

Compilateur pour langage while
Devoir de programmation, partie I : Analyse
Université d'Orléans, Master 1 Informatique, 2012-2013

L'objectif de ce devoir est de réaliser la phase d'analyse pour la compilation des programmes d'un langage simplifié. La production de code sera réalisée dans un second devoir. Vous devrez mettre en place les phases d'analyse lexicale, syntaxique et sémantique pour ce langage. La partie 1 décrit la syntaxe du langage, la partie 2 donne les règles de typage nécessaires à la phase d'analyse sémantique et la partie 3 décrit le travail à réaliser.

1 Syntaxe

Expressions. Les expressions du langage sont divisées en deux catégories, expressions arithmétiques A et expressions booléennes B .

$$A ::= n \mid d \mid x \mid x[A] \mid \text{unop } A \mid A \text{ binop } A \mid (A)$$

$$B ::= \text{true} \mid \text{false} \mid A \text{ cmp } A \mid ! B \mid B \&\& B \mid B \parallel B \mid (B)$$

où n dénote une constante entière, d une constante décimale, x un identificateur et

$$\text{unop} \in \{-\} \quad \text{binop} \in \{+, -, *, /\} \quad \text{cmp} \in \{<, >, <=, >=, =, <>\}$$

Un identificateur est une suite de caractères alpha-numériques commençant par une lettre minuscule. Une expression $x[A]$ représente l'accès à l'indice dénoté par A du tableau x . Les règles de priorités sont les règles usuelles des opérateurs arithmétiques et booléens.

Types. Les types du langage sont donnés par la règle $T ::= \text{int} \mid \text{decimal} \mid \text{array}(T)$. Il n'est pas nécessaire de distinguer un type pour les expressions booléennes.

Instructions. Les instructions du langage sont définies de la manière suivante

$$LHS ::= x \mid x[A]$$

$$I ::= \begin{array}{l} \text{return} \mid \text{return } A \\ \mid LHS := A \\ \mid LHS = \text{new } T[A] \mid \text{free}(A) \\ \mid f(A, \dots, A) \mid LHS := f(A, \dots, A) \end{array}$$

Une instruction peut être une affectation de la forme $LHS := A$ où LHS représente l'emplacement mémoire où sera stockée la valeur de l'expression A . La partie gauche LHS peut être une variable ou une case d'un tableau. L'instruction $LHS := \text{new } T[A]$ représente l'allocation d'un tableau de type T dont la taille est donnée par le résultat de l'évaluation de A . L'instruction $\text{free}(A)$, où A représente l'adresse d'un tel tableau, permet de récupérer la mémoire ainsi allouée. On dispose également d'instructions d'appels de procédure/fonction et de retour.

Statements. Les statements du langage sont définis à partir des instructions de la manière suivante

$$\begin{aligned}
S ::= & \mid I \\
& \mid \text{if } b \text{ then } S \text{ else } S \\
& \mid \text{while } b \text{ do } S \\
& \mid \text{begin } VDECL_1 \dots VDECL_m \ S_1; \dots; S_n \text{ end} \quad m \geq 0, n \geq 1
\end{aligned}$$

Un statement peut être une instruction, une conditionnelle, une boucle ou un block. Un block est une suite de déclaration de variables locales, données par la règle ci-dessous, suivie d'une séquence de statements.

$$VDECL ::= T \ x_1, \dots, x_n;$$

Programmes. Un programme est constitué d'une suite de déclarations de procédures et fonctions et d'un statement dénotant le corps du programme (équivalent de la fonction *main* en Java).

$$\begin{aligned}
P & ::= PDECL_1 \dots PDECL_n \ S & n \geq 0 \\
PDECL & ::= \mid \text{function } f(x_1 : T_1, \dots, x_n : T_n) : T \text{ is } S & n \geq 0 \\
& \mid \text{procedure } f(x_1 : T_1, \dots, x_n : T_n) \text{ is } S & n \geq 0
\end{aligned}$$

2 Typage

Typage des expressions On note Γ une suite d'éléments $x : \tau$. Intuitivement, Γ , que l'on appelle environnement de typage, joue le rôle de table des symboles. L'environnement se lit de droite à gauche, i.e s'il y a plusieurs occurrences d'une même variable, l'occurrence courante est celle qui se trouve la plus à droite. On note $\Gamma(x)$ le type associé à cette occurrence de x dans Γ . On définit ci-dessous les règles d'inférence des relation $\Gamma \vdash_a e : \tau$ et $\Gamma \vdash_b B$ qui caractérisent respectivement le bon typage d'une expression arithmétique et

d'une expression booléenne.

$$\begin{array}{c}
\frac{}{\Gamma \vdash_a n : \text{integer}} \quad \frac{}{\Gamma \vdash_a d : \text{decimal}} \quad \frac{\Gamma(x) = \tau}{\Gamma \vdash_a x : \tau} \\
\\
\frac{\Gamma \vdash_a A : \text{integer} \quad \Gamma(x) = \text{array}(\tau)}{\Gamma \vdash_a x[A] : \tau} \\
\\
\frac{\Gamma \vdash A : \tau \quad \tau \in \{\text{integer}, \text{decimal}\}}{\Gamma \vdash_a \text{unop } A : \tau} \\
\\
\frac{\Gamma \vdash_a A_1 : \tau \quad \Gamma \vdash_a A_2 : \tau \quad \tau \in \{\text{integer}, \text{decimal}\}}{\Gamma \vdash_a A_1 \text{ binop } A_2 : \tau} \\
\\
\frac{}{\Gamma \vdash_b \text{true}} \quad \frac{}{\Gamma \vdash_b \text{false}} \quad \frac{\Gamma \vdash_a A_1 : \tau \quad \Gamma \vdash_a A_2 : \tau \quad \tau \in \{\text{integer}, \text{decimal}\}}{\Gamma \vdash_b A_1 \text{ cmp } A_2} \\
\\
\frac{\Gamma \vdash B}{\Gamma \vdash_b ! B} \quad \frac{\Gamma \vdash B_1 \quad \Gamma \vdash B_2}{\Gamma \vdash_b B_1 \ \&\& \ B_2} \quad \frac{\Gamma \vdash B_1 \quad \Gamma \vdash B_2}{\Gamma \vdash_b B_1 \ \&\& \ B_2}
\end{array}$$

Typage des instructions On définit maintenant les règles permettant de vérifier qu'une instruction est correctement typée. Pour cela, on ajoute à l'environnement Γ , un environnement Δ qui associe un schéma à chaque symbole de fonction. Les règles sont maintenant de la forme

$$\Gamma; \Delta; t \vdash_i I$$

qui indique que l'instruction I est correctement typé pour les environnement Γ et Δ pour un type de retour t où t est soit \bullet dans le cas d'une procédure soit un type τ dans le cas

d'une fonction.

$$\begin{array}{c}
\frac{}{\Delta, \Gamma; \bullet \vdash_i \text{return}} \quad \frac{\Gamma \vdash_a e : \tau}{\Delta; \Gamma; \tau \vdash_i \text{return } e} \\
\\
\frac{\Gamma \vdash_a LHS : \tau \quad \Gamma \vdash_a A : \tau}{\Delta; \Gamma; t \vdash_i LHS := A} \quad \frac{\Gamma \vdash_a LHS : \text{array}(T) \quad \Gamma \vdash_a A : \text{integer}}{\Delta; \Gamma; t \vdash_i LHS := \text{new } T [A]} \\
\\
\frac{\Gamma \vdash_a A : \text{array}(\tau)}{\Delta; \Gamma; t \vdash_i \text{free}(A)} \\
\\
\frac{\Delta(f) = \tau_1 \times \dots \times \tau_n \quad \Delta; \Gamma; t \vdash_a A_i : \tau_i, \ i = 1, \dots, n}{\Delta; \Gamma; t \vdash_i f(A_1, \dots, A_n)} \\
\\
\frac{\Delta(f) = \tau_1 \times \dots \times \tau_n \rightarrow \tau \quad \Gamma \vdash_a LHS : \tau \quad \Delta; \Gamma; t \vdash_a A_i : \tau_i, \ i = 1, \dots, n}{\Delta; \Gamma; t \vdash_i LHS := f(A_1, \dots, A_n)}
\end{array}$$

Typage des statements Le typage des statements suit le même schéma que celui des instructions, on définit la relation \vdash_s par les règles suivantes.

$$\begin{array}{c}
\frac{\Delta; \Gamma; t \vdash_i I}{\Delta; \Gamma; t \vdash_s I} \\
\\
\frac{\Gamma \vdash_b b \quad \Delta; \Gamma; t \vdash_s S_1 \quad \Delta; \Gamma; t \vdash_s S_2}{\Delta; \Gamma; t \vdash_s \text{if } b \text{ then } S_1 \text{ else } S_2} \quad \frac{\Gamma \vdash_b b \quad \Delta; \Gamma; t \vdash_s S}{\Delta; \Gamma; t \vdash_s \text{while } b \text{ do } S} \\
\\
\frac{\Delta; \Gamma, x_1^1 : T^1, \dots, x_{k_1}^1 : T^1, \dots, x_1^m : T^1, \dots, x_{k_m}^m : T^m; t \vdash_s S_i \quad i = 1, \dots, n}{\Delta; \Gamma; t \vdash_s \text{begin } T^1 \ x_1^1, \dots, x_{k_1}^1; \dots; T^m \ x_1^m, \dots, x_{k_m}^m; S_1; \dots; S_n \text{end}}
\end{array}$$

Typage des programmes

$$\begin{array}{c}
\Delta; ; \bullet \vdash_s S \\
\\
\Delta; x_1 : \tau_1 \dots; \tau_n; \bullet \vdash_s S \quad \text{pour chaque } PDECL_i = \text{procedure } f(x_1 : \tau_1, \dots, x_n : \tau_n) \text{ is } S \\
\\
\Delta; x_1 : \tau_1 \dots; \tau_n; \tau \vdash_s S \quad \text{pour chaque } PDECL_i = \text{function } f(x_1 : \tau_1, \dots, x_n : \tau_n) : \tau \text{ is } S \\
\\
\frac{\Delta = \left(\bigcup_1^n \{f : \tau_1 \times \dots \times \tau_n \mid PDECL_i = \text{procedure } f(x_1 : \tau_1, \dots, x_n : \tau_n) \text{ is } S\} \right) \cup \left(\bigcup_1^n \{f : \tau_1 \times \dots \times \tau_n \rightarrow \tau \mid PDECL_i = \text{function } f(x_1 : \tau_1, \dots, x_n : \tau_n) : \tau \text{ is } S\} \right)}{\vdash PDECL_1 \dots PDECL_n S}
\end{array}$$

3 Travail à faire

Objectifs : Définir les différentes phases d'analyse pour le langage considéré ici, en utilisant `flex` et `bison`.

- Attention la description du langage n'est peut être pas adaptée à la reconnaissance syntaxique et il sera peut être nécessaire de redéfinir la grammaire pour lever les ambiguïtés.
- le système de type assure que les instructions `return` e et `return`, lorsqu'elles existent, sont légales. En revanche rien ne garantit qu'elles sont bien présentes dans tous les cas. Vous devrez proposer une solution pour régler ce problème, par exemple en modifiant les règles de typage.

Comportement de votre programme : votre programme devra

- prendre en entrée un fichier source et générer l'arbre de syntaxe correspondant (en utilisant une grammaire attribuée pour faire remonter l'information de typage). Il sera nécessaire de construire une table des symboles pour gérer le problème de portée des variables.
- générer en sortie une représentation graphique du fichier source au format `dot` de `graphviz`
- générer une erreur compréhensible si le programme est incorrect.

Constitution des groupes : groupes de 3 ou 4 étudiants.

Format de rendu : Une archive au format `Nom1_Nom2_Nom3.tar.gz` contenant les sources de votre programme, une série de fichiers sources de test un fichier `README` décrivant l'utilisation du programme, et un document au format `pdf` expliquant votre solution pour gérer le problème des instructions de retour absentes.

Remarque : Il vous est demandé d'apporter une attention particulière aux structures de données utilisées pour la représentation de l'arbre qui sera utilisé dans la phase de production de code.