

Mini-compileur C¹

première partie : analyse syntaxique

deuxième partie : typage

Le langage considéré dans ce TER est un sous-ensemble de C. Vos fichiers pourront être compilés aussi bien par votre compilateur que par tout autre compilateur C.

Notation Dans la suite, les symboles « \star », « $+$ » et « $?$ » indiquent la répétition du symbole les précédant (zéro fois ou plus pour « \star », zéro ou une fois pour « $?$ » et une fois ou plus pour « $+$ »). Attention de ne pas confondre « \star » et « $+$ » avec « $*$ » et « $+$ » qui sont des symboles du langage C.

1 Analyse syntaxique

1.1 Analyse lexicale

Espaces, tabulations et retour-chariots constituent les blancs. Les commentaires peuvent prendre deux formes :

- débutant par `/*` et s'étendant jusqu'à `*/`, et ne pouvant être imbriqués ;
- débutant par `//` et s'étendant jusqu'à la fin de la ligne (ou du fichier s'il s'agit de la dernière ligne).

De plus on ignorera les lignes suivantes :

```
#include <stdlib.h>
#include <stdio.h>
```

quelle que soit leur position dans le fichier.

Les identificateurs obéissent à l'expression régulière $\langle ident \rangle$ suivante :

$$\begin{aligned} \langle chiffr \rangle &::= 0-9 \\ \langle alpha \rangle &::= a-z \mid A-Z \\ \langle ident \rangle &::= (\langle alpha \rangle \mid _)(\langle alpha \rangle \mid _ \mid \langle chiffr \rangle)^* \end{aligned}$$

Les identificateurs suivants sont des mots clés :

```
char    double else    extern  for      if      int    long
return short sizeof struct unsigned void while
```

Les constantes entières obéissent à l'expression régulière $\langle entier \rangle$ suivante :

$$\langle entier \rangle ::= \langle chiffr \rangle + (u|U)?(l|L)?$$

1. Ce projet reprend certains aspects du projet de compilation donné il y a trois ans, ce dernier lui-même inspiré d'un projet similaire donné par Jean-Christophe FILLIATRE dans le cadre de son cours de compilation à l'École Normale Supérieure <http://www.lri.fr/~filliatr>.

Les constantes flottantes obéissent à l'expression régulière $\langle double \rangle$ suivante :

$$\langle double \rangle ::= ((\langle chiffre \rangle + . \langle chiffre \rangle \star | . \langle chiffre \rangle +)((e|E) - ? \langle chiffre \rangle +) ?$$

les littéraux caractères sont délimités par des *quotes* simples ('). Les caractères dont le code ASCII est celui du blanc (code ASCII décimal 32) ou supérieur sont représentés tels quels, à l'exception des caractères « ' », « " » et « \ ». Ces derniers peuvent être représentés en utilisant les séquences d'échappement « \' », « \" » et « \\ » respectivement. Il est aussi possible d'utiliser la séquence d'échappement « \xxx » où x est un chiffre hexadécimal (« 0-9 », « a-f » ou « A-F »). Dans ce cas, la séquence représente le caractère de code ASCII hexadécimal xx . Enfin, la séquence « \n » représente un retour chariot.

Enfin les chaînes de caractères littérales sont des séquences de caractères (tels que définis ci-dessus) délimitées par des *quotes* doubles : « " ». Les chaînes ne peuvent pas s'étendre sur plusieurs lignes. Dans la suite, les caractères et les chaînes de caractères sont représentés par les non-terminaux $\langle car \rangle$ et $\langle chaîne \rangle$ respectivement.

1.2 Analyse syntaxique

La grammaire des fichiers sources considérée est donnée figure 1. Le point d'entrée est le non-terminal $\langle fichier \rangle$.

Les associativités et précédences des divers opérateurs sont données par la table suivante, de la plus faible à la plus forte précedence :

opérateur	associativité	précédence
=	à droite	faible
	à gauche	
&&	à gauche	
== !=	à gauche	
< <= > >=	à gauche	↓
+ -	à gauche	
* / %	à gauche	
! ++ -- & * (unaire) + (unaire) - (unaire)	à droite	
) [-> .	à gauche	forte

1.3 Localisation des erreurs

Lorsqu'une erreur lexicale ou syntaxique est détectée par votre compilateur, elle doit être signalée le plus précisément possible, par sa nature et sa localisation dans le fichier source. On adoptera le format suivant pour cette signalisation :

```
File "test.c", line 4, characters 5-6:
syntax error
```

La première ligne du message d'erreur (« File "test.c" ... ») doit respecter *strictement* ce format. La seconde ligne est libre, vous pouvez vous en servir pour donner un message d'explication (éventuellement en français). Les localisations peuvent être obtenues grâce aux fonctions `Lexing.lexeme_start_p`, `Lexing.lexeme_end_p` pour l'analyseur lexical et `Parsing.rhs_start_pos`, `Parsing.rhs_end_pos` pour l'analyseur syntaxique.

$\langle \text{fichier} \rangle$	$::=$	$\langle \text{decl} \rangle^* \text{EOF}$
$\langle \text{decl} \rangle$	$::=$	$\langle \text{decl_vars} \rangle \mid \langle \text{decl_typ} \rangle \mid \langle \text{decl_fct} \rangle \mid \langle \text{decl_ext} \rangle$
$\langle \text{decl_vars} \rangle$	$::=$	$\langle \text{type} \rangle \langle \text{var} \rangle ;$
$\langle \text{decl_typ} \rangle$	$::=$	struct $\langle \text{ident} \rangle \{ \langle \text{decl_vars} \rangle^* \} ;$
$\langle \text{decl_fct} \rangle$	$::=$	$\langle \text{type} \rangle^* \langle \text{ident} \rangle (\langle \text{arguments} \rangle ?) \langle \text{bloc} \rangle$
$\langle \text{decl_ext} \rangle$	$::=$	extern $\langle \text{type} \rangle^* \langle \text{ident} \rangle (\langle \text{arguments} \rangle ?) ;$
$\langle \text{type} \rangle$	$::=$	void $\mid \langle \text{int_type} \rangle \mid \text{double} \mid \text{struct } \langle \text{ident} \rangle$
$\langle \text{int_type} \rangle$	$::=$	unsigned? (char \mid short \mid int \mid long)
$\langle \text{arguments} \rangle$	$::=$	$\langle \text{type} \rangle \langle \text{var} \rangle \mid \langle \text{type} \rangle \langle \text{var} \rangle , \langle \text{arguments} \rangle$
$\langle \text{var} \rangle$	$::=$	$\langle \text{ident} \rangle \mid * \langle \text{var} \rangle$
$\langle \text{expr} \rangle$	$::=$	$\langle \text{entier} \rangle \mid \langle \text{car} \rangle \mid \langle \text{chaîne} \rangle \mid \langle \text{double} \rangle$ $\mid \langle \text{ident} \rangle$ $\mid * \langle \text{expr} \rangle$ $\mid \langle \text{expr} \rangle [\langle \text{expr} \rangle]$ $\mid \langle \text{expr} \rangle . \langle \text{ident} \rangle$ $\mid \langle \text{expr} \rangle \rightarrow \langle \text{ident} \rangle$ $\mid \langle \text{expr} \rangle = \langle \text{expr} \rangle$ $\mid \langle \text{ident} \rangle (\langle \text{l_expr} \rangle ?)$ $\mid ++ \langle \text{expr} \rangle \mid -- \langle \text{expr} \rangle \mid \langle \text{expr} \rangle ++ \mid \langle \text{expr} \rangle --$ $\mid \& \langle \text{expr} \rangle \mid ! \langle \text{expr} \rangle \mid - \langle \text{expr} \rangle \mid + \langle \text{expr} \rangle$ $\mid \langle \text{expr} \rangle \langle \text{opérateur} \rangle \langle \text{expr} \rangle$ $\mid \text{sizeof } (\langle \text{cplx_type} \rangle)$ $\mid (\langle \text{cplx_type} \rangle) \langle \text{expr} \rangle$ $\mid (\langle \text{expr} \rangle)$
$\langle \text{opérateur} \rangle$	$::=$	== \mid != \mid < \mid <= \mid > \mid >= \mid + \mid - \mid * \mid / \mid % \mid && \mid
$\langle \text{cplx_type} \rangle$	$::=$	$\langle \text{type} \rangle^*$
$\langle \text{l_expr} \rangle$	$::=$	$\langle \text{expr} \rangle \mid \langle \text{expr} \rangle , \langle \text{l_expr} \rangle$
$\langle \text{instruction} \rangle$	$::=$; $\mid \langle \text{expr} \rangle ;$ $\mid \text{if } (\langle \text{expr} \rangle) \langle \text{instruction} \rangle$ $\mid \text{if } (\langle \text{expr} \rangle) \langle \text{instruction} \rangle \text{ else } \langle \text{instruction} \rangle$ $\mid \text{while } (\langle \text{expr} \rangle) \langle \text{instruction} \rangle$ $\mid \text{for } (\langle \text{l_expr} \rangle ? ; \langle \text{expr} \rangle ? ; \langle \text{l_expr} \rangle ?) \langle \text{instruction} \rangle$ $\mid \langle \text{bloc} \rangle$ $\mid \text{return } \langle \text{expr} \rangle ? ;$
$\langle \text{bloc} \rangle$	$::=$	$\{ \langle \text{decl_vars} \rangle^* \langle \text{instruction} \rangle^* \}$

FIGURE 1 – Grammaire des fichiers C

2 Typage

2.1 Types et environnements de typage

Dans tout ce qui suit, les expressions de types sont de la forme suivante :

$$\tau ::= \text{void} \mid \iota \mid \text{double} \mid \text{struct } id \mid \tau^* \mid \text{typenull}$$

$$\begin{aligned} \iota &::= \text{unsigned } \nu \mid \nu \\ \nu &::= \text{char} \mid \text{short} \mid \text{int} \mid \text{long} \end{aligned}$$

où *id* désigne un identificateur de structure. Le type `typenull` est un type particulier introduit pour pouvoir typer la constante 0 qui en C représente soit l'entier 0 soit le pointeur nul (la macro `NULL`, définie dans le fichier `stdlib.h` n'est en général qu'un alias pour `((void *)0)`). Il s'agit là d'une notation pour la syntaxe *abstraite* des expressions de types. On introduit la relation \approx sur les types comme la plus petite relation réflexive et symétrique telle que

$$\frac{\tau_1, \tau_2 \in \{\iota, \text{typenull}, \text{double}\}}{\tau_1 \approx \tau_2} \quad \frac{}{\text{typenull} \approx \tau*} \quad \frac{}{\text{void}* \approx \tau*}$$

Attention, la relation *n'est pas transitive*. En effet, on a $\text{int} \approx \text{typenull}$ (par la première règle) et $\text{typenull} \approx \text{char}*$ (par la deuxième règle) mais on n'a surtout pas $\text{int} \approx \text{char}*$. Un environnement de typage Γ est une suite de déclarations de variables de la forme τx , de déclarations de structures de la forme `struct` $S \{\tau_1 x_1 \cdots \tau_n x_n\}$, et de déclarations de profils de fonctions de la forme $\tau f(\tau_1, \dots, \tau_n)$. On notera `struct` $S \{\tau x\}$ pour indiquer que la structure S contient un champ x de type τ .

Dans la suite, on dira qu'un type est *numérique*, et on notera $\text{num}(\tau)$ si τ n'est compatible ni avec `void` ni avec un type de structure (autrement dit, τ doit être compatible avec `typenull`, un type d'entier ι `double` ou un type pointeur $\tau'*$) et on dira qu'un type est *arithmétique*, et on notera $\text{arith}(\tau)$ si τ est numérique et n'est pas un type pointeur.

On définit ensuite une fonction **rank** sur les types arithmétiques. Cette dernière sert à exprimer de manière compacte les conversions entre les différents types d'entiers (et elle représente, intuitivement le nombre de bits utilisés par le type, hors bit de signe).

$$\begin{array}{ll} \text{rank}(\text{char}) &= 7 & \text{rank}(\text{unsigned } \nu) &= 1 + \text{rank}(\nu) \\ \text{rank}(\text{short}) &= 15 & \text{rank}(\text{double}) &= 100 \\ \text{rank}(\text{int}) &= 31 & \text{rank}(\text{typenull}) &= 0 \\ \text{rank}(\text{long}) &= 63 \end{array}$$

On définit à partir de *rank*, la relation d'ordre partielle sur les types :

$$\tau_1 < \tau_2 \Leftrightarrow \text{rank}(\tau_1) < \text{rank}(\tau_2)$$

et on note $\tau_1 \leq \tau_2$ le fait que $\tau_1 < \tau_2$ ou $\tau_1 = \tau_2$ et on définit :

$$\begin{aligned} \max(\tau_1, \tau_2) &= \tau_2 \quad \text{si } \tau_1 < \tau_2 \\ \max(\tau_1, \tau_2) &= \tau_1 \quad \text{si } \tau_2 < \tau_1 \end{aligned}$$

2.2 Typage

Bonne formation des types On dit qu'un type τ est *bien formé* dans un environnement Γ , et on note $\Gamma \vdash \tau$ bf, si tous les identificateurs de structures apparaissant dans τ correspondent à des structures déclarées dans Γ .

Typage des expressions En C, il existe une classe d'expressions particulières nommées « *valeurs gauches* ». Intuitivement, les valeurs gauches sont les expressions que l'on a le droit de placer à *gauche* d'une affectation (d'où leur nom). À l'inverse d'autres langages, le fait d'être une valeur gauche en C dépend du type de l'expression. C'est pourquoi

on introduit simultanément le jugement $\Gamma \vdash e : \tau$ signifiant « dans l'environnement Γ , l'expression e est bien typée de type τ » et le jugement $\Gamma \vdash_l e : \tau$ signifiant « dans l'environnement Γ , l'expression e est une valeur gauche bien typée de type τ ». Ces jugements sont définis par les règles d'inférence suivantes :

(constantes)

$$\frac{}{\Gamma \vdash 0 : \text{typenull}} \quad \frac{c \text{ constante de type } \tau}{\Gamma \vdash c : \tau} \quad \frac{\Gamma \vdash \tau \text{ bf } \tau \neq \text{void}}{\Gamma \vdash \text{sizeof}(\tau) : \text{unsigned long}}$$

(affectations)

$$\frac{\Gamma \vdash_l e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \approx \tau_2}{\Gamma \vdash e_1 = e_2 : \tau_1}$$

$$\frac{\Gamma \vdash_l e : \tau \quad \text{num}(\tau) \quad op \in \{++, --\}}{\Gamma \vdash op \ e : \tau} \quad \frac{\Gamma \vdash_l e : \tau \quad \text{num}(\tau) \quad op \in \{++, --\}}{\Gamma \vdash e \ op : \tau}$$

(comparaisons)

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \text{num}(\tau_1) \quad \tau_1 \approx \tau_2 \quad op \in \{==, !=, <, <=, >, >=\}}{\Gamma \vdash e_1 \ op \ e_2 : \text{int}}$$

(arithmétique et logique)

$$\frac{\Gamma \vdash e : \tau \quad \text{arith}(\tau) \quad op \in \{+, -\}}{\Gamma \vdash op \ e : \tau} \quad \frac{\Gamma \vdash e : \tau \quad \text{num}(\tau)}{\Gamma \vdash !e : \text{int}}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \max(\tau_1, \tau_2) \leq \text{unsigned long}}{\Gamma \vdash e_1 \% e_2 : \max(\tau_1, \tau_2)}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \approx \tau_2 \quad \tau_1 \approx \text{double} \quad op \in \{+, -, *, /\}}{\Gamma \vdash e_1 \ op \ e_2 : \max(\tau_1, \tau_2)}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \approx \tau_2 \quad \tau_1 \approx \text{double} \quad op \in \{||, \&\&\}}{\Gamma \vdash e_1 \ op \ e_2 : \text{int}}$$

(arithmétique de pointeurs)

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \tau'_1 * \quad \max(\tau_2, \text{int}) \leq \text{unsigned long} \quad op \in \{+, -\}}{\Gamma \vdash e_1 \ op \ e_2 : \tau'_1 *}$$

$$\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2 \quad \tau_1 \equiv \tau'_2 * \quad \tau_2 \equiv \tau'_2 *}{\Gamma \vdash e_1 - e_2 : \text{long}} \quad \frac{\Gamma \vdash e_2 + e_1 : \tau}{\Gamma \vdash e_1 + e_2 : \tau} \quad \frac{\Gamma \vdash e : \tau' \quad \text{num}(\tau) \quad \text{num}(\tau') \quad \Gamma \vdash \tau \text{ bf}}{\Gamma \vdash (\tau)e : \tau}$$

(appel de fonction)

$$\frac{\Gamma \vdash e_i : \tau_i \quad \tau \ f(\tau'_1, \dots, \tau'_n) \in \Gamma \quad \tau_i \approx \tau'_i}{\Gamma \vdash f(e_1, \dots, e_n) : \tau}$$

(accès)

$$\frac{\Gamma \vdash e : \text{struct } S \quad \text{struct } S \{ \tau x \} \in \Gamma}{\Gamma \vdash e.x : \tau}$$

(adresse)

$$\frac{\Gamma \vdash_l e : \tau}{\Gamma \vdash \&e : \tau^*}$$

(valeurs gauches)

$$\frac{\frac{\Gamma \vdash_l e : \tau}{\Gamma \vdash e : \tau} \quad \frac{\tau x \in \Gamma}{\Gamma \vdash_l x : \tau} \quad \frac{\Gamma \vdash e : \tau^*}{\Gamma \vdash_l *e : \tau}}{\Gamma \vdash_l e : \text{struct } S \quad \text{struct } S \{ \tau x \} \in \Gamma} \Gamma \vdash_l e.x : \tau$$

Typage des instructions On introduit le jugement $\Gamma, \tau_0 \vdash_i i$ signifiant « dans l'environnement Γ , l'instruction i est bien typée, pour un type de retour τ_0 ». Intuitivement, τ_0 représente le type de retour de la fonction dans la quelle se trouve l'instruction i . Ce jugement est établi par les règles d'inférence suivantes :

$$\begin{array}{c} \frac{}{\Gamma, \tau_0 \vdash_i ;} \quad \frac{\Gamma, \tau_0 \vdash e : \tau}{\Gamma, \tau_0 \vdash_i e;} \quad \frac{}{\Gamma, \text{void} \vdash_i \text{return};} \quad \frac{\Gamma \vdash e : \tau_0}{\Gamma, \tau_0 \vdash_i \text{return } e;} \\[10pt] \frac{\Gamma \vdash e : \tau \quad \text{num}(\tau) \quad \Gamma, \tau_0 \vdash_i i_1 \quad \Gamma, \tau_0 \vdash_i i_2}{\Gamma, \tau_0 \vdash_i \text{if } (e) i_1 \text{ else } i_2} \\[10pt] \frac{\Gamma, \tau_0 \vdash_i i_1 \quad \Gamma \vdash e : \tau \quad \text{num}(\tau) \quad \Gamma, \tau_0 \vdash_i i_2 \quad \Gamma, \tau_0 \vdash_i i_3}{\Gamma, \tau_0 \vdash_i \text{for}(i_1; e; i_2) i_3} \\[10pt] \frac{\forall j \leq k, \Gamma \vdash \tau_j \text{ bf } \tau_j \neq \text{void} \quad \forall j \leq n, \{ \tau_1 x_1, \dots, \tau_k x_k \} \cup \Gamma \vdash_i i_j}{\Gamma, \tau_0 \vdash_i \{ \tau_1 x_1 \dots \tau_k x_k; i_1 \dots i_n \}} \end{array}$$

Cette dernière règle signifie que pour typer un *bloc* constitué de k déclarations de variables (locales au bloc) et de n instructions, on vérifie d'abord la bonne formation des déclarations puis on type chacune des instructions dans l'environnement augmenté des nouvelles déclarations. On rappelle aussi que le *parsing* a supprimé des instructions superflues (**while**, **if** sans **else**, ...).

Typage des fichiers On rappelle qu'un fichier est une liste de déclarations. On introduit le jugement « $\Gamma \vdash d \rightarrow \Gamma'$ » qui signifie « dans l'environnement Γ , la déclaration d est bien formée et produit un environnement Γ' ». Ce jugement est dérivable grâce aux règles suivantes :

Déclarations de variables (globales)

$$\frac{\Gamma \vdash \tau \text{ bf } \tau \neq \text{void}}{\Gamma \vdash \tau x \rightarrow \{ \tau x \} \cup \Gamma}$$

Déclarations de structures

$$\frac{\Gamma, \mathbf{struct\ id\ } \{\tau_1\ x_1 \cdots \tau_n\ x_n\} \vdash \tau_i\ \mathbf{bf}}{\Gamma \vdash \mathbf{struct\ id\ } \{\tau_1\ x_1; \cdots \tau_n\ x_n; \} \rightarrow \{\mathbf{struct\ id\ } \{\tau_1\ x_1 \cdots \tau_n\ x_n\}\} \cup \Gamma}$$

On vérifiera d'autre part que les types de champs τ_i ne font référence à la structure *id* elle-même que sous un pointeur.

Déclarations de fonctions et de fonction externes

$$\frac{\Gamma \vdash \tau_i\ \mathbf{bf} \quad \forall i > 0, \tau_i \neq \mathbf{void} \quad \{\tau_0\ f(\tau_1, \dots, \tau_n), \tau_1\ x_1, \dots, \tau_n\ x_n\} \cup \Gamma, \tau_0 \vdash_i b}{\Gamma \vdash \tau_0\ f(\tau_1\ x_1, \dots, \tau_n\ x_n)\ b \rightarrow \{\tau_0\ f(\tau_1, \dots, \tau_n)\} \cup \Gamma}$$

$$\frac{\Gamma \vdash \tau_i\ \mathbf{bf} \quad \forall i > 0, \tau_i \neq \mathbf{void}}{\Gamma \vdash \tau_0\ f(\tau_1\ x_1, \dots, \tau_n\ x_n)\ b \rightarrow \{\tau_0\ f(\tau_1, \dots, \tau_n)\} \cup \Gamma}$$

On remarque que le prototype d'une fonction est ajouté à l'environnement pour le typage de cette dernière, cela dans le but de typer les fonctions récursives.

Fichiers On introduit finalement le jugement $\Gamma \vdash_f d_1 \cdots d_n$ signifiant « dans l'environnement Γ le fichier constitué par la suite de déclarations d_1, \dots, d_n est bien formé ». Le typage d'un fichier consiste à typer successivement les déclarations dans le contexte étendu par chaque nouvelle déclaration, d'où les règles :

$$\frac{}{\Gamma \vdash_f \emptyset} \quad \frac{\Gamma \vdash d_1 \rightarrow \Gamma' \quad \Gamma' \vdash_f d_2 \cdots d_n}{\Gamma \vdash_f d_1\ d_2 \cdots d_n}$$

Règles d'unicité Enfin, on vérifiera l'unicité :

- des identificateurs de structures sur l'ensemble du fichier ;
- des champs de structure à l'intérieur d'une *même* structure ;
- des symboles (variables *globales* et fonctions) sur l'ensemble du fichier.

Fonction principale On vérifiera la présence d'une fonction **main** avec l'un des deux prototypes suivants :

```
int main();
int main(int argc, char** argv);
```

dans le fichier.