

# Mini-compileur C<sup>1</sup>

première partie : analyse syntaxique  
deuxième partie : typage  
troisième partie : production de code

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 :

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

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 :

$$\langle entier \rangle ::= \langle chiffre \rangle + (\text{u|U}) ? (\text{l|L}) ?$$


---

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 FILLIÂTRE 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.

$\langle \text{fichier} \rangle$	$::=$	$\langle \text{decl} \rangle^* \text{EOF}$
$\langle \text{decl} \rangle$	$::=$	$\langle \text{decl\_vars} \rangle \mid \langle \text{decl\_typ} \rangle \mid \langle \text{decl\_fct} \rangle \mid \langle \text{decl\_ext} \rangle$
$\langle \text{decl\_vars} \rangle$	$::=$	$\langle \text{type} \rangle \langle \text{var} \rangle ;$
$\langle \text{decl\_typ} \rangle$	$::=$	<b>struct</b> $\langle \text{ident} \rangle \{ \langle \text{decl\_vars} \rangle^* \} ;$
$\langle \text{decl\_fct} \rangle$	$::=$	$\langle \text{type} \rangle^* \langle \text{ident} \rangle ( \langle \text{arguments} \rangle ? ) \langle \text{bloc} \rangle$
$\langle \text{decl\_ext} \rangle$	$::=$	<b>extern</b> $\langle \text{type} \rangle^* \langle \text{ident} \rangle ( \langle \text{arguments} \rangle ? ) ;$
$\langle \text{type} \rangle$	$::=$	<b>void</b> $\mid \langle \text{int\_type} \rangle \mid \text{double} \mid \text{struct } \langle \text{ident} \rangle$
$\langle \text{int\_type} \rangle$	$::=$	<b>unsigned?</b> ( <b>char</b> $\mid$ <b>short</b> $\mid$ <b>int</b> $\mid$ <b>long</b> )
$\langle \text{arguments} \rangle$	$::=$	$\langle \text{type} \rangle \langle \text{var} \rangle \mid \langle \text{type} \rangle \langle \text{var} \rangle , \langle \text{arguments} \rangle$
$\langle \text{var} \rangle$	$::=$	$\langle \text{ident} \rangle \mid * \langle \text{var} \rangle$
$\langle \text{expr} \rangle$	$::=$	$\langle \text{entier} \rangle \mid \langle \text{car} \rangle \mid \langle \text{chaîne} \rangle \mid \langle \text{double} \rangle$ $\mid \langle \text{ident} \rangle$ $\mid * \langle \text{expr} \rangle$ $\mid \langle \text{expr} \rangle [ \langle \text{expr} \rangle ]$ $\mid \langle \text{expr} \rangle . \langle \text{ident} \rangle$ $\mid \langle \text{expr} \rangle \rightarrow \langle \text{ident} \rangle$ $\mid \langle \text{expr} \rangle = \langle \text{expr} \rangle$ $\mid \langle \text{ident} \rangle ( \langle \text{l\_expr} \rangle ? )$ $\mid ++ \langle \text{expr} \rangle \mid -- \langle \text{expr} \rangle \mid \langle \text{expr} \rangle ++ \mid \langle \text{expr} \rangle --$ $\mid \& \langle \text{expr} \rangle \mid ! \langle \text{expr} \rangle \mid - \langle \text{expr} \rangle \mid + \langle \text{expr} \rangle$ $\mid \langle \text{expr} \rangle \langle \text{opérateur} \rangle \langle \text{expr} \rangle$ $\mid \text{sizeof } ( \langle \text{cplx\_type} \rangle )$ $\mid ( \langle \text{cplx\_type} \rangle ) \langle \text{expr} \rangle$ $\mid ( \langle \text{expr} \rangle )$
$\langle \text{opérateur} \rangle$	$::=$	<b>==</b> $\mid$ <b>!=</b> $\mid$ <b>&lt;</b> $\mid$ <b>&lt;=</b> $\mid$ <b>&gt;</b> $\mid$ <b>&gt;=</b> $\mid$ <b>+</b> $\mid$ <b>-</b> $\mid$ <b>*</b> $\mid$ <b>/</b> $\mid$ <b>%</b> $\mid$ <b>&amp;&amp;</b> $\mid$ <b>  </b>
$\langle \text{cplx\_type} \rangle$	$::=$	$\langle \text{type} \rangle^*$
$\langle \text{l\_expr} \rangle$	$::=$	$\langle \text{expr} \rangle \mid \langle \text{expr} \rangle , \langle \text{l\_expr} \rangle$
$\langle \text{instruction} \rangle$	$::=$	$;$ $\mid \langle \text{expr} \rangle ;$ $\mid \text{if } ( \langle \text{expr} \rangle ) \langle \text{instruction} \rangle$ $\mid \text{if } ( \langle \text{expr} \rangle ) \langle \text{instruction} \rangle \text{ else } \langle \text{instruction} \rangle$ $\mid \text{while } ( \langle \text{expr} \rangle ) \langle \text{instruction} \rangle$ $\mid \text{for } ( \langle \text{l\_expr} \rangle ? ; \langle \text{expr} \rangle ? ; \langle \text{l\_expr} \rangle ? ) \langle \text{instruction} \rangle$ $\mid \langle \text{bloc} \rangle$ $\mid \text{return } \langle \text{expr} \rangle ? ;$
$\langle \text{bloc} \rangle$	$::=$	$\{ \langle \text{decl\_vars} \rangle^* \langle \text{instruction} \rangle^* \}$

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  $\Gamma$  est une suite de déclarations de variables de la forme  $\tau \ x$ , 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  $\tau$ .

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

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

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  $\tau$  est *bien formé* dans un environnement  $\Gamma$ , et on note  $\Gamma \vdash \tau \text{ bf}$ , si tous les identificateurs de structures apparaissant dans  $\tau$  correspondent à des structures déclarées dans  $\Gamma$ .

**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  $\Gamma$ , l'expression  $e$  est bien typée de type  $\tau$  » et le jugement  $\Gamma \vdash_l e : \tau$  signifiant « dans l'environnement  $\Gamma$ , l'expression  $e$  est une valeur gauche bien typée de type  $\tau$  ». 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)

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

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

(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  $\Gamma$ , l'instruction  $i$  est bien typée, pour un type de retour  $\tau_0$  ». Intuitivement,  $\tau_0$  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 :

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

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  $\Gamma$ , la déclaration  $d$  est bien formée et produit un environnement  $\Gamma'$  ». 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  $\tau_i$  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  $\Gamma$  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.

## 3 Production de code

L'objectif est de réaliser un compilateur simple mais correct. En particulier, on ne cherchera pas à faire d'allocation de registres mais on se contentera d'utiliser la pile pour stocker les éventuels calculs intermédiaires. Bien entendu, il est possible, et même souhaitable, d'utiliser localement les registres de x86-64.

Les difficultés liées à la production de code dans ce projet sont les suivantes :

- les données n'ont pas toutes la même taille; en particulier, le compilateur doit savoir calculer la taille de la représentation mémoire de chaque type.
- certaines valeurs doivent être correctement *alignées*, c'est-à-dire stockées à des adresses multiples de 8.
- les structures sont affectées, passées en argument et renvoyées comme résultat *par valeur*, ce qui nécessite de copier tout un bloc d'octets.
- Il y a beaucoup de types différents et de combinaisons à gérer! (addition d'un `char` avec un double, d'un `unsigned long` avec un `int`, ...)

### 3.1 Représentation des valeurs

On commence par définir la notion d'*alignement*. On dit qu'une donnée est alignée sur  $n$  octets si l'adresse à laquelle elle se trouve est un multiple de  $n$  octets. Même si l'architecture x86-64 supporte les accès mémoire non alignés (à l'inverse d'ARM et MIPS), la norme C impose certains alignements. Afin de simplifier la disposition des objets en mémoire on appliquera la règle qu'un type simple (entier, double ou pointeur) doit toujours être aligné sur sa taille. (*i.e* les `long` et pointeurs se trouvent à une adresse multiple de 8, les `int` de 4, les `short` de 2).

Une valeur de type `char` sera stockée sur 8 bits non signés et une valeur de type `int` ou de type pointeur sur 64 bits signés. Une structure est représentée comme la suite ordonnée de ses champs. Les champs sont éventuellement séparés par des octets inutilisés, dits de remplissage (*padding* en anglais), lorsque l'alignement des champs l'exige. Ainsi la structure

```
struct S { int a; char b; char c; char *p; };
```

doit être alignée sur 8 octets et sera représentée sur 16 octets, dont 2 octets de remplissage marqués X, de la manière suivante :

a	a	a	a	b	c	X	X	p	p	p	p	p	p	p	p
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

Une structure doit éventuellement comporter des octets terminaux de remplissage pour que la structure suivante dans un tableau soit alignée. Ainsi la structure

```
struct T { int x; char y; };
```

occupe 8 octets (dont 3 de remplissage après `y`). Ceci permet notamment d'allouer un « tableau » de `n` telles structures avec l'idiome

```
struct T *a = malloc(n * sizeof(struct T));
```

tout en garantissant l'alignement de chacune d'elles.

La construction `sizeof( $\tau$ )` doit renvoyer le nombre d'octets occupés par la représentation du type  $\tau$ . En particulier on a `sizeof(struct S) = 16` et `sizeof(struct T) = 8` pour les deux structures ci-dessus.



## 3.2 Schéma de compilation

Les variables globales seront allouées sur le segment de données. Les variables locales seront allouées sur la pile (dans le tableau d'activation). Les arguments et résultat d'une fonction seront passés sur la pile (les données étant de taille variable, c'est le plus simple). On pourra adopter le schéma suivant :

	⋮	
	résultat	
	argument $n$	
	⋮	
	argument 1	
appelant	adresse de retour	
appelé	ancien %rbp	←%rbp
	locale 1	
	⋮	
	locale $m$	
	calculs	
	⋮	
	calculs	←%esp
	⋮	

Pour la copie de structure (passage d'argument, valeur de retour, affectation), on pourra écrire une fonction de type `memmove`, par exemple directement en assembleur.

## 3.3 Retour sur le typage

Afin de faciliter (grandement) la production de code, il est utile de revenir sur la passe de typage pour stocker les informations suivantes :

- pour chaque fonction gardée dans l'environnement global `fun_env`, il faut garder en plus un booléen indiquant si c'est une fonction `extern` ou non. Les déclarations `extern` nous permettent de déclarer une fonction de la bibliothèque standard C que `gcc` saura résoudre. Par exemple :

```
//affiche un caractère sur le terminal
extern int putchar(int c);
//affiche une chaine suivie de \n dans le terminal
extern int puts(char * s);
//renvoie un pointeur vers une zone mémoire de taille l
void * malloc(unsigned long l);
//libere une zone mémoire renvoyée par malloc
void free(void * p);

int main() {
    //le code peut utiliser toutes les fonctions
    //ci-dessus !
}
```

- pour tous les endroits où la relation de compatibilité a été utilisée, insérer un **cast** approprié. Le cas le plus complexe est celui des opérations binaires, traité ci dessous. Concrètement, étant donné un appel de fonction de la forme :  $f(e_1, \dots, e_n)$  où les arguments attendus sont de type  $\tau_1, \dots, \tau_n$ , on le transformera en un appel de fonction  $f((\tau_1)e_1, \dots, (\tau_n)e_n)$ , ce qui permet par exemple de donner des **int** à une fonction qui attend des **long** en argument. La production de code s'en trouvera alors simplifiée car le code de conversion sera écrit une fois pour toute pour le *cast* et vous pourrez supposer que les fonctions ont toujours des arguments du bon type en paramètre. De même une affectation  $e_1 = e_2$  où  $e_1$  est une valeur gauche de type  $\tau_1$ , sera remplacée par une affectation  $e_1 = (\tau_1)e_2$

**Opérateurs binaires.** Les opérateurs binaires peuvent impliquer plusieurs combinaisons de types différents. La norme C est particulièrement sordide de ce point de vue là. Une fois simplifiée (pour l'architecture qui nous interesse) les conversions peuvent être exprimées comme suit. Soit une expression *bien typée*  $e_1 o e_2$  où  $o$  est un opérateur binaire (addition, soustraction, comparaison, ...). On appelle  $\tau_1$  et  $\tau_2$  les types de  $e_1$  et  $e_2$ .

- si  $\text{arith}(\tau_1)$  et  $\text{arith}(\tau_2)$  et  $\tau_1 \neq \tau_2$  (i.e. se sont tous les deux des types doubles ou entiers, pas de pointeur)
  - si  $\tau_1$  (resp.  $\tau_2$ ) est **double**,  $e_2$  est remplacé par  $(\text{double})e_2$  (resp.  $(\text{double})e_1$ )
  - sinon si  $\tau_1 \leq \text{int}$ , alors  $e_1$  est remplacé par  $e'_1 = (\text{int})e_1$  (idem pour  $e_2$ ).
  - On appelle  $\tau'_1$  et  $\tau'_2$  les nouveaux types de  $e'_1$  et  $e'_2$  (qui sont au moins **(int)**).
    - si  $\tau'_1$  vaut **unsigned long** alors  $e'_2$  est remplacé par  $(\text{unsigned long})e'_2$  (idem pour  $e'_1$ )
    - sinon, si  $\tau'_1$  vaut **long** alors  $e'_2$  est remplacé par  $(\text{long})e'_2$  (idem pour  $e'_1$ )
    - sinon, si  $\tau'_1$  vaut **unsigned int** alors  $e'_2$  est remplacé par  $(\text{unsigned int})e'_2$  (idem pour  $e'_1$ )
    - sinon pas de conversion
- sinon, pas de conversion

Pour donner un exemple, avec les définitions suivantes :

```
long x;
char y;
double z;
```

l'expression  $x * y$  est équivalent à  $x * ((\text{long}) ((\text{int}) y))$  et l'expression  $y - z$  est équivalente à  $((\text{double}) y) - z$ .

### 3.4 Spécificité de l'architecture x86-64

**Appels de fonction externes.** L'ABI x86-64 n'utilise pas (que) la pile pour passer des paramètres lors d'appels de fonctions, mais plutôt les registres **rdi**, **rsi**, **rdx**, **rcx**, **r8** et **r9** dans cet ordre, pour les 6 premiers arguments entiers et rendent leur résultat (s'il est entier) dans **rax**. Afin de pouvoir appeler des fonctions externes, il faudra donc compiler spécialement leur appel et retour en utilisant ces registres.

**Registre spéciaux** certains registres ont une utilisation particulière. Le registre **%esp** contient l'adresse de la dernière valeur posée sur la pile et est modifié par les instructions

`pushq` et `popq`. On peut cependant le manipuler comme tout autre registre. Le registre `%ebp` contient le pointeur vers le *tableau d'activation*. Il doit être correctement sauvegardé et initialisé lors de l'entrée dans une fonction. L'accès aux paramètres et aux variables locales se fait alors en additionnant leur décalage par rapport à `%ebp`.

Il faudra en particulier tenir compte de cela pour compiler correctement le retour de la fonction `main`. Il ne faudra pas écrire son résultat sur la pile, mais plutôt le placer dans le registre `rax`. Ce dernier sera le code de sortie de votre programme.

On lira aussi de manière attentive les ressources disponibles sur la page du cours ainsi que la documentation du module `Amd64` fourni, qui permet l'écriture de code assembleur en OCaml.

### 3.5 Test des programmes

Une fois le fichier assembleur (extension `.s`) généré par votre compilateur, vous pouvez le compiler avec :

```
gcc -o prog.exe -g prog.s
```

Vous pouvez ensuite l'exécuter directement sur votre machine (`./prog.exe`) ou l'exécuter à travers le debugger graphique `nemiver` installé au PUIO.

## 4 Rendu

La date final du rendu de rapport pour le TER est fixée au :

dimanche 23 avril, 2017, 21h00

via le formulaire de rendu en ligne sur la page du cours. Les soutenances auront lieu la semaine suivante (planning à venir).

### 4.1 Rapport et code

Un rapport de 15 pages maximum (police 11 points) au format PDF ou `odt` doit être rendu. Ce dernier doit être synthétique (pas besoin de 2 pages d'introduction ou conclusions pour un rapport de ce type). Il doit comporter deux parties, d'égale importance.

#### 4.1.1 Partie conception du compilateur

Vous devez y décrire votre code, en particulier tous les fichiers de votre projet ainsi que leur dépendances (e.g. `typing.ml` dépend de `ast.mli`). Pour chacune des quatre phases de compilation (analyse lexicale, analyse syntaxique, typage et production de code), vous choisirez un exemple de code C et indiquerez comment il est traité par la phase en question dans votre compilateur (évidemment vous devez prendre un exemple non trivial). Vous vous servirez de cet exemple pour décrire *succinctement l'organisation du code* (par exemple les fonctions principales du module *typing* pour le typage) ainsi que les structures de données utilisées.

### 4.1.2 Partie recherche

Vous devez choisir l'un des concepts proposé en fin d'énoncé et le décrire précisément via une recherche bibliographique (attention, un copier collé de wikipedia ne suffit pas et vous prendrez soin de citer vos sources). Vous illustrez ce concept en fournissant quelques scénarios d'utilisation permettant de le « tester » avec le compilateur GCC. Vous devrez aussi en implanter une version simple dans votre compilateur.

### 4.1.3 Rendu du code

Vous devez rendre le code. Ce dernier sera évalué sur sa forme (commentaires, indentation, élégance générale et efficacité) ainsi que sur son comportement. Pour la partie typage et analyse syntaxique, votre compilateur doit implanter fidèlement l'énoncé. Pour la partie production de code, les structures ne seront pas testée (i.e. on ne testera pas le code produit pour des fichiers sources contenant des structures. Cependant, votre compilateur ne doit pas planter dans ce cas – vous pouvez par exemple générer des `nop` pour tout code utilisant des structures).

### 4.1.4 Soutenance

Les soutenances durent 10 minutes par étudiant (donc 20 minutes pour un binôme). Vous consacrerez 1/3 du temps expliquer ce que contient votre compilateur (si possible avec une démo) et 2/3 du temps sur le concept de recherche.

### 4.1.5 Concepts

**Profilage de code** . Expliquer en quoi consiste le profilage de code, les différents outils (sous Unix ou autre) permettant de profiler du code (code C en particulier, mais d'autres exemples par exemple pour du code Java ou Javascript sont les bienvenus), quelles métriques le profilage de code permet de collecter et la manière de les analyser. Votre compilateur devra avoir un mode profilage permettant de profiler le code généré.

**Optimisations simples** . Expliquer en quoi consiste l'*inlining* de fonctions, la technique du *constant folding* et la technique du *constant propagation*. Montrez (en utilisant GCC) comment ces différentes optimisations sont complémentaires. Vous expliquerez aussi ce qui peut bloquer ces optimisations (en particulier dans les langages de haut niveau). Vous implémenterez soit un *inliner* simple, soit un évaluateur partiel d'expressions arithmétiques.

**Test boîte blanche et couverture de code** . Expliquer en quoi consiste le test « boîte blanche » d'un logiciel (et toutes les techniques associées). Votre compilateur devra posséder un mode de génération de code où le programme exécuté affiche, en fin d'exécution la liste de tous les blocs de code non exécutés.

**Allocateur de mémoire** . Expliquer comment la mémoire est gérée sur un système d'exploitation moderne (mémoire virtuelle et *paging*) puis détailler les différentes manières d'implémenter `malloc` et `free`, avec les compromis que fait chaque méthode et

leur complexité. Votre compilateur doit implémenter (en assembleur) `malloc` et `free` en réutilisant l'appel système `sbrk`.