR, BLEU, VERT, ROSE) sera encodé (1, 2, 4, 8, 16) et VERT sera encodé 8.
<i>pointer</i> (τ)	$\mathbb{Z}/2^{32}\mathbb{Z}$	4	Encodage d'une adresse sur 4 octets.
<i>array</i> (τ, min, max)	$\mathbb{Z}/n\mathbb{Z}$	$(max - min + 1) \times taille(\tau)$	Encodage d'un tableau d'octets.

Sémantique des opérateurs

La sémantique des opérateurs arithmétiques sur les entiers est une application $(\mathbb{Z}/2^{32}\mathbb{Z})^k \rightarrow \mathbb{Z}/2^{32}\mathbb{Z}$, où k est l'arité de l'opérateur.

La sémantique des opérateurs logiques sur les booléens est une application $(\mathbb{Z}/2^8\mathbb{Z})^k \rightarrow \mathbb{Z}/2^8\mathbb{Z}$, où k est l'arité de l'opérateur.

La sémantique des opérateurs de comparaison est une application $(\mathbb{Z}/2^n\mathbb{Z})^2 \rightarrow \mathbb{Z}/2^8\mathbb{Z}$. Où n dépend des opérandes pour l'égalité et la différence : $n = 2^{32}$ pour des entiers, 2^8 pour des booléens et 2^{512} pour des chaînes de caractères.

Tous ces opérateurs ont une sémantique habituelle ; celle que l'on trouve dans le langage C.

Sémantique des expressions et instructions

- La procédure `println x` permet de provoquer un effet vers la sortie standard. L'affichage sera alors la valeur de l'expression `x` selon son type, suivi du caractère `'\n'`. Les valeurs sont
 - Entier relatif décimal signé pour **integer**,

- « true » ou « false » pour *boolean*,
 - La chaîne de caractères UTF-8 pour *string*,
 - « array » pour $array(\tau, min, max)$,
 - « pointer » pour $pointer(\tau)$.
- La procédure **readln** *x* lit depuis l'entrée standard une valeur entière. *x* doit être de type *integer*, *string* ou *boolean*.
 - Si *p* est de type $pointer(\tau)$, l'expression p^\wedge désigne la valeur pointée à l'adresse *p*. Si *p* a comme valeur *null*, le programme produit un message d'erreur.
 - Si *t* est de type $array(\tau, min, max)$, et *e* de type *integer* ou valeur d'énuméré, l'expression $t[e]$ désigne la valeur de $t + (\sigma(e) - min)$ où $\sigma(e)$ est la valeur entière de l'expression *e*.
 - Si *o* est de type $struct(\tau_1, \tau_2, \dots \tau_k)$, et *n* un nom de champ, l'expression *o.n* désigne la valeur de $o + \omega$ où ω désigne le décalage du champ désigné par *n*.
 - Les procédures et fonctions utilisent exclusivement un passage de paramètres par valeur et les valeurs ne peuvent être que de type *integer*, *boolean*, *string*, $pointer(\tau)$.
 - Les affectations ne peuvent se faire que sur les expressions dites *accessibles*, c'est-à-dire, les variables de type simple, ou les expressions $t[e]$ où *t* est de type $array(\tau, min, max)$, ou les expressions p^\wedge où *p* est de type $pointer(\tau)$.
 - L'instruction **switch** *x* **begin** ... **end**
es affectations ne peuvent se faire que sur les expressions dites *accessibles*, c'est-à-dire, les variables de type simple, ou les expressions $t[e]$ où *t* est de type $array(\tau, min, max)$, ou les expressions p^\wedge où *p* est de type $pointer(\tau)$.

8 Langage intermédiaire

Le langage intermédiaire de notre compilateur est le langage à 3 opérateurs vu en cours et décrit pp. 123 et suivantes de (*Modern Compiler Implementation in Java, Second Edition*, Andrew W. Appel and Jens Palsberg, Cambridge University Press, 2002) donné en référence bibliographique du cours.

Nous citons le début du passage ici :

«

Here is a description of the meaning of each tree operator. First, the expressions (**Exp**), which stand for the computation of some value (possibly with side effects) :

«

- **CONST**(*i*) The integer constant *i*.
- **NAME**(*n*) The symbolic constant *n* (corresponding to an assembly language label).
- **TEMP**(*t*) Temporary *t*. A temporary in the abstract machine is similar to a register in a real machine. However, the abstract machine has an infinite number of temporaries.
- **BINOP**(*o*, *e*₁, *e*₂) The application of binary operator *o* to operands *e*₁, *e*₂. Subexpression *e*₁ is evaluated before *e*₂. The integer arithmetic operators are **PLUS**, **MINUS**, **MUL**, **DIV**; the integer bitwise logical operators are **AND**, **OR**, **XOR**; the integer logical shift operators are **LSHIFT**, **RSHIFT**; the integer arithmetic right-shift is **ARSHIFT**. The MiniJava language has only one logical operator, but the intermediate language is meant to be independent of any source language; also, the logical operators might be used in implementing other features of MiniJava.
- **MEM**(*e*) The contents of *wordSize* bytes of memory starting at address *e* (where *wordSize* is defined in the **Frame** module). Note that when **MEM** is used as the left child of a **MOVE**, it means "store", but anywhere else it means "fetch".
- **CALL**(*f*, *l*) A procedure call : the application of function *f* to argument list *l*. The subexpression *f* is evaluated before the arguments which are evaluated left to right.
- **ESEQ**(*s*, *e*) The statement *s* is evaluated for side effects, then *e* is evaluated for a result.

«

The statements (**stm**) of the tree language perform side effects and control flow :

- **MOVE**(**TEMP** *t*, *e*) Evaluate *e* and move it into temporary *t*.
- **MOVE**(**MEM**(*e*₁), *e*₂) Evaluate *e*₁, yielding address *a*. Then evaluate *e*₂, and store the result into *wordSize* bytes of memory starting at *a*.
- **EXP**(*e*) Evaluate *e* and discard the result.
- **JUMP**(*e*, *labs*) Transfer control (jump) to address *e*. The destination *e* may be a literal label, as in **NAME**(*lab*), or it may be an address calculated by any other kind of expression. For example, a C-language **switch**(*i*) statement may be implemented by doing arithmetic on *i*. The list of labels *labs* specifies all the possible locations that the expression *e* can evaluate to; this is necessary for dataflow analysis later. The common case of jumping to a known label *l* is written as **JUMP**(**NAME**(*l*), *new LabelList*(*l*, *null*)), but the **JUMP** class has an extra constructor so that this can be abbreviated as **JUMP**(*l*).
- **CJUMP**(*o*, *e*₁, *e*₂, *t*, *f*) Evaluate *e*₁, *e*₂ in that order, yielding values *a*, *b*. Then compare *a*, *b* using the relational operator *o*. If the result is true, jump to *t*; otherwise jump to *f*. The relational operators are **EQ** and **NE** for integer equality and nonequality (signed or unsigned); signed integer inequalities **LT**, **GT**, **LE**, **GE**; and unsigned integer inequalities **ULT**, **ULE**, **UGT**, **UGE**.

- **SEQ**(s_1, s_2) The statement s_1 followed by s_2 .
- **LABEL**(n) Define the constant value of name n to be the current machine code address. This is like a label definition in assembly language. The value **NAME**(n) may be the target of jumps, calls, etc.