

1 Lexique

1.1 “Atomes” du langage

Unit : ‘unit’
Vrai : ‘vrai’
Faux : ‘faux’
Entier : ‘[0-9][0-9_]*’
Flottant : ‘([0-9][0-9_]*)?\.[0-9_]*’
Chaine : ‘"[\.]*”’
IdfMin : ‘[a-z][a-zA-Z_]*’

1.2 Opérateurs

Et : ‘et’
Ou : ‘ou’
Egale : ‘==’
Different : ‘<>’
Inf : ‘<’
InfEgale : ‘<=’
Sup : ‘>’
SupEgale : ‘>=’
Plus : ‘+’
Moins : ‘-’
Mul : ‘*’
Div : ‘/’
Modulo : ‘%’

1.3 Genre (Types)

Pipe : ‘|’
Appli : ‘->’
FnGenre : ‘Fn’
IdfMaj : ‘[A-Z][a-zA-Z_]*’

1.4 Affectations

Aff : ‘=’
Incr : ‘++’
Decr : ‘--’
AffPlus : ‘+=’
AffMoins : ‘-=’
AffMul : ‘*=’
AffDiv : ‘/=’
AffModulo : ‘%=’

1.5 Séparateurs & Délimiteurs

ParO : '('

ParF : ')'

AccolO : '{'

AccolF : '}'

Comment : '#'

Virgule : ','

DeuxPoints : ':'

FinLigne : '[\CR\LF]+'

1.6 Structure de contrôle

Fn : 'fn'

Retour : 'retour'

Si : 'si'

Sinon : 'sinon'

Tq : 'tq'

Boucle : 'boucle'

Remarques :

- En cas de 'conflit de caractère' (eg : '<' et '<=') le lexer reconnaît l'élément le plus long (ie : '<=')
- Les mots-clés sont reconnus comme des identifiants dans un premier temps et distingués a posteriori
- Le lexer produit un lexème **ChaineNonFermee** le cas échéant
- Le lexer produit un lexème **Erreur** pour tout caractère non reconnu

2 Grammaire

2.1 Expression

Expression \rightarrow Disj
Disj \rightarrow Conj DisjSuite
DisjSuite \rightarrow ‘*Ou*’ Disj
DisjSuite \rightarrow ‘ ϵ ’
Conj \rightarrow Rel ConjSuite
ConjSuite \rightarrow ‘*Et*’ Conj
ConjSuite \rightarrow ‘ ϵ ’
RelOpde \rightarrow Somme
Rel \rightarrow RelOpde
Rel \rightarrow RelOpde ‘*Egale*’ RelOpde
Rel \rightarrow RelOpde ‘*Different*’ RelOpde
Rel \rightarrow RelOpde ‘*Inf*’ RelOpde
Rel \rightarrow RelOpde ‘*InfEgale*’ RelOpde
Rel \rightarrow RelOpde ‘*Sup*’ RelOpde
Rel \rightarrow RelOpde ‘*SupEgale*’ RelOpde
Somme \rightarrow Terme SommeSuite
SommeSuite \rightarrow ‘*Plus*’ Somme
SommeSuite \rightarrow ‘*Moins*’ Somme
SommeSuite \rightarrow ‘ ϵ ’
Terme \rightarrow Facteur TermeSuite
TermeSuite \rightarrow ‘*Mul*’ Terme
TermeSuite \rightarrow ‘*Div*’ Terme
TermeSuite \rightarrow ‘ ϵ ’
Facteur \rightarrow Unaire Appel
Appel \rightarrow ‘*ParO*’ AppelArgs ‘*ParF*’ Appel
Appel \rightarrow ‘*ParO*’ ‘*ParF*’ Appel
Appel \rightarrow ‘ ϵ ’
AppelArgs \rightarrow Expression AppelArgsSuite
AppelArgsSuite \rightarrow ‘*Virgule*’ AppelArgs
AppelArgsSuite \rightarrow ‘ ϵ ’
Unaire \rightarrow ‘*Plus*’ Unaire
Unaire \rightarrow ‘*Moins*’ Unaire
Unaire \rightarrow Atome
Atome \rightarrow ‘*ParO*’ Expression ‘*ParF*’
Atome \rightarrow ‘*Unit*’
Atome \rightarrow ‘*Vrai*’
Atome \rightarrow ‘*Faux*’
Atome \rightarrow ‘*Entier*’
Atome \rightarrow ‘*Flottant*’
Atome \rightarrow ‘*Chaine*’
Atome \rightarrow ‘*IdfMin*’
Atome \rightarrow ‘*Fn*’ FnArgs FnRetour Bloc

FnArgs \rightarrow '*ParO*' FnArg FnArgsSuite '*ParF*'
FnArgs \rightarrow '*ParO*' '*ParF*'
FnArgs \rightarrow ' ϵ '
FnArg \rightarrow '*IdfMin*' '*DeuxPoints*' Genre
FnArgSuite \rightarrow '*Virgule*' FnArg FnArgSuite
FnArgSuite \rightarrow ' ϵ '
FnRetour \rightarrow '*Appli*' Genre
FnRetour \rightarrow ' ϵ '

2.2 Genre (Déclaration des types)

Genre \rightarrow Union
Union \rightarrow GenreTerme UnionSuite
UnionSuite \rightarrow '*Pipe*' Union
UnionSuite \rightarrow ' ϵ '
GenreTerme \rightarrow GenreAtome
GenreAtome \rightarrow '*IdfMaj Tous*'
GenreAtome \rightarrow '*IdfMaj Unit*'
GenreAtome \rightarrow '*IdfMaj Vrai*'
GenreAtome \rightarrow '*IdfMaj Faux*'
GenreAtome \rightarrow '*IdfMaj Booleen*'
GenreAtome \rightarrow '*IdfMaj Entier*'
GenreAtome \rightarrow '*IdfMaj Flottant*'
GenreAtome \rightarrow '*IdfMaj Nombre*'
GenreAtome \rightarrow '*IdfMaj Chaine*'
GenreAtome \rightarrow '*ParO*' Genre '*ParF*'
GenreAtome \rightarrow '*FnGenre*' GenreFnArgs GenreFnRetour
GenreFnArgs \rightarrow '*ParO*' Genre GenreFnArgsSuite '*ParF*'
GenreFnArgs \rightarrow '*ParO*' '*ParF*'
GenreFnArgs \rightarrow ' ϵ '
GenreFnArgsSuite \rightarrow '*Virgule*' Genre GenreFnArgsSuite
GenreFnRetour \rightarrow '*Appli*' Genre

2.3 Instruction

Programme \rightarrow Corps '*FinSequence*'
Bloc \rightarrow '*AccolO*' Corps '*AccolF*'
Corps \rightarrow Instruction CorpsSuite
CorpsSuite \rightarrow Comment FinLigne Corps
Comment \rightarrow '*Comment*' [Tous lèxemes sauf '*FinLigne*']+
CorpsSuite \rightarrow ' ϵ '
Instruction \rightarrow Retour
Instruction \rightarrow Si
Instruction \rightarrow Sinon
Instruction \rightarrow Tq
Instruction \rightarrow Boucle
Instruction \rightarrow '*Idf*' IdfInstruction

Instruction \rightarrow '*FinLigne*'
Retour \rightarrow '*Retour*'
Retour \rightarrow '*Retour*' Expression
Si \rightarrow '*Si*' Expression Bloc
Sinon \rightarrow '*Sinon*' Si
Sinon \rightarrow '*Sinon*' Bloc
Tq \rightarrow '*Tq*' Expression Bloc
Boucle \rightarrow '*Boucle*' Bloc '*Tq*' Expression
IdfInstruction \rightarrow Appel
IdfInstruction \rightarrow '*Aff*' Expression
IdfInstruction \rightarrow '*Incr*'
IdfInstruction \rightarrow '*Decr*'
IdfInstruction \rightarrow '*AffPlus*' Expression
IdfInstruction \rightarrow '*AffMoins*' Expression
IdfInstruction \rightarrow '*AffMul*' Expression
IdfInstruction \rightarrow '*AffDiv*' Expression
IdfInstruction \rightarrow '*AffModulo*' Expression

Remarques :

- Les fins de lignes à l'exception du cas où elles sont rencontrées en fin d'instruction sont gérés en dehors de la grammaire.
- La vérification qu'un sinon apparait à la suite d'un si n'est pas géré par la définition de la grammaire mais comme un cas particulier dans le code.

3 Sémantique

3.1 Affectations

La gestion des affectations par le langage est extrêmement laxiste, il est possible d'affecter un entier à une variable puis de lui affecter une chaîne de caractères plus loin voire de lui assigner deux valeurs de types différents dans deux branches alternatives de code. Néanmoins l'analyseur sémantique garde trace du type de chaque variable au fur et à mesure du flot du programme et effectue une vérification précise de la compatibilité entre les types des variables et leurs utilisations.

3.2 Fonction de première classe et fermeture

Le langage supporte la définition de nouvelles fonctions. De plus, ces fonctions sont de première classe (manipulables comme des valeurs au même titre qu'un entier ou une chaîne de caractères) et supporte le principe de fermeture (par portée lexicale). Elles sont définies exclusivement anonymement (ce n'est qu'en étant assignée à une variable qu'elles obtiennent un nom). La fermeture des variables au moment de la définition de la fonction est réalisée par 'copie'. Autrement dit, les nouvelles affectations a posteriori de la variable enfermée n'affectent pas la fonction :

```
x = 1
f = fn -> Entier {
    retour x
}

x = 2
assert(f() == 1)
```

3.3 Système de Type

Les types sont définis récursivement à partir des types de base. Les types de base sont les types concrets présents à l'exécution :

- Unit** : Type singleton de la valeur unit
- Vrai** : Type singleton de la valeur vrai
- Faux** : Type singleton de la valeur faux
- Entier** : Type des nombres entiers
- Flottant** : Type des nombres réels représentés avec virgule flottante
- Chaîne** : Type des chaînes de caractères

À partir de ces types peut être construit de nouveaux types avec les deux constructions suivantes :

- Union** : Etant donné deux types A et B leur union (notée 'A|B') représente un nouveau type. Dans la correspondance de Curry-Howard, il représente le 'ou' logique ; ie. avec P(T) le prédicat associé à un type T, on a " $P(A|B) \iff P(A) \cup P(B)$ ".
- Fonction** : Ce type est nécessaire pour représenter les fonctions qui sont de première classe. Il est notée $\text{Fn}(\text{Type_Arg1}, \dots) \rightarrow \text{Type_Retour}$ et encode le nombre d'arguments, leur types et celui du retour de la fonction.

Pour raisonner quant au types des variables le langage utilise la notion de généralisation. Un type A généralise un autre type B si toute valeur assignable au type B est assignable au type A. On remarque :

- $\forall A$, A généralise A
- $\forall A, B$, A|B généralise A et B
- $\text{Fn}(A_1, \dots, A_n) \rightarrow R$ généralise $\text{Fn}(B_1, \dots, B_n) \rightarrow S$ si et seulement si, $\forall i$, B_i généralise A_i et R généralise S. (Notons l'inversion pour le type des arguments. Si l'on assimile le type fonction à un type paramétrique, on retrouve le problème de la variance, i.e : Une fonction est covariante en son type de retour et contravariante en les types de ses arguments)

A partir de ce principe de généralisation est défini un nouveau type, le type Tous qui est le type qui généralise tout les autres types. (Ou de manière équivalente, le type qui est l'union de tous les autres types).

Les déclarations de types ne sont nécessaires que dans les déclarations de fonctions. Tout les autres types peuvent être inférés à partir des valeurs littérales et des types des fonctions.

Problème des fonctions récursives : Tel qu'est défini le langage les fonctions récursives posent problèmes. En effet, l'analyse sémantique est purement séquentielle et les fonctions n'ont un nom qu'après avoir été analysée (si elles sont assignées à une variable). Pour supporter les fonctions récursive l'analyseur est obligé de considérer le fait qu'une fonction est sur le point d'être affectée à une variable et de faire passer le nom de la variable en question à l'analyse de la fonction pour que cette dernière puissent reconnaître les occurrences d'appels récursifs. Une fonction peut ainsi être flaggée comme récursive et être considérée comme s'enfermant elle-même. Il faut noter que cette solution ne fonctionne pas pour les fonctions mutuellement récursives, par exemple :

```
pair = fn(n: Entier) -> Booleen {
  si n == 0 { retour vrai }
  retour impair(n - 1)
}

impair = fn(n: Entier) -> Booleen {
  si n == 0 { retour faux }
  retour pair(n - 1)
}
```

La première occurrence de la variable 'impair' induit une erreur car elle est considérée non définie par le compilateur.

4 Langage Intermédiaire & "Machine Virtuelle"

Après l'analyse sémantique et avant l'exécution, le programme est compilé dans un langage intermédiaire semblable à un langage d'assemblage ou à un bytecode. Ce langage intermédiaire est ensuite interprété par une "machine virtuelle" fonctionnant sur une base de pile. (Plus exactement deux piles mais on aurait très bien pu se ramener à une unique pile).

En effet, on peut grossièrement assimiler le langage intermédiaire à une transformation en notation postfixée du programme source facilement interpretable par une machine à pile. Par exemple l'expression :

```
x = 5 * 6 + 3
```

transformée en un arbre syntaxique abstrait :

```
Affectation
  x
  Somme
    Produit
      5
      6
    3
```

peut ensuite être compilée sous la forme

```
empiler 5
empiler 6
faire Produit
empiler 3
faire Somme
depiler/stocker dans x
```

en assumant que chaque opération (et par extension chaque fonction) prend ses arguments en sommet de pile et y met à la place sa valeur de retour.

Le jeu d'instruction du langage intermédiaire/de la machine virtuelle est défini sur ce principe :

- EPL** <valeur littéral> : Met la <valeur littéral> au sommet de la pile d'appel
- DPL** : Dépile le sommet de la pile d'appel
- CHR** <indice variable> : Charge la variable à <indice variable> et sur la pile d'appel
- STK** <indice variable> : Stocker le sommet de la pile (en dépilant) à <indice variable>
- PRM** <nom primitive> : Met primitive <nom primitive> au sommet de la pile d'appel
- CLT** <nom de la routine> [~]!<nombre de variables à cloturer> : Définition d'une fonction mise au sommet de la pile d'appel
- STI** <decalage> : Saut inconditionnel
- STV** <decalage> : Saut si vrai dépilé de la pile d'appel
- STF** <decalage> : Saut si faux dépilé de la pile d'appel
- APP** : Appel de la fonction sur la pile
- RET** : Fin de la fonction en cours et retour à la fonction appelante

Pour représenter le code des fonctions sous la forme de ce langage intermédiaire on utilise le concept de routine qui est la suite statique d'instructions d'une fonction. Les routines sont définies par un bloc commençant par le nom entre crochet suivi d'un nombre représentant l'allocation nécessaire à la fonction.

Une fonction est ainsi la donnée dynamique constituée d'une référence vers une routine et d'un ensemble de variables enfermées. Elles sont créées par l'instruction **CLT** qui associe ces deux composantes. Le caractère '~' devant le nombre de variables à capturer indique que la fonction est recursive.

Les variables locales de chaque fonction sont représentées par la deuxième pile fonctionnant par frame. Autrement dit, chaque fonction définit un nombre de variables locales dont elle a besoin (une allocation) et au moment de l'appel de la fonction la machine virtuelle déplace le curseur de la pile de la taille demandée et la fonction peut référencer ses variables en fonction d'indices déterminés statiquement vers cette frame.

L'argument des instructions de saut est défini en terme de nombre d'instructions de décalage par rapport à l'instruction de saut en elle-même. Elles sont suffisantes pour représenter les trois structures de contrôle et les deux opérateurs paresseux 'et' et 'ou' (uniquement STV aurait même suffi).

Le langage intermédiaire contient aussi des annotations informatives pour le nom des variables et les positions correspondantes du code source.

Exemple de compilation

```
[ ] 4
EPL 1 @1:14
STK 0 &fact_borne @1:12
EPL 1 @2:15
STK 1 &fact_default @2:13
EPL 5 @3:13
STK 2 &fact_test @3:11
CHR 0 &fact_borne @5:8
CHR 1 &fact_default @5:8
CLT fact0 ~2 @5:8
STK 3 &fact @5:6
CHR 2 &fact_test @10:16
CHR 3 &fact @10:11
APP @10:15
PRM ecrire_ln @10:10
APP @10:10
EPL unit @1:1
RET @1:1
```

```
[fact0] 4
```



```
STK 3 &fact0 @5:8
STK 2 &fact_default @5:8
STK 1 &fact_borne @5:8
STK 0 &n @5:8
CHR 0 &n @6:7
CHR 1 &fact_borne @6:12
PRM inferieur_egale @6:9
APP @6:9
STF 4 @6:4
CHR 2 &fact_default @6:32
RET @6:25
STI 1 @6:4
CHR 0 &n @7:11
CHR 0 &n @7:20
EPL 1 @7:24
PRM moins @7:22
APP @7:22
CHR 3 &fact @7:15
APP @7:19
PRM multiplier @7:13
APP @7:13
RET @7:4
```