

Le langage C, et sa sémantique

“C is quirky, flawed and an enormous success.”

Dennis Ritchie, *The Development of the C Language*, ACM SIGPLAN Notices 28(3):201–208, mars 1993.

Une vision simplifiée de C: 1. les expressions

Les expressions calculent des **valeurs**:

$e ::= x$	variables
1 2 ... -1 -2 ... 'a' 'b' ...	constantes
$e+e$ $e-e$ $e*e$ e/e $e\%e$ $-e$	arithmétique
$e\&e$ $e e$ $e^{\wedge}e$ $e>>e$ $e<<e$ $\sim e$	opérations bit à bit
$e==e$ $e<e$ $e<=e$ $e>e$ $e>=e$	comparaisons
$e\&\&e$ $e e$ $!e$	logique
$*e$ $e[e]$ $\&e$	pointeurs, tableaux
$e(e, \dots, e)$	appel de fonction

Une vision simplifiée de C: 2. les commandes (“statement”)

Les commandes **modifient la mémoire** et forment une **relation de transition**:

<code>c ::= e=e;</code>	affectation (attn! le test c'est ==!)
<code>i</code>	ne fait rien
<code>c c...c</code>	séquence
<code>{<déclarations de variables> c}</code>	bloc
<code>if (e) c [else c]</code>	conditionnelle (test)
<code>while (e) c</code>	boucle while
<code>do c while (e)</code>	boucle repeat
<code>for ([e];[e];[e]) c</code>	boucle for
<code>switch (e) {</code>	analyse de cas
<code>(case e: c break;)*</code>	
<code>[default: c break;] }</code>	
<code>return [e];</code>	retour de fonction

Une vision simplifiée de C: 3. les fonctions

Un programme C est une collection de déclarations de **variables** globales et de **fonctions** (\sim procédures).

Ces déclarations sont regroupées en fichiers (“modules”)

\langle nom-de-fichier \rangle .c, qui seront compilés sous forme de fichiers **objet**

\langle nom-de-fichier \rangle .o.

Chaque module peut référencer une variable ou une fonction d’un autre fichier en utilisant le mot-clé `extern`. Le programme démarre toujours en appelant la fonction `main`.

Sémantique simplifiée de C: 1. expressions

On va décrire la sémantique $\llbracket e \rrbracket \rho$ des expressions e dans un **environnement** ρ ; ρ est une fonction qui à chaque variable x associe sa valeur.

$$\llbracket x \rrbracket \rho = \rho(x) \quad \text{variables}$$

$$\llbracket c \rrbracket \rho = c \quad \text{constantes}$$

Note: les constantes caractères (ex: ' a ') représentent leur code ASCII (ex: ' a ' = 97 = 0x61).

Sémantique simplifiée de C: 1.1. arithmétique

$$\llbracket e_1 + e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho + \llbracket e_2 \rrbracket \rho \bmod (\text{maxint} + 1)$$

$$\llbracket e_1 - e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho - \llbracket e_2 \rrbracket \rho \bmod (\text{maxint} + 1)$$

$$\llbracket e_1 * e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho \times \llbracket e_2 \rrbracket \rho \bmod (\text{maxint} + 1)$$

$$\llbracket e_1 / e_2 \rrbracket \rho = \lfloor \llbracket e_1 \rrbracket \rho / \llbracket e_2 \rrbracket \rho \rfloor \quad (\text{quotient})$$

$$\llbracket e_1 \% e_2 \rrbracket \rho = \text{reste de la division entière ci-dessus}$$

$$\llbracket -e \rrbracket \rho = - \llbracket e \rrbracket \rho \bmod (\text{maxint} + 1)$$

Sémantique simplifiée de C: 1.2. opérations bit à bit

$\llbracket e_1 \& e_2 \rrbracket \rho =$ et bit à bit de $\llbracket e_1 \rrbracket \rho$ et $\llbracket e_2 \rrbracket \rho$

$\llbracket e_1 | e_2 \rrbracket \rho =$ ou —

$\llbracket e_1 \wedge e_2 \rrbracket \rho =$ ou exclusif —

$\llbracket \sim e \rrbracket \rho =$ non —

$\llbracket e_1 \gg e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho / 2^{\llbracket e_2 \rrbracket \rho}$

$\llbracket e_1 \ll e_2 \rrbracket \rho = \llbracket e_1 \rrbracket \rho * 2^{\llbracket e_2 \rrbracket \rho} \bmod (maxint + 1)$

Sémantique simplifiée de C: 1.3. comparaisons

$$\llbracket e_1 == e_2 \rrbracket \rho = \delta(\llbracket e_1 \rrbracket \rho = \llbracket e_2 \rrbracket \rho)$$

$$\llbracket e_1 < e_2 \rrbracket \rho = \delta(\llbracket e_1 \rrbracket \rho < \llbracket e_2 \rrbracket \rho)$$

$$\llbracket e_1 \leq e_2 \rrbracket \rho = \delta(\llbracket e_1 \rrbracket \rho \leq \llbracket e_2 \rrbracket \rho)$$

$$\llbracket e_1 > e_2 \rrbracket \rho = \delta(\llbracket e_1 \rrbracket \rho > \llbracket e_2 \rrbracket \rho)$$

$$\llbracket e_1 \geq e_2 \rrbracket \rho = \delta(\llbracket e_1 \rrbracket \rho \geq \llbracket e_2 \rrbracket \rho)$$

où $\delta(\text{vrai}) = 1$, $\delta(\text{faux}) = 0$.

Sémantique simplifiée de C: 1.4. logique

$$\begin{aligned} \llbracket e_1 \&\& e_2 \rrbracket \rho &= \begin{cases} 0 & \text{si } \llbracket e_1 \rrbracket \rho = 0 \\ \llbracket e_2 \rrbracket \rho & \text{sinon} \end{cases} \\ \llbracket e_1 \|\| e_2 \rrbracket \rho &= \begin{cases} \llbracket e_1 \rrbracket \rho & \text{si non nul} \\ \llbracket e_2 \rrbracket \rho & \text{sinon} \end{cases} \\ \llbracket !e \rrbracket \rho &= \begin{cases} 1 & \text{si } \llbracket e \rrbracket \rho = 0 \\ 0 & \text{sinon} \end{cases} \end{aligned}$$

Sémantique simplifiée de C: 1.5. pointeurs, tableaux

Idée: $*e$ retourne la valeur stockée à l'adresse [qui est la valeur de] e .

Corriger la sémantique! Elle doit prendre un deuxième argument

$\sigma : \mathbb{N} \rightarrow \mathbb{N}$ représentant la mémoire (comme en assembleur).

$$\llbracket *e \rrbracket \rho \sigma = \sigma(\llbracket e \rrbracket \rho \sigma)$$

$$\llbracket e_1[e_2] \rrbracket \rho \sigma = \sigma(\llbracket e_1 \rrbracket \rho \sigma + N \times \llbracket e_2 \rrbracket \rho \sigma)$$

où N est la taille d'un élément quelconque du tableau e_1 . Par exemple, si e_1 est un `int []` (tableau d'entiers), alors $N = \text{sizeof}(\text{int}) = 4$ (typiquement).

Note: en C, une variable tableau (e.g. `int []`) est de type pointeur (e.g. `int *`); $e[0]$ est synonyme de $*e$; $e_1[e_2]$ est synonyme de $*(e_1 + e_2)$ [modifier la sémantique de $+$ lorsqu'on additionne un pointeur et un entier!]

Sémantique simplifiée de C: 1.6. opérateur &

Slogan: $\&e$ est l'**adresse** de e .

On doit avoir: $*(\&e) = e$, $\&(*e) = e$.

Note: si $\&x$ a un sens, par contre $\&3$ n'en a pas. Les expressions pour lesquelles il est sensé de calculer leur adresse sont appelées les **lvalues**.

$lval$	$::=$	x	variables
		$*e$ $e[e]$	pointeurs, tableaux

Note: les affectations $e_1 = e_2$ n'ont de sens que lorsque e_1 est une lvalue; le résultat est alors de stocker la valeur de e_2 à l'adresse [valeur de] e_1 .

Sémantique simplifiée de C: 1.7. appel de fonction

Idée: $e(e_1, \dots, e_n)$ évalue e, e_1, \dots, e_n ; e doit être un nom de fonction f (ou un pointeur $*x$, où $x = \&f$), qui est appelée avec paramètres [les valeurs de] e_1, \dots, e_n ; la valeur retournée (par `return`) est celle de $e(e_1, \dots, e_n)$.

Pb.: la fonction f peut faire des **effets de bord** (modifications de la mémoire, entrées-sorties, etc.), qui ne sont pas prises en compte par la sémantique.

⇒ compliquer la sémantique des expressions. On utilise des **jugements**:

$$\sigma \vdash e \Rightarrow V, \sigma'$$

où V est une valeur, σ et σ' des mémoires.

Note: on n'a plus besoin de ρ : pourquoi?

Exercice: réécrire la sémantique des expressions...

Approchons-nous de la vraie sémantique de C...

Tant qu'à avoir des expressions qui ont des effets de bord, les **affectations** vont être considérées comme des expressions et non des commandes:

la **valeur** de $lval=e$ est celle de e

... et on considère que toute expression peut servir de commande, dont la valeur est ignorée.

Nouvelles abréviations:

$lval+=e \Leftrightarrow lval=lval + e$ (sauf que $lval$ n'est évaluée qu'1 fois)

(pareil pour les autres opérateurs, $-$, $*$, \dots , $>>$, $<<$, \dots)

$++lval \Leftrightarrow lval+=1$

$--lval \Leftrightarrow lval-=1$

$lval++ \Leftrightarrow lval+=1$ (mais retourne l'ancienne val. de $lval$)

$lval-- \Leftrightarrow lval-=1$ (mais retourne l'ancienne val. de $lval$)

La sémantique des commandes de C (1)

$$\frac{\sigma \vdash e \Rightarrow V, \sigma'}{\sigma \vdash ei \Rightarrow \sigma'}$$

$$\frac{}{\sigma \vdash i \Rightarrow \sigma}$$

$$\frac{\sigma \vdash c_1 \Rightarrow \sigma_1 \dots \sigma_{n-1} \vdash c_n \Rightarrow \sigma_n}{\sigma \vdash c_1 \dots c_n \Rightarrow \sigma_n}$$

$$\frac{\sigma \vdash e \Rightarrow 0, \sigma' \quad \sigma' \vdash c_2 \Rightarrow \sigma''}{\sigma \vdash \text{if } (e) \ c_1 \ \text{else } c_2 \Rightarrow \sigma''} \quad \frac{\sigma \vdash e \Rightarrow n, \sigma' \quad \sigma' \vdash c_1 \Rightarrow \sigma''}{\sigma \vdash \text{if } (e) \ c_1 \ \text{else } c_2 \Rightarrow \sigma''}$$

où $n \neq 0$.

Note: $\text{if } (e) \ c_1 \Leftrightarrow \text{if } (e) \ c_1 \ \text{else } i$.

La sémantique des commandes de C (2)

$$\frac{\sigma \vdash e \Rightarrow 0, \sigma'}{\sigma \vdash \text{while } (e) \ c \Rightarrow \sigma'} \quad \frac{\sigma \vdash e \Rightarrow n, \sigma' \quad \sigma' \vdash \text{while } (e) \ c \Rightarrow \sigma''}{\sigma \vdash \text{while } (e) \ c \Rightarrow \sigma''}$$

où $n \neq 0$.

Note: $\text{do } c \ \text{while } (e) \Leftrightarrow c \ \text{while } (e) \ c$

Note: $\text{for } (e_1; e_2; e_3) \ c \Leftrightarrow e_1; \text{while } (e_2) \ \{ c \ e_3; \}$

Exercice: écrire la sémantique du `switch`, en se basant sur celle du `if`.

Retour sur les expressions; évaluation court-circuit

Opérateurs logiques évalués de **gauche à droite**.

$$\frac{\sigma \vdash e_1 \Rightarrow 0, \sigma'}{\sigma \vdash e_1 \&\& e_2 \Rightarrow 0, \sigma'} \quad \frac{\sigma \vdash e_1 \Rightarrow n, \sigma' \quad \sigma' \vdash e_2 \Rightarrow n', \sigma''}{\sigma \vdash e_1 \&\& e_2 \Rightarrow n', \sigma''}$$

Note: e_2 n'est même pas évalué si e_1 vaut 0 (règle de gauche).

→ important notamment lorsque e_2 ne termine pas!

Exercice: écrire la sémantique de $||$, $!$.

Retour sur les expressions: l'indéterminisme

Pour les autres opérations, par ex. +, l'ordre d'évaluation est **non spécifié**.

En première approximation:

$$\frac{\sigma \vdash e_1 \Rightarrow V_1, \sigma' \quad \sigma' \vdash e_2 \Rightarrow V_2, \sigma''}{\sigma \vdash e_1 + e_2 \Rightarrow V_1 + V_2 \bmod (\text{maxint} + 1), \sigma''}$$

$$\frac{\sigma \vdash e_2 \Rightarrow V_2, \sigma' \quad \sigma' \vdash e_1 \Rightarrow V_1, \sigma''}{\sigma \vdash e_1 + e_2 \Rightarrow V_1 + V_2 \bmod (\text{maxint} + 1), \sigma''}$$

Exercice: comment s'évalue `i = 3 + (i=4); i = i++?`

Note: en réalité, c'est plus compliqué: C a le droit d'**entrelacer** les lectures et écritures mémoire de e_1 et e_2 dans un ordre quelconque!

C, informellement

C'est un "assembleur portable".

Assembleur car il manipule des objets au même niveau d'abstraction.

Portable car le même source fonctionne (grosso modo) sur différentes machines.

C \approx assembleur

<code>x = y;</code>	<code>mov y, x</code>
<code>x = y+z;</code>	<code>mov z, x; add y, x</code>
<code>if (x==y) ...</code>	<code>cmp x, y; je ...</code>
<code>goto toto;</code>	<code>jmp toto</code>
<code>x = f ();</code>	<code>call f; mov %eax, x</code>
<code>return 42;</code>	<code>mov \$42, %eax; ret</code>
<code>x = f (y);</code>	<code>push y; call f; add \$4,%esp; mov %eax, x</code>
<code>x = *p;</code>	<code>mov p, %reg; mov (%reg), x</code>
<code>*p = x;</code>	<code>mov p, %reg; mov x, (%reg)</code>

C contient aussi des constructions d'un peu plus haut niveau...

La sémantique formelle (?) de C

- Sémantique définie en langage naturel: norme ISO/IEC 9899-1999 (“ANSI C”).
- Un gros effort:
Nikolaos S. Papaspyrou, *Denotational Semantics of ANSI C*, to appear in *Computer Standards and Interface*, 2001. (Regarder particulièrement la section 3.4.) <http://www.cs.yale.edu/homes/nickie/>.
Version complète: <ftp://ftp.softlab.ntua.gr/pub/users/nickie/papers/thesis.ps.gz>
- Autres références: tapez "langage C" sous Google, vous n'aurez que l'embaras du choix! Une accessible est:
<http://diwww.epfl.ch/w3lsp/teaching/coursC/>