



Analyse statique de logiciel système par typage statique fort

— Application au noyau Linux —

ÉTIENNE MILLON

sous la direction d'Emmanuel Chailloux et de Sarah Zennou

THÈSE

pour obtenir le titre de
Docteur en Sciences
mention Informatique

Soutenue le xx-yy-2013 devant un jury composé de

aaa

aaa

aaa

aaa

aaa

Abstract

ABSTRACT

i

Dédicace

TABLE DES MATIÈRES

Table des matières	iv
1 Introduction	1
I Méthodes formelles pour la sécurité	3
2 Systèmes d'exploitation	5
2.1 Rôle d'un système d'exploitation	5
2.2 Architecture Intel	7
2.2.1 Assembleur	7
2.2.2 Fonctions et conventions d'appel	10
2.2.3 Tâches, niveaux de privilèges	10
2.2.4 Mémoire virtuelle	11
2.3 Cas de Linux	12
2.3.1 Appels système	13
2.4 Sécurité des appels système	15
3 État de l'art	19
II Typage statique de langages impératifs	21
4 Sémantique d'un langage impératif	23
4.1 Principes	23
4.2 Syntaxe	23
4.3 Définitions préliminaires	24
4.4 Mémoire	26
4.5 Accesseurs	29
4.6 Contextes d'évaluation	31
4.7 Expressions	32
4.8 Instructions	39
4.9 Phrases	41
4.10 Exécution	42

TABLE DES MATIÈRES

v

4.11 Exemple : l'algorithme d'Euclide	42
5 Typage	45
5.1 Principe	45
5.2 Définitions	46
5.3 Expressions	48
5.4 Instructions	51
5.5 Fonctions	51
5.6 Phrases	52
5.7 Sûreté du typage	52
6 Qualificateurs de type	59
6.1 Provenance des pointeurs	59
6.1.1 Éditions et ajouts	59
6.1.2 Propriété d'isolation mémoire	60
6.2 Analyse de terminaison des chaînes C	60
6.2.1 But	60
6.2.2 Approche	62
6.2.3 Annotation de <code>string.h</code>	62
6.2.4 Typage des primitives	65
6.2.5 Extensions au système de types	65
6.2.6 Résultats	65
III Expérimentation	67
7 Implantation	69
7.1 Langages intermédiaires	69
7.2 Newspeak	71
7.3 Chaîne de compilation	71
7.3.1 Prétraitement	71
7.3.2 Compilation (levée des ambiguïtés)	71
7.3.3 Annotations	72
7.3.4 Implantation de l'algorithme de typage	73
7.3.5 Algorithme d'unification	81
8 Étude de cas : un pilote de carte graphique	85
8.1 Description du problème	85
8.2 Principes de l'analyse	86
8.3 Implantation	87

8.4 Conclusion	87
9 Conclusion	89
9.1 Limitations	89
9.2 Perspectives	89
A Code du module noyau	91
B Règles d'évaluation	95
C Règles de typage	99
Table des figures	103
Références web	105
Bibliographie	107

CHAPITRE



1

INTRODUCTION

Première partie

Méthodes formelles pour la sécurité

SYSTÈMES D'EXPLOITATION

Le système d'exploitation est le programme qui permet à un système informatique d'exécuter d'autres programmes. Son rôle est donc capital et ses responsabilités multiples. Dans ce chapitre, nous allons voir quel est son rôle, et comment il peut être implanté. Pour ce faire, nous étudierons l'exemple d'une architecture Intel 32 bits, et d'un noyau Linux 2.6.

Pour une description plus détaillée des rôles d'un système d'exploitation ainsi que plusieurs cas d'étude détaillés, on pourra se référer à [Tan07].

2.1 Rôle d'un système d'exploitation

Un ordinateur est constitué de nombreux composants matériels : microprocesseur, mémoire, et divers périphériques. Pourtant, au niveau de l'utilisateur, des dizaines de logiciels permettent d'effectuer toutes sortes de calculs et de communications. Le système d'exploitation permet de faire l'interface entre ces niveaux d'abstraction.

Au cours de l'histoire des systèmes informatiques, la manière de les programmer a beaucoup évolué. Au départ, les programmeurs avaient accès au matériel dans son intégralité : toute la mémoire pouvait être accédée, toutes les instructions pouvaient être utilisées.

Néanmoins c'est un peu restrictif, puisque cela ne permet qu'à une personne d'interagir avec le système. Dans la seconde moitié des années 60, sont apparus les premiers systèmes "à temps partagé", permettant à plusieurs utilisateurs de travailler en même temps.

Permettre l'exécution de plusieurs programmes en même temps est une idée révolutionnaire, mais elle n'est pas sans difficultés techniques : en effet les ressources de la machine doivent être aussi partagées entre les utilisateurs et les programmes. Par exemple, plusieurs programmes vont utiliser le processeur les uns à la suite des

autres (partage *temporel*) ; et chaque programme aura à sa disposition une partie de la mémoire principale, ou du disque dur (partage *spatial*).

Si deux programmes (ou plus) s'exécutent de manière concurrente sur le même matériel, il faut s'assurer que l'un ne puisse pas écrire dans la mémoire de l'autre, ou que les deux utilisent la carte réseau les uns à la suite des autres. Ce sont des rôles du système d'exploitation.

Cela passe donc par une limitation des possibilités du programme : plutôt que de permettre n'importe quel type d'instruction, il communique avec le système d'exploitation. Celui-ci centralise donc les appels au matériel, ce qui permet d'abstraire certaines opérations.

Par exemple, si un programme veut copier des données depuis un cédérom vers la mémoire principale, il devra interroger le bus SATA, interroger le lecteur sur la présence d'un disque dans le lecteur, activer le moteur, calculer le numéro de trame des données sur le disque, demander la lecture, puis déclencher une copie de la mémoire.

Si dans un autre cas il désire récupérer des données depuis une mémoire USB, il devrait interroger le bus USB, rechercher le bon numéro de périphérique, le bon numéro de canal dans celui-ci, lui appliquer une commande de lecture au bon numéro de bloc, puis copier la mémoire.

Ces deux opérations, bien qu'elles aient le même but (copier de la mémoire depuis un périphérique amovible), ne sont pas effectuées en pratique de la même manière. C'est pourquoi le système d'exploitation fournit les notions de fichier, lecteur, etc : le programmeur n'a plus qu'à utiliser des commandes de haut niveau ("monter un lecteur", "ouvrir un fichier", "lire dans un fichier") et selon le type de lecteur, le système d'exploitation effectuera les actions appropriées.

En résumé, un système d'exploitation est l'intermédiaire entre le logiciel et le matériel, et en particulier assure les rôles suivants :

- Gestion des processus : un système d'exploitation peut permettre d'exécuter plusieurs programmes à la fois. Il faut alors orchestrer ces différents processus et les séparer en terme de temps et de ressources partagées.
- Gestion de la mémoire : chaque processus, en plus du noyau, doit disposer d'un espace mémoire différent. C'est-à-dire qu'un processus ne doit pas pouvoir interférer avec un autre.
- Gestion des fichiers : les processus peuvent accéder à une hiérarchie de fichiers, indépendamment de la manière d'y accéder.
- Gestion des périphériques : le noyau étant le seul code ayant des privilèges, c'est lui qui doit communiquer avec les périphériques matériels.
- Abstractions : le noyau fournit aux programmes une interface unifiée, permettant de stocker des informations de la même manière sur un disque dur ou une clef USB (alors que l'accès se déroulera de manière très différente en pratique).

2.2 Architecture Intel

L'implantation d'un système d'exploitation est très proche du matériel sur lequel il s'exécute. Pour étudier une implantation en particulier, voyons ce que permet le matériel lui-même.

Dans cette section nous décrivons le fonctionnement d'un processeur utilisant une architecture Intel 32 bits. Les exemples de code seront écrits en syntaxe AT&T, celle que comprend l'assembleur GNU.

La référence pour la description de l'assembleur Intel est la documentation du constructeur [Int]; une bonne explication de l'agencement dans la pile peut aussi être trouvée dans [One96].

2.2.1 Assembleur

Pour faire des calculs, le processeur est composé de registres, qui sont des petites zones de mémoire interne, et peut accéder à la mémoire principale.

La mémoire principale contient divers types des données :

- le code des programmes à exécuter
- les données à disposition des programmes
- la pile d'appels

La pile d'appels est une zone de mémoire qui est notamment utilisée pour tenir une trace des calculs en cours. Par exemple, c'est ici que seront stockées les données propres à chaque fonction appelée : paramètres, adresse de retour et variables locales. La pile est manipulée par un pointeur de pile (*stack pointer*), qui est l'adresse du "haut de la pile". On peut la manipuler en empilant des données (les placer au niveau du pointeur de pile et déplacer celui si) ou dépilant des données (déplacer le pointeur de pile dans l'autre sens et retourner la valeur présente à cet endroit).

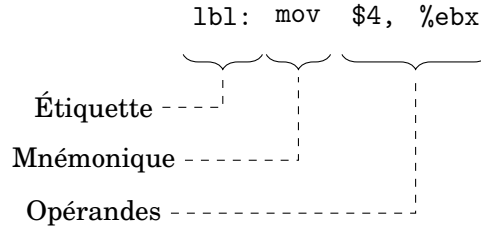
L'état du processeur est défini par la valeur de ses registres, qui sont des petites zones de mémoire interne (quelques bits chacun). Par exemple, la valeur du pointeur de pile est stockée dans ESP. Le registre EBP, couplé à ESP sert à adresser les variables locales et paramètres d'une fonction, comme ce sera expliqué dans la section 2.2.2.

L'adresse de l'instruction courante est accessible dans le registre EIP.

En plus de ces registres spéciaux, le processeur possède de nombreux registres génériques, qui peuvent être utilisés pour réaliser des calculs intermédiaires. Ils sont nommés EAX, EBX, ECX, EDX, ESI et EDI. Ils peuvent être utilisés pour n'importe quel type d'opération, mais certains sont spécialisés : par exemple il est plus efficace d'utiliser EAX en accumulateur, ou ECX en compteur.

Les calculs sont décrits sous forme d'une suite d'instructions. Chaque instruction est composée d'un mnémonique et d'une liste d'opérandes. Les mnémoniques (*mov*,

call, sub, etc) définissent un type d'opération à appliquer sur les opérandes. L'instruction peut aussi être précédée d'une étiquette, qui correspondra à son adresse.



Ces opérandes peuvent être de plusieurs types :

- un nombre, noté \$4
- le nom d'un registre, noté %eax
- une opérande mémoire, c'est à dire le contenu de la mémoire à une adresse effective. Cette adresse effective peut être exprimée de plusieurs manières :
 - directement : addr
 - indirectement : (%ecx). L'adresse effective est le contenu du registre.
 - "base + déplacement" : 4(%ecx). L'adresse effective est le contenu du registre plus le déplacement (4 ici).

En pratique il y a des modes d'adressage plus complexes, et toutes les combinaisons ne sont pas possibles, mais ceux-ci suffiront à décrire les exemples suivants :

- `mov src, dst` copie le contenu de `src` dans `dst`.
- `add src, dst` calcule la somme des contenus de `src` et `dst` et place ce résultat dans `dst`.
- `push src` place `src` sur la pile, c'est à dire que cette instruction enlève au pointeur de pile ESP la taille de `src`, puis place `src` à l'adresse mémoire de la nouvelle valeur ESP.
- `pop src` réalise l'opération inverse : elle charge le contenu de la mémoire à l'adresse ESP dans `src` puis incrémente ESP de la taille correspondante.
- `jmp addr` saute à l'adresse `addr` : c'est l'équivalent de `mov addr, %eip`.
- `call addr` sert aux appels de fonction : cela revient à `push %eip` puis `jmp addr`.
- `ret` sert à revenir d'une fonction : c'est l'équivalent de `pop %eip`.

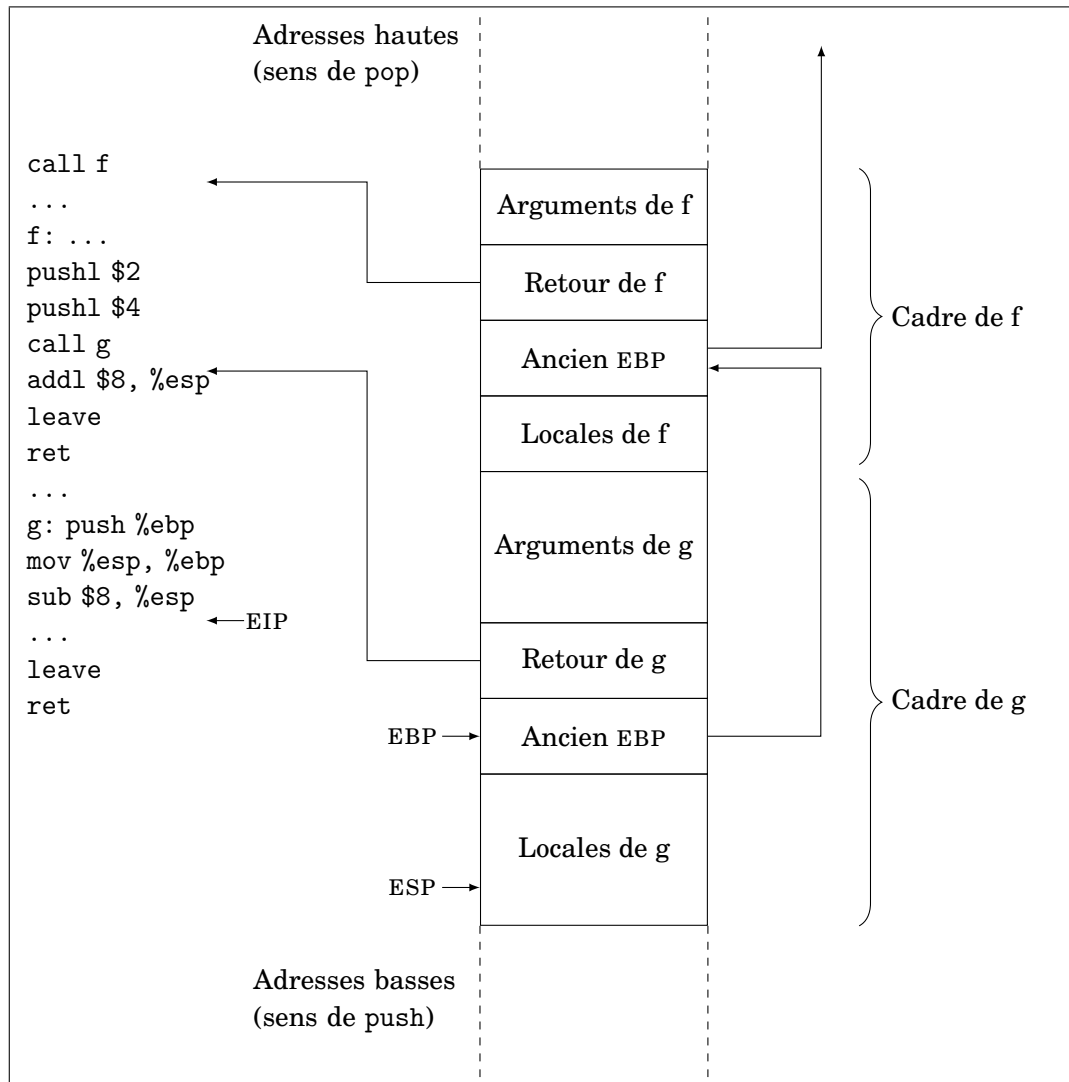


FIGURE 2.1 – Cadres de pile.

2.2.2 Fonctions et conventions d'appel

Dans le langage d'assemblage, il n'y a pas de notion de fonction ; mais `call` et `ret` permettent de sauvegarder et de restaurer une adresse de retour, ce qui permet de faire un saut et revenir à l'adresse initiale. Ce système permet déjà de créer des procédures, c'est-à-dire des fonctions sans arguments ni valeur de retour.

Pour gérer ceux-ci, il faut que les deux morceaux (appelant et appelé) se mettent d'accord sur une convention d'appel commune. La convention utilisée sous GNU/Linux est appelée *cdecl* et possède les caractéristiques suivantes :

- la valeur de retour d'une fonction est stockée dans EAX
- EAX, ECX et EDI peuvent être écrasés sans avoir à les sauvegarder
- les arguments sont placés sur la pile (et enlevés) par l'appelant. Les paramètres sont empilés de droite à gauche.

Pour accéder à ses paramètres, une fonction peut donc utiliser un adressage relatif à ESP. Cela peut fonctionner, mais cela complique les choses si elle contient aussi des variables locales. En effet, les variables locales sont placées sur la pile, au dessus des (c'est à dire, empilées après) paramètres, augmentant le décalage.

Pour simplifier, la pile est organisée en cadres logiques : chaque cadre correspond à un niveau dans la pile d'appels de fonctions. Si `f` appelle `g`, qui appelle `h`, il y aura dans l'ordre sur la pile le cadre de `f`, celui de `g` puis celui de `h`.

Ces cadres sont chaînés à l'aide du registre EBP : à tout moment, EBP contient l'adresse du cadre de l'appelant.

Prenons exemple sur la figure 2.1 : pour appeler `g(4,2)`, `f` empile les arguments de droite à gauche. L'instruction `call g` empile l'adresse de l'instruction suivante sur la pile puis saute au début de `g`.

Au début de la fonction, les trois instructions permettent de sauvegarder l'ancienne valeur de EBP, faire pointer EBP à une endroit fixe dans le cadre de pile, puis allouer 8 octets de mémoire pour les variables locales.

Dans le corps de la fonction `g`, on peut donc se référer aux variables locales par `-4(%ebp)`, `-8(%ebp)`, etc, et aux arguments par `8(%ebp)`, `12(%ebp)`, etc.

À la fin de la fonction, l'instruction `leave` est équivalente à `mov %ebp, %esp` puis `pop %ebp` et permet de défaire le cadre de pile, laissant l'adresse de retour en haut de pile. Le `ret` final la dépile et y saute.

2.2.3 Tâches, niveaux de privilèges

Sans mécanisme particulier, le processeur exécuterait uniquement une suite d'instruction à la fois. Pour lui permettre d'exécuter plusieurs tâches, un système de partage du temps existe.



FIGURE 2.2 – Les différents *rings*. Seul le *ring 0* a accès au hardware via des instructions privilégiées. Pour accéder aux fonctionnalités du noyau, les programmes utilisateur doivent passer par des appels système.

À des intervalles de temps réguliers, le système est programmé pour recevoir une interruption. C'est une condition exceptionnelle (au même titre qu'une division par zéro) qui fait sauter automatiquement le processeur dans une routine de traitement d'interruption. À cet endroit le code peut sauvegarder les registres et restaurer un autre ensemble de registres, ce qui permet d'exécuter plusieurs tâches de manière entrelacée. Si l'alternance est assez rapide, cela peut donner l'illusion que les programmes s'exécutent en parallèle. Comme l'interruption peut survenir à tout moment, on parle de multitâche préemptif.

En plus de cet ordonnancement de processus, l'architecture Intel permet d'affecter des niveaux de privilège à ces tâches, en restreignant le type d'instructions exécutables, ou en donnant un accès limité à la mémoire aux tâches de niveaux moins élevés.

Il y a 4 niveaux de privilèges (nommés aussi *rings* - figure 2.2) : le *ring 0* est le plus privilégié, le *ring 3* le moins privilégié. Dans l'exemple précédent, on pourrait isoler l'ordonnanceur de processus en le faisant s'exécuter en *ring 0* alors que les autres tâches seraient en *ring 3*.

2.2.4 Mémoire virtuelle

À partir du moment où plusieurs processus s'exécutent de manière concurrente, un problème d'isolation se pose : si un processus peut lire dans la mémoire d'un autre, des informations peuvent fuir ; et s'il peut y écrire, il peut en détourner l'exécution.

Le mécanisme de mémoire virtuelle permet de donner à deux tâches une vue différente de la mémoire : c'est à dire que vue de tâches différentes, une adresse contiendra une valeur différente.

Ce mécanisme est contrôlé par valeur du registre CR3 : les 10 premiers bits d'une adresse virtuelle sont un index dans le répertoire de pages qui commence à l'adresse contenue dans CR3. À cet index, se trouve l'adresse d'une table de pages. Les 10 bits suivants de l'adresse sont un index dans cette page, donnant l'adresse d'une page de 4 kibi-octets (figure 2.3).

En ce qui concerne la mémoire, les différentes tâches ont une vision différente de la mémoire physique : c'est-à-dire que deux tâches lisant à une même adresse peuvent

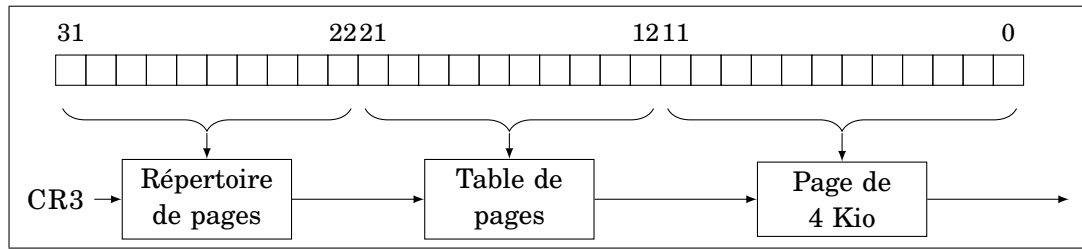


FIGURE 2.3 – Implantation de la mémoire virtuelle

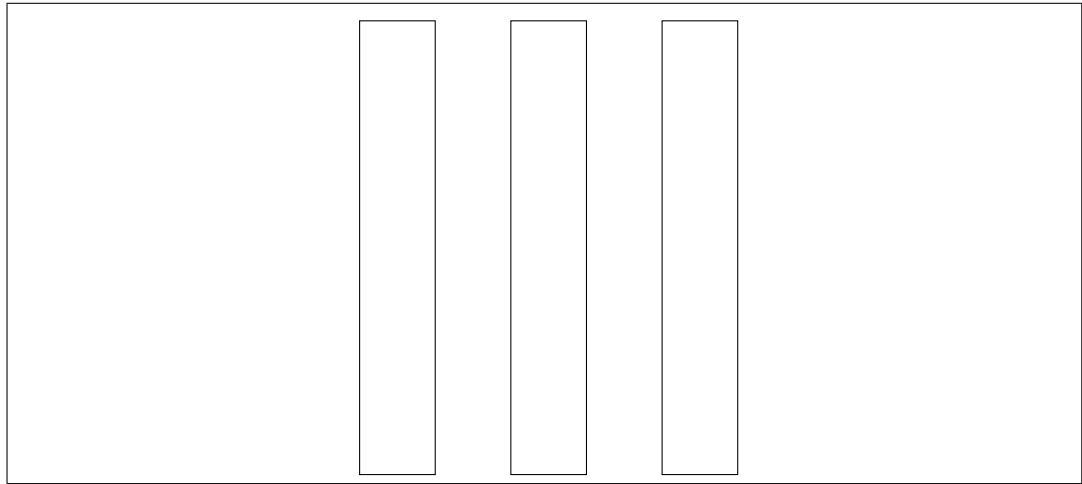


FIGURE 2.4 – Mécanisme de mémoire virtuelle.

avoir un résultat différent. C'est le concept de mémoire virtuelle (fig 2.4).

2.3 Cas de Linux

Dans cette section, nous allons voir comment ces mécanismes sont implantés dans le noyau Linux. Une description plus détaillée pourra être trouvée dans [BC05], ou pour le cas de la mémoire virtuelle, [Gor04].

Deux rings sont utilisés : en *ring* 0, le code noyau et en *ring* 3, le code utilisateur.

Une notion de tâche similaire à celle décrite dans la section 2.2.3 existe : elles s'exécutent l'une après l'autre, le changement s'effectuant sur interruptions.

Pour faire appel aux services du noyau, le code utilisateur doit faire appel à des appels systèmes, qui sont des fonctions exécutées par le noyau. Chaque tâche doit donc avoir deux piles : une pile "utilisateur" qui sert pour l'application elle-même, et une pile "noyau" qui sert aux appels système.

Grâce à la mémoire virtuelle, chaque processus possède sa propre vue de la mémoire dans son espace d'adressage (figure 2.5), et donc chacun un ensemble de tables

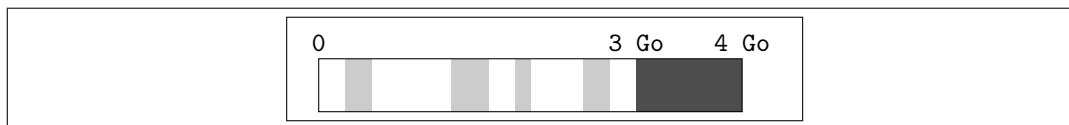


FIGURE 2.5 – L’espace d’adressage d’un processus. En gris clair, les zones accessibles à tous les niveaux de privilèges : code du programme, bibliothèques, tas, pile. En gris foncé, la mémoire du noyau, réservée au mode privilégié.

de pages et une valeur de CR3 associée. Au moment de changer le processus en cours, l’ordonnanceur charge donc le CR3 du nouveau processus.

Les adresses basses (inférieures à `PAGE_OFFSET = 3 Gio = 0xc0000000`) sont réservées à l’utilisateur. On y trouvera par exemple :

- le code du programme
- les données du programmes (variables globales)
- la pile utilisateur
- le tas (mémoire allouée par `malloc` et fonctions similaires)
- les bibliothèques partagées

Au dessus de `PAGE_OFFSET`, se trouve la mémoire réservée au noyau. Cette zone contient le code du noyau, les piles noyau des processus, etc.

2.3.1 Appels système

Les programmes utilisateur s’exécutant en *ring* 3, ils ne peuvent pas contenir d’instructions privilégiées, et donc ne peuvent pas accéder directement au matériel (c’était le but !). Pour qu’ils puissent interagir avec le système (afficher une sortie, écrire sur le disque...), le mécanisme des appels système est nécessaire. Il s’agit d’une interface de haut niveau entre les *rings* 3 et 0. Du point de vue du programmeur, il s’agit d’un ensemble de fonctions C “magiques” qui font appel au système d’exploitation pour effectuer des opérations.

Voyons ce qui se passe derrière la magie apparente. Une explication plus détaillée est disponible dans la documentation fournie par Intel [Int].

Dans la bibliothèque C

Il y a bien une fonction `getpid` présente dans la bibliothèque C du système. C’est la fonction qui est directement appelée par le programme. Cette fonction commence par placer le numéro de l’appel système (noté `__NR_getpid`, valant 20 ici) dans `EAX`, puis les arguments éventuels dans les registres (`EBX`, `ECX`, `EDX`, `ESI` puis `EDI`). Une interruption logicielle est ensuite déclenchée (`int 0x80`).

Dans la routine de traitement d'interruption

Étant donné la configuration du processeur¹, elle sera traitée en *ring* 0, à un point d'entrée prédéfini (`arch/x86/kernel/entry_32.S`, `ENTRY(system_call)`).

```

    # system call handler stub
ENTRY(system_call)
    RINGO_INT_FRAME                # can't unwind into user space anyw
    pushl %eax                    # save orig_eax
    CFI_ADJUST_CFA_OFFSET 4
    SAVE_ALL
    GET_THREAD_INFO(%ebp)

                                # system call tracing in operation / emula
    testl $_TIF_WORK_SYSCALL_ENTRY, TI_flags(%ebp)
    jnz syscall_trace_entry
    cmpl $(nr_syscalls), %eax
    jae syscall_badsys
syscall_call:
    call *sys_call_table(,%eax,4)
    movl %eax,PT_EAX(%esp)        # store the return value
    # ...
    INTERRUPT_RETURN

```

L'exécution reprend donc en *ring* 0, avec dans ESP le pointeur de pile noyau du processus. Les valeurs des registres ont été préservées, la macro `SAVE_ALL` les place sur la pile. Ensuite, à l'étiquette `syscall_call`, le numéro d'appel système (toujours dans EAX) sert d'index dans le tableau de fonctions `sys_call_table`.

Dans l'implantation de l'appel système

Puisque les arguments sont en place sur la pile, comme dans le cas d'un appel de fonction “classique”, la convention d'appel *cdecl* est respectée. La fonction implantant l'appel système, nommée `sys_getpid`, peut donc être écrite en C.

On trouve cette fonction dans `kernel/timer.c` :

```

SYSCALL_DEFINE0(getpid)
{
    return task_tgid_vnr(current);
}

```

1. Il est impropre de dire que le processeur est configuré — tout dépend uniquement de l'état de certains registres, ici la *Global Descriptor Table* et les *Interrupt Descriptor Tables*.

La macro `SYSCALL_DEFINE0` nomme la fonction `sys_getpid`, et définit entre autres des points d'entrée pour les fonctionnalités de débogage du noyau. À la fin de la fonction, la valeur de retour est placée dans `EAX`, conformément à la convention *cdecl*.

Retour vers le ring 3

Au retour de la fonction, la valeur de retour est placée à la place de `EAX` là où les registres ont été sauvegardés sur la pile noyau (`PT_EFLAGS(%esp)`). L'instruction `iret` (derrière la macro `INTERRUPT_RETURN`) permet de restaurer les registres et de repasser en mode utilisateur, juste après l'interruption. La fonction de la bibliothèque C peut alors retourner au programme appelant.

2.4 Sécurité des appels système

On a vu que les appels systèmes permettent aux programmes utilisateur d'accéder aux services du noyau. Ils forment donc une interface particulièrement sensible aux problèmes de sécurité.

Comme pour toutes les interfaces, on peut être plus ou moins fin. D'un côté, une interface pas assez fine serait trop restrictive et ne permettrait pas d'implémenter tout type de logiciel. De l'autre, une interface trop laxiste ("écrire dans tel registre matériel") empêche toute isolation. Il faut donc trouver la bonne granularité.

Nous allons présenter ici une difficulté liée à la manipulation de mémoire au sein de certains types d'appels système.

Il y a deux grands types d'appels systèmes : d'une part, ceux qui renvoient un simple nombre, comme `getpid` qui renvoie le numéro du processus appelant.

```
pid_t pid = getpid();  
printf("%d\n", pid);
```

Ici, pas de difficulté particulière : la communication entre le *ring* 0 et le *ring* 3 est faite uniquement à travers les registres, comme décrit dans la section 2.3.1.

Mais la plupart des appels systèmes communiquent de l'information de manière indirecte, à travers un pointeur. L'appellant alloue une zone mémoire dans son espace d'adressage et passe un pointeur à l'appel système. Ce mécanisme est utilisé par exemple par la fonction `gettimeofday` (figure 2.6).

Considérons une implémentation naïve de cet appel système qui écrirait directement à l'adresse pointée. La figure 2.7a présente ce qui se passe lorsque le pointeur fourni est dans l'espace d'adressage du processus : c'est le cas d'utilisation normal et l'écriture est donc possible.

Si l'utilisateur passe un pointeur dont la valeur est supérieure à `0xc0000000` (figure 2.7b), que se passe-t-il ? Comme le déréférencement est fait dans le code du

```

struct timeval tv;
struct timezone tz;
int z = gettimeofday(&tv, &tz);
if (z == 0) {
    printf( "tv.tv_sec = %ld\ntv.tv_usec = %ld\n"
           "tz.tz_minuteswest = %d\ntz.tz_dsttime = %d\n",
           tv.tv_sec, tv.tv_usec,
           tz.tz_minuteswest, tz.tz_dsttime
    );
}

```

FIGURE 2.6 – Appel de gettimeofday

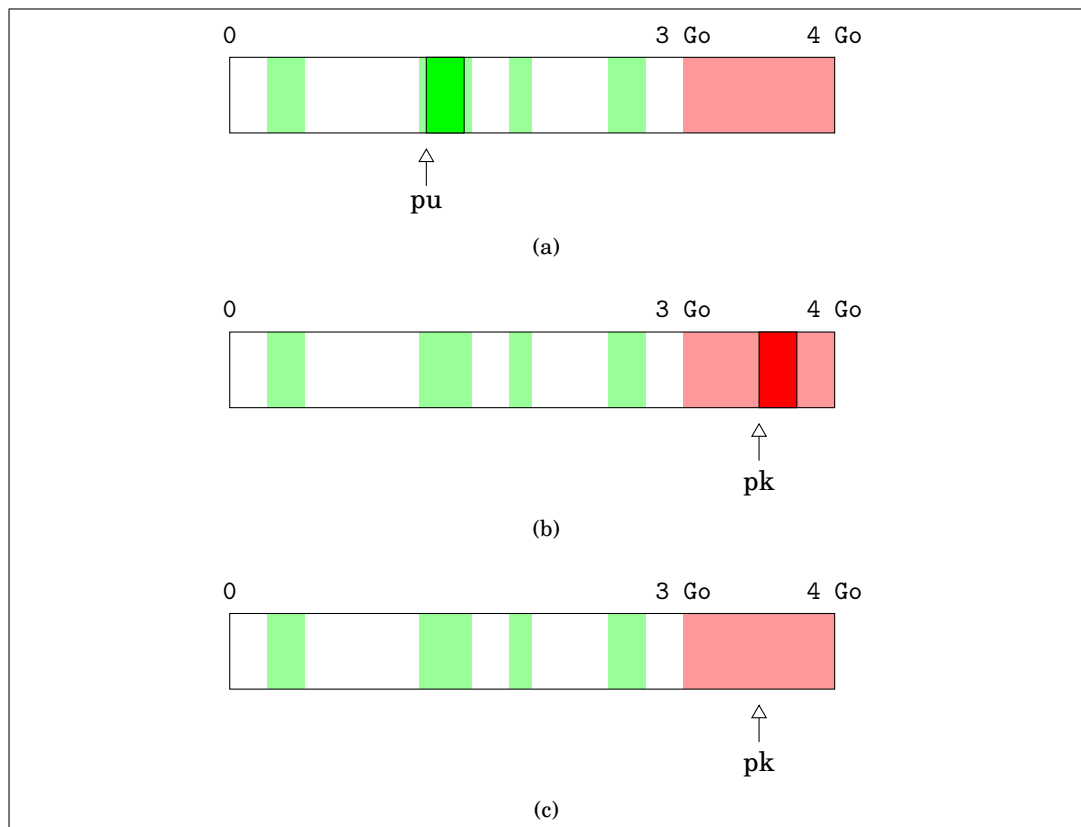


FIGURE 2.7 – Zones mémoire


```
SYSCALL_DEFINE2(gettimeofday, struct timeval __user *, tv,
                struct timezone __user *, tz)
{
    if (likely(tv != NULL)) {
        struct timeval ktv;
        do_gettimeofday(&ktv);
        if (copy_to_user(tv, &ktv, sizeof(ktv)))
            return -EFAULT;
    }
    if (unlikely(tz != NULL)) {
        if (copy_to_user(tz, &sys_tz, sizeof(sys_tz)))
            return -EFAULT;
    }
    return 0;
}
```

FIGURE 2.8 – Implantation de l'appel système gettimeofday

noyau, il est également fait en *ring 0*, et va pouvoir être réalisé sans erreur : l'écriture se fait et potentiellement une structure importante du noyau est écrasée.

Un utilisateur malicieux peut donc utiliser cet appel système pour écrire à n'importe quelle adresse dans l'espace d'adressage du noyau. Ce problème vient du fait que l'appel système utilise les privilèges du noyau au lieu de celui qui contrôle la valeur des paramètres sensibles. Cela s'appelle le *Confused Deputy Problem*[Har88].

La bonne solution est de tester dynamiquement la valeur du pointeur : si la valeur du pointeur est supérieure à `0xc0000000`, il faut indiquer une erreur avant d'écrire (figure 2.7c. Sinon, cela ne veut pas dire que le déréférencement se fera sans erreur, mais au moins le noyau est protégé.

Dans le noyau, un ensemble de fonctions permet d'effectuer des copies sûres. La fonction `access_ok` réalise le test décrit précédemment. Les fonctions `copy_from_user` et `copy_to_user` réalisent une copie de la mémoire après avoir fait ce test. Ainsi, l'implantation correcte de l'appel système `gettimeofday` fait appel à celle-ci (figure 2.8).

Pour préserver la sécurité du noyau, il est donc nécessaire de vérifier la valeur de tous les pointeurs dont la valeur est contrôlée par l'utilisateur. Cette conclusion est assez contraignante, puisqu'il existe de nombreux endroits dans le noyau où des données proviennent de l'utilisateur. Il est donc raisonnable de vouloir vérifier automatiquement et statiquement l'absence de tels défauts.

ÉTAT DE L'ART

L'analyse statique de programmes est un sujet de recherche actif depuis l'apparition de la science informatique.

L'analyse la plus simple consiste à traiter un programme comme du texte, et à y rechercher des motifs dangereux. Ainsi, utiliser des outils comme `grep` permet parfois de trouver un grand nombre de vulnérabilités[Spe05].

On peut continuer cette approche en recherchant des motifs mais en étant sensible à la syntaxe et au flot de contrôle du programme. Cette notion de *semantic grep* est présente dans l'outil Coccinelle [BDH⁺09, PTS⁺11] : on peut définir des *patches sémantiques* pour détecter ou modifier des constructions particulières.

Dans le cas particulier des vulnérabilités liées à une mauvaise utilisation de la mémoire, les développeurs du noyau Linux ont ajouté un système d'annotations au code source. Un pointeur peut être décoré d'une annotation `__kernel` ou `__user` selon s'il est sûr ou pas. Celle-ci sont ignorées par le compilateur, mais un outil d'analyse statique ad-hoc nommé Sparse [S⁺05] peut être utilisé pour détecter les cas les plus simples d'erreurs.

Ce système d'annotations sur les types a été formalisé sous le nom de *qualificateurs de types* : chaque type peut être décoré d'un ensemble de qualificateurs (à la manière de `const`), et des règles de typage permettent d'établir des propriétés sur le programme. Ces analyses ont été implantées dans l'outil CQual [FFA99, STFW01, FTA02, JW04, FJKA06]. .

L'interprétation abstraite est une technique d'analyse générique qui permet de simuler statiquement tous les comportements d'un programme Cousot [CC77, CC92]. Un exemple d'application est de calculer les bornes de variations des variables pour s'assurer qu'aucun débordement de tableau n'est possible. Cette technique est très puissante mais possède plusieurs inconvénients. D'une part, pour réaliser une analyse interprocédurale il faut partir d'un point en particulier du programme (comme

la fonction `main`). Cette hypothèse n'est pas facilement satisfaite dans un noyau de système d'exploitation, qui possède de nombreux points d'entrée. D'autre part, il est très difficile de faire passer à l'échelle un interpréteur abstrait [CCF⁺09, BBC⁺10].

L'approche par typage, plus légère, est séduisante. Pour les différents enjeux des systèmes de types statiques, on pourra se référer à [Pie02]. Il est possible d'encoder ce genre de propriétés dans un système de types, cf. [KcS07] et [LZ06].

On peut aller plus loin que les simples types et utiliser un langage de contrats : chaque fonction est annotée par des pré- et post-conditions sur la mémoire [DRS00].

Du côté de l'analyse dynamique, [SAB10].

Ce que nous voulons vérifier peut être vue comme une propriété de flot. Un *survey* des problèmes et techniques existantes peut être trouvé dans [SM03].

(wip)

Interprétation Abstraite : widening [Gra92], CGS [VB04], Astrée : presentation [Mau04, CCF⁺05],

Divers : Taint sequences [CMP10],

Frama-C ?

CCured ?

Deuxième partie

Typage statique de langages impératifs

SÉMANTIQUE D'UN LANGAGE IMPÉRATIF

On définit ici un langage impératif servant de base à nos analyses. Un premier système de types sera donné dans le chapitre 5.

4.1 Principes

Nous voulons capturer l'essence de C. Les traits principaux sont les suivants :

- les types de données sont très simples : entiers machine, flottants, pointeurs et types composés (structures et tableaux) composés de ceux-ci.
- les variables sont mutables, et on peut passer des données par valeur ou par pointeur.
- le flôt de contrôle repose sur les constructions "if" et "while".
- le code est organisé en fonctions "simples", c'est-à-dire qui ne sont pas des fermatures.

4.2 Syntaxe

Les figures 4.2, 4.1 et 4.3 présentent notre langage intermédiaire. Il contient la plupart des fonctionnalités présentes dans les langages impératifs comme C.

Un programme est organisé en fonctions, qui contiennent des instructions, qui elles-même manipulent des expressions.

Le flot de contrôle est simplifié par rapport à C : il ne contient que l'alternative ("if") et la boucle "while". Les autres formes de boucle ("do/while" et "for") peuvent être émulées par une boucle "while".

Les fonctionnalités manquantes, et comment les émuler, seront discutés dans le chapitre??.

Constantes	$c ::= i$	Entier
	f	Flottant
	NULL	Pointeur nul
Expressions	$e ::= c$	Constante
	lv	Accès mémoire
	$\Box e$	Opération unaire
	$e \boxplus e$	Opération binaire
	$\&lv$	Pointeur
	$lv \leftarrow e$	Affectation
	$\{f_1 : e_1; \dots; f_n : e_n\}$	Structure
	$\{e_1; \dots; e_n\}$	Tableau
	f	Fonction
	$e(e_1, \dots, e_n)$	Appel de fonction
Left-values	$lv ::= x$	Variable
	$*lv$	Déréférencement
	$lv.f$	Accès à un champ
	$lv[e]$	Accès à un élément

FIGURE 4.1 – Syntaxe - expressions

Pour l'alternative, on introduit également la forme courte $\text{IF}(e)\{i\} = \text{IF}(e)\{i\}\text{ELSE}\{\text{PASS}\}$.
Les opérateurs sont donnés dans la figure 4.3.

4.3 Définitions préliminaires

On suppose avoir à notre disposition un ensemble infini dénombrable d'identificateurs ID (par exemple des chaînes de caractères).

X^* est l'ensemble des suites finies de X , indexées à partir de 1. Si $u \in X^*$, on note $|u|$ le nombre d'éléments de u (le cardinal de son ensemble de définition). Pour $1 \leq i \leq |u|$, on note $u_i = u(i)$ le i -ème élément de la suite.

On peut aussi voir les suites comme des listes : on note $[]$ la suite vide, telle que $|[]| = 0$. On définit en outre la construction de suite de la manière suivante : si $x \in X$ et $u \in X^*$, la liste $x :: u \in X^*$ est la liste v telle que :

Instructions	$i ::= \text{PASS}$	Instruction vide
	$i; i$	Séquence
	e	Expression
	$\text{IF}(e)\{i\}\text{ELSE}\{i\}$	Alternative
	$\text{WHILE}(e)\{i\}$	Boucle
	$\text{RETURN}(e)$	Retour de fonction
Fonctions	$f ::= \text{fun}(x_1, \dots, x_n)$	Arguments
	$((x'_1, e_1), \dots, (x'_p, e_p))$	Variables locales
	$\{i\}$	Corps
Phrases	$p ::= x = e$	Variable globale
	e	Évaluation d'expression
	$\text{struct } s\{x_1 : t_1; \dots; x_n : t_n\}$	Déclaration de structure
Programme	$P ::= (p_1, \dots, p_n)$	Phrases

FIGURE 4.2 – Syntaxe - instructions

Opérateurs binaires	$\boxplus ::= +, -, \times, /$	Arithmétique entière
	$+., -., \times., /.$	Arithmétique flottante
	$+_p, -_p$	Arithmétique de pointeurs
	$=, \neq, \leq, \geq, <, >$	Comparaisons
	$\&, , ^$	Opérateurs bit à bit
	$\&\&, $	Opérateurs logiques
	\ll, \gg	Décalages
Opérateurs unaires	$\boxminus ::= +, -$	Arithmétique entière
	$+., -.$	Arithmétique flottante
	\sim	Négation bit à bit
	$!$	Négation logique

FIGURE 4.3 – Syntaxe - opérateurs

$$v_1 = x$$

$$\forall i \in [1; |u|], v_{i+1} = u_i$$

La concaténation des listes u et v est la liste $u@v = w$ telle que :

$$|w| = |u| + |v|$$

$$\forall i \in [1; |u|], w_i = u_i$$

$$\forall j \in [1; |v|], w_{|u|+j} = v_j$$

4.4 Mémoire

L'interprète que nous nous apprêtons à définir manipule des valeurs qui sont associées aux variables du programme.

$$\begin{aligned} \text{VAL} &= \text{INT} \uplus \text{FLOAT} \uplus \{\text{NULL}\} \uplus \Phi \\ \text{INT} &= \mathbb{Z}/2^{32}\mathbb{Z} - 2^{31} \\ \text{FLOAT} &= \text{IEEEFLOAT}(32) \\ \Phi &= \text{fix}(\Phi') \\ \Phi'(X) &= X \\ &\quad \cup \text{ADDR} \\ &\quad \cup \{*\varphi/\varphi \in \Phi\} \\ &\quad \cup \{\varphi.f/\varphi \in \Phi, f \in \text{ID}\} \\ &\quad \cup \{\varphi[n]/\varphi \in \Phi, n \in \text{INT}\} \\ \text{ADDR} &= \text{ID} \uplus (\mathbb{N} \times \text{ID}) \end{aligned}$$

$\text{IEEEFLOAT}(n)$ correspond à l'ensemble des flottants IEEE 754 de n bits[oEE08]. Ici, INT est choisi pour représenter les nombres entiers de -2^{31} à $2^{31} - 1$, mais ce choix est arbitraire : de la même manière, on aurait pu choisir des nombres à 64 bits ou même de précision arbitraire.

L'ensemble des états mémoire est :

$$\text{MEM} = ((\text{ID} \times \text{VAL})^*)^* \times (\text{ID} \times \text{VAL})^*$$

Un état mémoire état mémoire $(s, g) \in \text{MEM}$ est composé :

Valeurs	$v ::= c$	Constante
	φ	Référence mémoire
	$\{f_1 : v_1; \dots; f_n : v_n\}$	Structure
	$\{v_1; \dots; v_n\}$	Tableau
	f	Fonction
Adresses	$a ::= (n, x)$	Variable locale
	x	Variable globale
Chemins	$\varphi ::= a$	Adresse
	$*\varphi$	Déréférencement
	$\varphi.f$	Accès à un champ
	$\varphi[n]$	Accès à un élément
Pile	$s ::= []$	Pile vide
	$\{x_1; \dots; x_n\} :: s$	Ajout d'un cadre
État mémoire	$m ::= (s,$	Pile
	$\{x_1; \dots; x_n\},$	Globales
	$\{a_1 \mapsto v_1; \dots; a_p \mapsto v_p\})$	Valeurs

FIGURE 4.4 – Composantes d'un état mémoire

- d'une part, d'une pile s de cadres, qui sont des listes d'association (nom de variable, valeur).
- d'autre part, une liste d'association qui représente les variables globales.

La structure de pile des locales permet de les organiser en niveaux indépendants : à chaque appel de fonction, un nouveau cadre de pile est créé, comprenant ses paramètres et ses variables locales.

Au contraire, pour les globales il n'y a pas de système d'empilement, puisque ces variables sont accessibles depuis tout point du programme.

Ces définitions sont résumées dans la figure 4.4

Définition 4.1 (Recherche de variable). *La recherche de variable permet d'associer à une variable x une adresse a .*

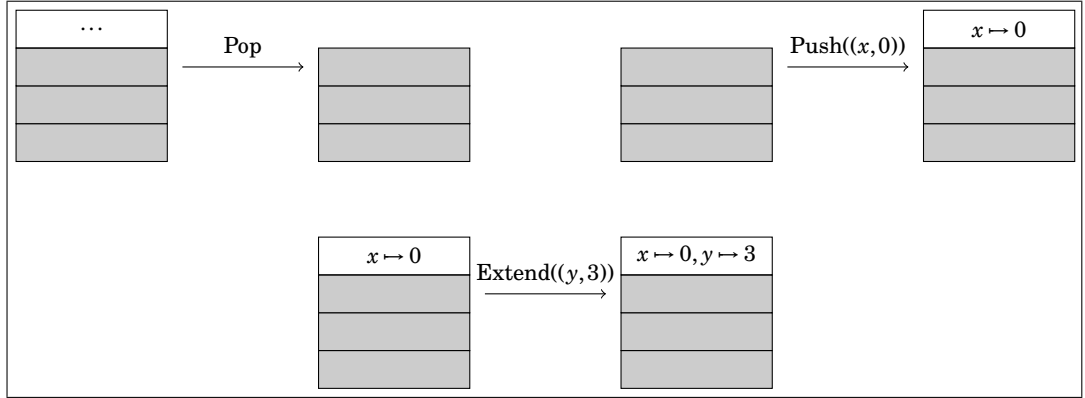


FIGURE 4.5 – Opérations de pile

Chaque fonction peut accéder aux variables locales de la fonction en cours, ainsi qu'aux variables globales.

$$\text{Lookup}((s, g), x) = (|s|, x) \text{ si } |s| > 0 \text{ et } \exists (x, v) \in s_1$$

$$\text{Lookup}((s, g), x) = x \text{ si } (x, v) \in g$$

En entrant dans une fonction, on rajoutera un cadre de pile qui contient les paramètres de la fonction ainsi que ses variables locales. En retournant à l'appelant, il faudra supprimer ce cadre de pile.

Définition 4.2 (Manipulations de pile). *On définit l'empilement d'un cadre de pile $f = ((x_1, v_1), \dots, (x_n, v_n))$ sur un état mémoire $m = (s, g)$:*

$$\text{Push}((s, g), f) = (f :: s, g)$$

On définit aussi l'extension du dernier cadre de pile :

$$\text{Extend}((f :: s, g), x) = (((x @ f) :: s), g)$$

De même on définit le dépilement :

$$\text{Pop}(f :: s, g) = (s, g)$$

Ces définitions sont illustrées dans la figure 4.5.

4.5 Accesseurs

On définit quelques accesseurs. Un accesseur $[\cdot]$ permet d'accéder à une structure s et d'obtenir un élément e à partir d'un indice i (noté $e = s[i]$) ou de modifier le sous-élément à l'indice i par e (noté $s[i \leftarrow e]$).

Définition 4.3 (Accès à une liste d'associations). *Une liste d'association est une liste de paires (clef, valeur) avec l'invariant supplémentaire que les clefs sont uniques. Il est donc possible de trouver au plus une valeur associée à une clef donnée. L'écriture est également possible, en remplaçant un couple par un couple avec une valeur différente.*

L'accesseur $[\cdot]_L$ est défini par :

$$l[x]_L = v \text{ où } \{v\} = \{y / (x, y) \in l\}$$

$$l[x \leftarrow v]_L = (x, v) :: \{(y, v) \in g(x) / y \neq x\}$$

Définition 4.4 (Accès par adresse). *Les états mémoire sont constitués des listes d'association (nom, valeur).*

L'accesseur par adresse $[\cdot]_A$ permet de généraliser l'accès à ces valeurs en utilisant comme clef non pas un nom mais une adresse.

Selon cette adresse, on accède soit à la liste des globales, soit à une des listes de la pile des locales.

Pour $m = (s, g)$,

$m[x]_A = g[x]_L$	<i>Lecture d'une globale</i>
$m[(n, x)]_A = s_{ l -n+1}[x]_L$	<i>Lecture d'une locale</i>
$m[x \leftarrow v]_A = (s, g[x \leftarrow v]_L)$	<i>Écriture d'une globale</i>
$m[(n, x) \leftarrow v]_A = (s', g)$	<i>Écriture d'une locale</i>
où $s'_{ l -n+1} = s_{ l -n+1}[x \leftarrow v]_L$	
$\forall i \neq l - n + 1, s'_i = s_i$	

Les numéros de cadre qui permettent d'identifier les globales (le n dans (n, x)) croissent avec la pile. D'autre part, l'empilement se fait en tête de liste (près de l'indice 1). Donc pour accéder aux plus vieilles locales (numérotées 1), il faut accéder au dernier élément de la liste. Ceci explique pourquoi un indice $|l| - n + 1$ apparaît dans la définition précédente.

Définition 4.5 (Accès par champ). *Les valeurs qui sont des structures possèdent des sous-valeurs, associées à des noms de champ.*

L'accesseur $[\cdot]_F$ permet de lire et de modifier un champ de ces valeurs.

C'est une erreur d'accéder à un champ d'une valeur non structure ($4[f]$ par exemple).

$$\begin{aligned} \{f_1 : v_1; \dots; f_n : v_n\}[f_i]_F &= v_i \\ \{f_1 : v_1; \dots; f_n : v_n\}[f_p \leftarrow v]_F &= \{f_1 : v'_1; \dots; f_n : v'_n\} \\ &\text{où } v'_p = v \\ &\quad \forall i \neq p, v'_i = v_i \end{aligned}$$

Définition 4.6 (Accès par indice). *On définit de même un accesseur $[\cdot]_I$ pour les accès par indice à des valeurs tableaux. Néanmoins le paramètre indice est toujours un entier et pas une expression arbitraire.*

$$\begin{aligned} \{v_1; \dots; v_n\}[i]_I &= e_i \\ \{v_1; \dots; v_n\}[i \leftarrow v]_I &= \{v'_1; \dots; v'_n\} \\ &\text{où } v'_i = v \\ &\quad \forall j \neq i, v'_j = v_j \end{aligned}$$

Définition 4.7 (Accès par chemin). *L'accès par chemin $[\cdot]_\Phi$ permet de lire et de modifier la mémoire en profondeur.*

$$\begin{aligned} m[a]_\Phi &= m[a]_A \\ m[*\varphi]_\Phi &= m[\varphi']_\Phi \text{ où } \varphi' = m[\varphi]_\Phi \\ m[\varphi.f]_\Phi &= m[\varphi]_\Phi[f]_F \\ m[\varphi[i]]_\Phi &= m[\varphi]_\Phi[i]_I \\ m[a \leftarrow v]_\Phi &= m[a \leftarrow v]_A \\ m[*\varphi \leftarrow v]_\Phi &= m[\varphi' \leftarrow v]_\Phi \text{ où } \varphi' = m[\varphi]_\Phi \\ m[\varphi.f \leftarrow v]_\Phi &= m[\varphi \leftarrow (m[\varphi]_\Phi[f \leftarrow v]_F)]_\Phi \\ m[\varphi[i] \leftarrow v]_\Phi &= m[\varphi \leftarrow (m[\varphi]_\Phi[i \leftarrow v]_I)]_\Phi \end{aligned}$$

Cette dernière définition mérite une explication. Dans le cas de la lecture, il suffit d'appliquer les bons accesseurs : $[\cdot]_F$ pour $\varphi.f$, etc.

En revanche, la modification est plus complexe. Les deux premiers cas ($\varphi = a$ et $\varphi = *\varphi'$) modifient directement une valeur complète (en modifiant une association), mais les deux suivants ($\varphi = \varphi'.f$ et $\varphi = \varphi'[i]$) ne font qu'altérer une sous-valeur existante. Il est donc nécessaire de procéder en 3 étapes :

- obtenir la valeur à modifier (soit $m[\varphi]_\Phi$)
- construire une valeur altérée (en appliquant par exemple $[f \leftarrow v]_F$)

- affecter cette valeur au même chemin (le $m[\varphi \leftarrow \dots]_\varphi$ externe)

Dans la suite, on notera uniquement $[\cdot]$ tous ces accesseurs lorsque ce n'est pas ambigu.

4.6 Contextes d'évaluation

L'évaluation des expressions repose sur la notion de contextes d'évaluation. L'idée est que si on peut évaluer une expression, alors on peut évaluer une expression qui contient celle-ci.

Par exemple, supposons que $\langle f(3), m \rangle \rightarrow \langle 2, m \rangle$. Alors on peut ajouter la constante 1 à gauche de chaque expression sans changer le résultat : $\langle 1 + f(3), m \rangle \rightarrow \langle 1 + 2, m \rangle$. On a utilisé le même contexte $C = 1 + \bullet$.

Pour pouvoir raisonner en termes de contextes, 3 points sont nécessaires :

- comment découper une expression selon un contextes
- comment appliquer une règle d'évaluation sous un contexte
- comment regrouper une expression et un contexte

Le premier point consiste à définir les contextes eux-mêmes (figure 4.6).

Le deuxième est résolu la règle d'inférence suivante :

$$\frac{\langle e, m \rangle \rightarrow \langle e', m' \rangle}{\langle C[e], m \rangle \rightarrow \langle C[e'], m' \rangle} \text{ (CTX)}$$

Enfin, le troisième revient à définir l'opérateur de substitution $\cdot(\cdot)$ présent dans la règle précédente.

Dans la définition de l'ensemble des contextes, chaque cas hormis le cas de base fait apparaître exactement un "C". Chaque contexte est donc constitué d'exactly un "trou" \bullet (une dérivation de C est toujours linéaire). L'opération de substitution consiste à remplacer ce trou, comme décrit dans la figure 4.7.

Exemple de dérivation.

$$\begin{array}{c} \frac{\langle e_1, m \rangle \rightarrow^* \langle v_1, m' \rangle \quad (C = \bullet \boxplus e_2)}{\langle \bullet \boxplus e_2, m \rangle \rightarrow^* \langle v_1 \boxplus e_2, m' \rangle} \text{ (CTX)} \quad \frac{\langle e_2, m' \rangle \rightarrow^* \langle v_2, m'' \rangle \quad (C = v_1 \boxplus \bullet)}{\langle v_1 \boxplus e_2, m \rangle \rightarrow^* \langle v_1 \boxplus v_2, m' \rangle} \text{ (CTX)} \\ \vdots \quad \vdots \\ \hline \langle e, m \rangle \rightarrow^* \langle v_1 \boxplus v_2, m'' \rangle \quad \langle v_1 \boxplus v_2, m'' \rangle \rightarrow^* \langle v_1 \hat{\boxplus} v_2, m'' \rangle \text{ (EXP-BINOP)} \\ \vdots \quad \vdots \\ \hline \langle e, m \rangle \rightarrow^* \langle v_1 \hat{\boxplus} v_2, m'' \rangle \end{array}$$

Contextes	$C ::= \bullet$
	$C \boxplus e$
	$v \boxplus C$
	$\boxminus C$
	$*C$
	$\varphi.C$
	$\varphi[C]$
	$C[e]$
	$C \leftarrow e$
	$\varphi \leftarrow C$
	$\{f_1 : v_1; \dots; f_i : C; \dots; f_n : e_n\}$
	$[v_1; \dots; C; \dots; e_n]$
	$C(e_1, \dots, e_n)$
	$f(v_1, \dots, C, \dots, e_n)$
	$C; i$
	$\text{IF}(C)\{i_1\}\text{ELSE}\{i_2\}$
	$\text{RETURN}(C)$

FIGURE 4.6 – Contextes d'exécution

4.7 Expressions

Définition 4.8 (Évaluation d'une expression). *L'évaluation d'une expression e se fait sous un état mémoire particulier m et est susceptible de modifier celui-ci en le transformant en un nouveau m' . Le résultat est toujours une valeur v , c'est à dire que nous présentons pour les expressions une sémantique à grands pas. Cette évaluation est notée :*

$$\langle e, m \rangle \rightarrow \langle v, m' \rangle$$

Définition 4.9 (Évaluation d'une left-value). *L'évaluation d'une left-value lv produit un "chemin" φ dans une variable, qui est en fait équivalent à une left-value dont toutes les sous-expressions (d'indices) ont été évaluées.*

On note :

$$\begin{aligned}
& \bullet \langle e_0 \rangle = e_0 \\
& (C \boxplus e) \langle e_0 \rangle = C \langle e_0 \rangle \boxplus e \\
& (v \boxplus C) \langle e_0 \rangle = v \boxplus C \langle e_0 \rangle \\
& (\boxplus C) \langle e_0 \rangle = \boxplus C \langle e_0 \rangle \\
& (*C) \langle e_0 \rangle = * C \langle e_0 \rangle \\
& (\varphi.C) \langle e_0 \rangle = \varphi.C \langle e_0 \rangle \\
& (\varphi[C]) \langle e_0 \rangle = \varphi[C \langle e_0 \rangle] \\
& (C[e]) \langle e_0 \rangle = C \langle e_0 \rangle [e] \\
& (C \leftarrow e) \langle e_0 \rangle = C \langle e_0 \rangle \leftarrow e \\
& (\varphi \leftarrow C) \langle e_0 \rangle = \varphi \leftarrow C \langle e_0 \rangle \\
& \{f_1 : v_1; \dots; f_i : C; \dots; f_n : e_n\} \langle e_0 \rangle = \{f_1 : v_1; \dots; f_i : C \langle e_0 \rangle; \dots; f_n : e_n\} \\
& [v_1; \dots; C; \dots; e_n] \langle e_0 \rangle = [v_1; \dots; C \langle e_0 \rangle; \dots; e_n] \\
& C(e_1, \dots, e_n) \langle e_0 \rangle = C \langle e_0 \rangle (e_1, \dots, e_n) \\
& f(v_1, \dots, C, \dots, e_n) \langle e_0 \rangle = f(v_1, \dots, C \langle e_0 \rangle, \dots, e_n) \\
& (C; i) \langle e_0 \rangle = C \langle e_0 \rangle; i \\
& (\text{IF}(C)\{i_1\}\text{ELSE}\{i_2\}) \langle e_0 \rangle = \text{IF}(C \langle e_0 \rangle)\{i_1\}\text{ELSE}\{i_2\} \\
& (\text{RETURN}(C)) \langle e_0 \rangle = \text{RETURN}(C \langle e_0 \rangle)
\end{aligned}$$

FIGURE 4.7 – Substitution dans les contextes d'évaluation

$$\langle lv, m \rangle \rightarrow \langle \varphi, m' \rangle$$

La sémantique présentée n'a pas d'erreur explicite. Si aucune règle ne peut s'appliquer, on considère que l'exécution est terminée.

Puisque des left-values peuvent apparaître dans les expressions, et des expressions dans les left-values (en indice de tableau), leurs règles d'évaluation sont mutuellement récursives.

Left-values

Obtenir un chemin à partir d'un nom de variable revient à résoudre le nom de cette variable : est-elle accessible? Le nom désigne-t'il une variable locale ou une

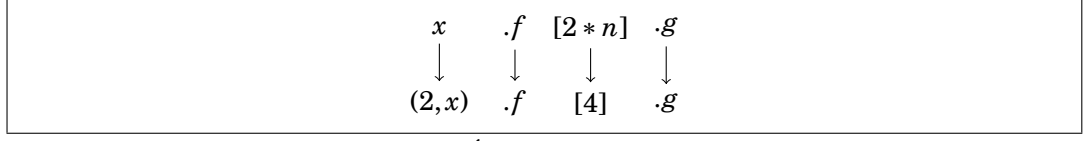


FIGURE 4.8 – Évaluation des left-values.

variable globale ?

$$\frac{a = \text{Lookup}(x, m)}{\langle x, m \rangle \rightarrow \langle a, m \rangle} \text{ (PHI-VAR)}$$

Les règles portant sur le déréférencement et l'accès à un champ de structure sont similaires : on commence par évaluer la left-value sur laquelle porte ce modificateur, et on place le même modificateur sur le chemin résultant.

$$\frac{}{\langle * \varphi, m \rangle \rightarrow \langle \widehat{* \varphi}, m \rangle} \text{ (PHI-DEREF)} \qquad \frac{}{\langle lv.f, m \rangle \rightarrow \langle lv.\widehat{f}, m \rangle} \text{ (PHI-STRUCT)}$$

Enfin, pour évaluer un chemin dans un tableau, on commence par procéder comme précédemment, c'est-à-dire en évaluant la left-value sur laquelle porte l'opération d'indexation. Puis on évalue l'expression d'indice en une valeur qui permet de construire le chemin résultant.

$$\frac{}{\langle \varphi[n], m \rangle \rightarrow \langle \varphi[\widehat{n}], m \rangle} \text{ (PHI-ARRAY)}$$

Notons qu'en procédant ainsi, on évalue les left-values de gauche à droite : dans l'expression $x[e_1][e_2][e_3]$, e_1 est évalué en premier, puis e_2 , puis e_3 .

Un exemple d'évaluation est donné dans la figure 4.8.

Expressions

Évaluer une constante est le cas le plus simple, puisqu'en quelque sorte celle-ci est déjà évaluée. À chaque constante syntaxique c , on peut associer une valeur sémantique \widehat{c} . Par exemple, au chiffre (symbole) 3, on associe le nombre (entier) $\widehat{3}$.

$$\frac{}{\langle c, m \rangle \rightarrow \langle \widehat{c}, m \rangle} \text{ (EXP-CST)}$$

De même, une fonction n'est pas évaluée plus :

$$\frac{}{\langle f, m \rangle \rightarrow \langle \hat{f}, m \rangle} \text{ (EXP-FUN)}$$

Pour lire le contenu d'un emplacement mémoire (left-value), il faut tout d'abord l'évaluer en un chemin.

$$\frac{}{\langle \varphi, m \rangle \rightarrow \langle m[\varphi], m \rangle} \text{ (EXP-LV)}$$

Pour évaluer une expression constituée d'un opérateur, on évalue une expression, puis l'autre (l'ordre d'évaluation, est encore imposé : de gauche à droite). À chaque opérateur \boxplus , correspond un opérateur sémantique $\hat{\boxplus}$ qui agit sur les valeurs. Par exemple, l'opérateur $\hat{+}$ est l'addition classique entre entiers. Afin d'interdire la division par zéro, celle-ci et le modulo sont traités dans une règle à part.

$$\begin{array}{c} \frac{}{\langle \boxminus v, m \rangle \rightarrow \langle \hat{\boxminus} v, m \rangle} \text{ (EXP-UNOP)} \qquad \frac{\boxplus \notin \{/, \%\}}{\langle v_1 \boxplus v_2, m \rangle \rightarrow \langle v_1 \hat{\boxplus} v_2, m \rangle} \text{ (EXP-BINOP)} \\[10pt] \frac{\boxplus \in \{/, \%\} \quad v_2 \neq \hat{0}}{\langle v_1 \boxplus v_2, m \rangle \rightarrow \langle v_1 \hat{\boxplus} v_2, m \rangle} \text{ (EXP-DIV)} \qquad \frac{\boxplus \in \{/, \%\}}{\langle v_1 \boxplus 0, m \rangle \rightarrow \Omega_{div}} \text{ (EXP-DIV-ZERO)} \end{array}$$

Il est nécessaire de dire un mot sur les opérations $\hat{+}_p$ et $\hat{-}_p$ définissant l'arithmétique des pointeurs. Celles-ci sont uniquement définies pour les références mémoire à un tableau, c'est à dire celles qui ont la forme $\varphi[n]$. On a alors :

$$\begin{aligned} \varphi[n] +_p m &= \varphi[n + m] \\ \varphi[n] -_p m &= \varphi[n - m] \end{aligned}$$

Pour prendre l'adresse d'une variable, il suffit de résoudre celle-ci dans l'état mémoire courant.

$$\frac{a = \text{Lookup}(x, m)}{\langle \&x, m \rangle \rightarrow \langle a, m \rangle} \text{ (EXP-ADDR OF)}$$

L'affectation se déroule 3 étapes : d'abord, l'expression est évaluée en une valeur v . Ensuite, la left-value est évaluée en un chemin φ . Enfin, un nouvel état mémoire est construit, où la valeur accessible par φ est remplacée par v . Comme dans le langage C, l'expression d'affectation produit une valeur, qui est celle qui a été affectée.

$$\frac{}{\langle \varphi \leftarrow v, m \rangle \rightarrow \langle v, m[\varphi \leftarrow v] \rangle} \text{ (EXP-SET)}$$

Expressions composées

On commence par définir une opération d'évaluation de plusieurs expressions à la fois : on note

$$\left\langle \begin{pmatrix} e_1 \\ \vdots \\ e_n \end{pmatrix}, m \right\rangle \rightarrow \left\langle \begin{pmatrix} v_1 \\ \vdots \\ v_n \end{pmatrix}, m' \right\rangle$$

si $\exists(m_1, \dots, m_n), \forall i \in [1; n-1], \langle e_i, m_i \rangle \rightarrow \langle e_{i+1}, m_{i+1} \rangle$ avec $m = m_1$ et $m' = m_n$.

Notons que l'évaluation se fait encore de gauche à droite. On utilise la notation vecteur colonne pour signifier qu'il s'agit ici de métasyntaxe (il n'y a pas de tuples dans le langage).

Cette évaluation chaînée est au coeur de la règle suivante qui permet d'évaluer les structures : à une structure (syntaxique) correspond une valeur structurale dont les champs sont ceux de la première structure évalués :

$$\frac{}{\langle \{f_1 : v_1; \dots; f_n : v_n\}, m \rangle \rightarrow \langle \{f_1 : v_1; \dots; f_n : v_n\}, m \rangle} \text{ (EXP-STRUCT)}$$

De même, l'évaluation d'un littéral de tableau se fait en évaluant de gauche à droite ses éléments :

$$\frac{}{\langle [v_1, \dots, v_n], m \rangle \rightarrow \langle [v_1, \dots, v_n], m \rangle} \text{ (EXP-ARRAY)}$$

L'appel de fonction repose également sur cette évaluation multiple. Tout d'abord, les arguments sont évalués et placés dans un nouveau cadre de pile. Puis les expressions qui initialisent les variables locales sont elle aussi évaluées et ajoutées à ce même cadre de pile (opérateur Extend). Ensuite, le corps de la fonction est évalué jusqu'à se réduire en une instruction `RETURN(v)`. Puis, le cadre précédemment utilisé est dépilé.

La dernière étape consiste à nettoyer la mémoire de références à l'ancien cadre de pile. En effet, si une référence au dernier cadre est toujours présente après le retour, elle pourra se résoudre en un objet différent plus tard dans l'exécution du programme.

La fonction Cleanup est donnée par :

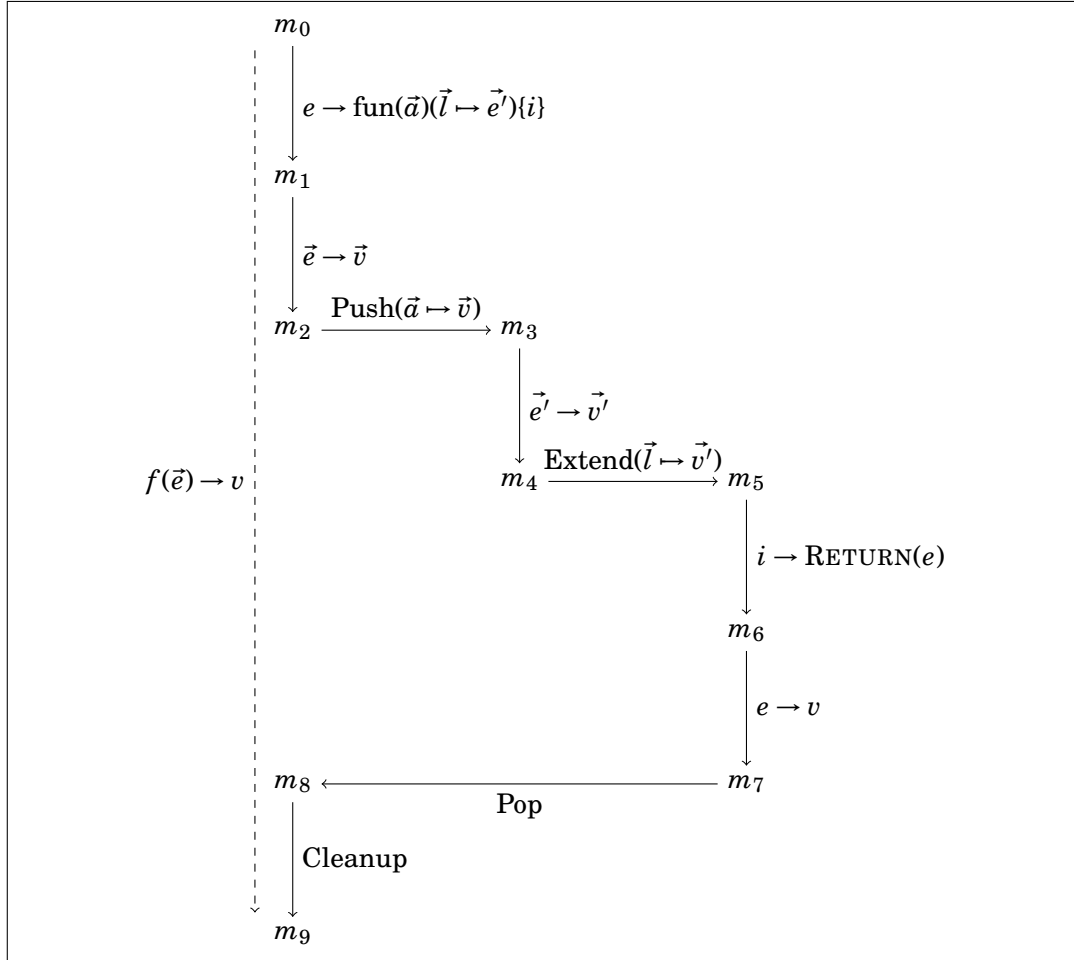


FIGURE 4.9 – L'appel d'une fonction. La taille de la pile croît de gauche à droite, et les réductions se font de haut en bas.

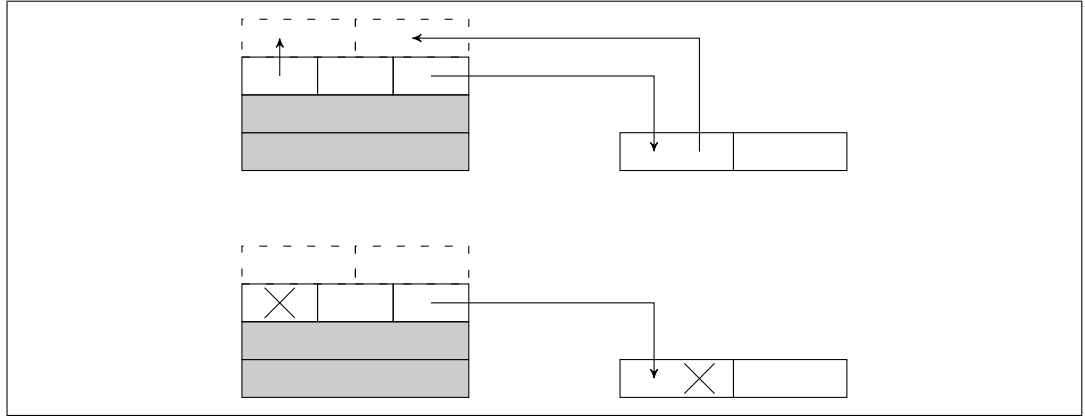


FIGURE 4.10 – Nettoyage d'un cadre de pile

$$\text{Cleanup}(s, g) = (s', g')$$

$$\text{où } g' = \text{CleanupList}(|s|, g)$$

$$s' = [\text{CleanupList}(|s|, s_1), \dots, \text{CleanupList}(|s|, s_n)]$$

$$\text{CleanupList}(p, u) = \{(x, v) \in u / v \text{ n'est pas une adresse}\}$$

$$\cup \{(x, \varphi) \in u / \text{Live}(p, \varphi)\}$$

$$\text{Live}(p, (n, x)) = n < p$$

$$\text{Live}(p, (x)) = \text{Vrai}$$

$$\text{Live}(p, * \varphi) = \text{Live}(p, \varphi)$$

$$\text{Live}(p, \varphi.f) = \text{Live}(p, \varphi)$$

$$\text{Live}(p, \varphi[n]) = \text{Live}(p, \varphi)$$

Sans cette règle, examinons le programme suivant :

$\begin{array}{l} f \ () \\ \quad (x=0) \\ \{ \\ \quad \text{return } (\&x); \\ \} \end{array}$	$\begin{array}{l} g \ (p) \\ \quad (x=0.0) \\ \{ \\ \quad *p = 1; \\ \} \end{array}$	$\begin{array}{l} h \ () \\ \quad (p=f()) \\ \{ \\ \quad g(p); \\ \} \end{array}$
---	--	---

L'exécution de $h()$ donne à p la valeur $(1, x)$. Donc en arrivant dans g , le déréférencement de p va modifier x .

$$\begin{array}{c}
f = \text{fun}(a_1, \dots, a_n)((l'_1, e'_1), \dots, (l'_p, e'_p))\{i\} \\
m_1 = \text{Push}(m_0, ((a_1, v_1), \dots, (a_n, v_n))) \\
\left\langle \begin{pmatrix} e'_1 \\ \vdots \\ e'_p \end{pmatrix}, m_1 \right\rangle \rightarrow \left\langle \begin{pmatrix} v'_1 \\ \vdots \\ v'_p \end{pmatrix}, m_2 \right\rangle \quad m_3 = \text{Extend}(m_2, ((l_1, v_1), \dots, (l_n, v_n))) \\
\frac{\langle i, m_3 \rangle \rightarrow \langle \text{RETURN}(v), m_4 \rangle \quad m_5 = \text{Pop}(m_4) \quad m_6 = \text{Cleanup}(m_5)}{\langle f(v_1, \dots, v_n), m_0 \rangle \rightarrow \langle v, m_6 \rangle} \text{ (EXP-CALL)}
\end{array}$$

Cette évaluation est décrite dans la figure 4.9.

4.8 Instructions

Contrairement à l'évaluation des expressions, on choisit une sémantique de réécriture à petits pas. La sémantique fonctionne de la manière suivante : partant d'un état mémoire m , on veut exécuter une instruction i . Les règles d'évaluation suivantes permettent de réduire le problème en se ramenant à l'exécution d'une instruction i' "plus simple" en partant d'un état mémoire m' . Un tel pas est noté :

$$\langle i, m \rangle \rightarrow \langle i', m' \rangle$$

Par exemple, exécuter $x \leftarrow 3; y \leftarrow x$ revient à évaluer $y \leftarrow x$ depuis un état mémoire dans lequel on a déjà réalisé la première affectation. La seconde affectation se réalise de même et permet de réécrire l'instruction restante en PASS :

$$\begin{aligned}
\langle (x \leftarrow 3; y \leftarrow x), m \rangle &\rightarrow \langle y \leftarrow x, m[x \mapsto \widehat{3}] \rangle \\
&\rightarrow \langle \text{PASS}, m[x \mapsto \widehat{3}][y \mapsto \widehat{3}] \rangle
\end{aligned}$$

Il n'est pas possible de réduire plus loin l'instruction PASS. Dans un tel cas, l'évaluation est terminée.

Les seuls cas terminaux sont PASS et RETURN(e).

Les cas de la séquence et de l'affectation ont été utilisés dans l'exemple ci-dessus.

$$\begin{array}{c}
\frac{\langle i, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle (i; i'), m \rangle \rightarrow \langle i', m' \rangle} \text{ (SEQ)} \quad \frac{}{\langle (\text{PASS}; i), m \rangle \rightarrow \langle i, m \rangle} \text{ (PASS)} \\
\frac{}{\langle v, m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{ (EXP)}
\end{array}$$

Pour traiter l'alternative, on a besoin de 2 règles. Elles commencent de la même manière, en évaluant la condition. Si le résultat est 0 (et seulement dans ce cas), c'est la règle IF-FALSE qui est appliquée et l'instruction revient à évaluer la branche "else". Dans les autres cas, c'est la règle IF-TRUE qui s'applique et la branche "then" qui est prise.

$$\frac{}{\langle \text{IF}(0)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_f, m \rangle} \text{ (IF-FALSE)}$$

$$\frac{v \neq 0}{\langle \text{IF}(v)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_t, m \rangle} \text{ (IF-TRUE)}$$

Pour traiter la boucle, on peut être tenté de procéder de la même manière :

$$\frac{v \neq 0}{\langle \text{WHILE}(v)\{i\}, m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{ (WHILE-FALSE-BAD)}$$

$$\frac{}{\langle \text{WHILE}(0)\{i\}, m \rangle \rightarrow \langle i; \text{WHILE}(e)\{i\}, m' \rangle} \text{ (WHILE-TRUE-BAD)}$$

Mais la seconde règle est impossible : puisque e a déjà été évaluée, il est impossible de la réintroduire non évaluée en partie droite.

À la place, on exprime la sémantique de la boucle comme une simple règle de réécriture :

$$\frac{}{\langle \text{WHILE}(e)\{i\}, m \rangle \rightarrow \langle \text{IF}(e)\{i; \text{WHILE}(e)\{i\}\}, m \rangle} \text{ (WHILE)}$$

Cette règle revient à dire qu'on peut dérouler une boucle. Pour la comprendre, on peut remarquer qu'une boucle "while" est en réalité équivalente une infinité de "if" imbriqués.

<pre>while(e) { i }</pre>	\cong	<pre>if(e) { i; if(e) { i; if(e) { i; if(e) { i; ... } } } }</pre>
-----------------------------	---------	--

Donc en remplaçant le second "if" par le "while", on obtient :

$$\begin{array}{ccc} \text{while}(e) \{ & & \text{if}(e) \{ \\ \quad i & \hat{=} & \quad i; \\ \} & & \text{while}(e) \{ \\ & & \quad i \\ & & \} \end{array}$$

Enfin, si un "return" apparaît dans une séquence, on peut supprimer la suite :

$$\frac{}{\langle \text{RETURN}(v); i, m \rangle \rightarrow \langle \text{RETURN}(e), m \rangle} \text{ (RETURN)}$$

4.9 Phrases

Un programme est constitué d'une suite de phrases : déclarations de fonctions, de variables et de types, et évaluation d'expressions.

Il est donc logique que l'évaluation d'une phrase fasse passer d'un état mémoire à un autre :

$$m \vdash p \rightarrow m'$$

La définition d'une structure est ignorée. En effet, celle-ci ne sert qu'au typage.

$$\frac{}{m \vdash \text{struct } s\{\dots\} \rightarrow m} \text{ (PH-STRUCT)}$$

L'évaluation d'une expression est uniquement faite pour ses effets de bord. Par exemple, après avoir défini les fonctions du programme, on pourra appeler main().

$$\frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{m \vdash e \rightarrow m'} \text{ (PH-EXP)}$$

La déclaration d'une variable globale (avec un initialiseur) consiste à évaluer cet initialiseur et à étendre l'état mémoire avec ce couple (variable, valeur).

$$\frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{(s, g) \vdash x = e \rightarrow (s, (x, v) :: g)} \text{ (PH-VAR)}$$

4.10 Exécution

L'exécution d'un programme est sans surprise l'exécution de ses phrases, les unes à la suite des autres.

On commence par étendre l'extension \rightarrow^* au listes de la relation \rightarrow :

$$\frac{}{m \vdash [] \rightarrow^* m} \text{ (PH* - NIL)} \qquad \frac{m \vdash p \rightarrow m' \quad m' \vdash ps \rightarrow^* m''}{m \vdash p :: ps \rightarrow^* m''} \text{ (PH* - CONS)}$$

L'exécution d'un programme est alors :

$$([], []) \vdash P \rightarrow^* m$$

4.11 Exemple : l'algorithme d'Euclide

Version par divisions successives :

```
function gcd(a, b)
  var t = 0;
  while b != 0
    t = b
    b = a mod b
    a = t
  return a
```

Soit :

$$f(a, b)(t = 0)\{\text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}$$

$$\langle f(1071, 462), m \rangle \rightarrow ?$$

$$\langle \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a), m[a \mapsto 1071][b \mapsto 462][t \mapsto 0] \rangle \rightarrow ?$$

(on notera cet état $s_0 = \langle i_0, m_0 \rangle$)

$$\langle a = 0, m_0 \rangle \rightarrow \langle 0, m_0 \rangle$$

donec

$$\langle \text{IF}(a = 0)\{\text{RETURN}(b)\}, m_0 \rangle \rightarrow \langle \text{PASS}, m[a \mapsto 1071][b \mapsto 462] \rangle$$

$$s_0 \rightarrow \langle \text{IF}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_0 \rangle \quad (4.1)$$

$$\rightarrow \langle t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_0 \rangle \quad (4.2)$$

$$\rightarrow \langle b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_0 \rangle \quad (4.3)$$

$$\rightarrow \langle a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_0'' \rangle \quad (4.4)$$

$$\rightarrow \langle \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_1 \rangle \quad (4.5)$$

$$\rightarrow \langle \text{IF}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_1 \rangle \quad (4.6)$$

$$\rightarrow \langle t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_1 \rangle \quad (4.7)$$

$$\rightarrow \langle \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_2 \rangle \quad (4.8)$$

$$\rightarrow \langle \text{IF}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_2 \rangle \quad (4.9)$$

$$\rightarrow \langle t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_2 \rangle \quad (4.10)$$

$$\rightarrow \langle \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_3 \rangle \quad (4.11)$$

$$\rightarrow \langle \text{IF}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t; \text{WHILE}(b \neq 0)\{t \leftarrow b; b \leftarrow a \% b; a \leftarrow t\}; \text{RETURN}(a)\}, m_3 \rangle \quad (4.12)$$

$$\rightarrow \langle \text{PASS}; \text{RETURN}(a)\}, m_3 \rangle \quad (4.13)$$

$$\rightarrow \langle \text{RETURN}(a)\}, m_3 \rangle \quad (4.14)$$

$$m_0' = m_0[t \mapsto 462] = m[a \mapsto 1071][b \mapsto 462][t \mapsto 462]$$

$$m_0'' = m_0'[b \mapsto 147] = m[a \mapsto 1071][b \mapsto 147][t \mapsto 462]$$

$$m_1 = m_0''[a \mapsto 462] = m[a \mapsto 462][b \mapsto 147][t \mapsto 462]$$

$$m_2 = m_1[t \mapsto 147][b \mapsto 21][a \mapsto 147] = m[a \mapsto 147][b \mapsto 21][t \mapsto 147]$$

$$m_3 = m_2[t \mapsto 21][b \mapsto 0][a \mapsto 21] = m[a \mapsto 21][b \mapsto 0][t \mapsto 21]$$

TODO

- changer les para de présentation des règles
- cas d'erreurs explicites Ω_{div} , Ω_{ptr} et Ω_{array}
- widehats sur les constantes ?
- liste d'assos -> fonction
- interdire d'avoir plusieurs variables qui ont le même nom dans un cadre
- dédupliquer la def de l'état mémoire
- syntaxe concrète
- procédures vs fonctions
- top ?
- return implicite en fin de fct

TYPAGE

Dans ce chapitre, nous enrichissons le langage défini dans le chapitre 4 d'un système de types. Celui-ci permet de séparer les programmes bien formés, comme celui de la figure 5.1a des programmes mal formés comme celui de la figure 5.1b.

Le but d'un tel système de types est de rejeter les programmes qui sont "évidemment faux", c'est à dire dont on peut prouver qu'il provoqueraient des erreurs à l'exécution dues à une incompatibilité entre valeurs. En ajoutant cette étape, on restreint la classe d'erreurs qui pourraient bloquer la sémantique.

5.1 Principe

Le principe est d'associer à chaque construction syntaxique une étiquette représentant le genre de valeurs qu'elle produira. Dans le programme de la figure 5.1a, la variable x est initialisée avec la valeur 0, c'est donc un entier. Cela signifie que dans

<pre>f () (x=0) { x = 1 return x }</pre>	<pre>f () (x=0) { x = 1 return (*x) }</pre>
(a) Programme bien formé	(b) Programme mal formé

FIGURE 5.1 – Programmes bien et mal formés

Type	$t ::= \text{INT}$	Entier
	FLOAT	Flottant
	$t[]$	Tableau
	$t*$	Pointeur
	$\text{struct } s$	Structure
	$(t_1, \dots, t_n) \rightarrow t_r$	Fonction
Environnement de typage	$\Gamma ::= []$	Environnement vide
	$(a, t) :: \Gamma'$	Extension

FIGURE 5.2 – Types et environnements de typage

tout le programme, toutes les instances de cette variable¹ porteront ce type. La première instruction est l'affectation de la constante 1 (entière) à x dont on sait qu'elle porte des valeurs entières, ce qui est donc correct. Le fait de rencontrer $\text{RETURN}(x)$ permet de conclure que le type de la fonction est $() \rightarrow \text{INT}$.

Dans la seconde fonction, au contraire, l'opérateur $*$ est appliqué à x (le début de l'analyse est identique et permet de conclure que x porte des valeurs entières). Or cet opérateur prend un argument d'un type pointeur de la forme $t*$ et renvoie alors une valeur de type t . Ceci est valable pour tout t (INT , FLOAT où même t' : le déréférencement d'un pointeur sur pointeur donne un pointeur), mais le type de x , INT , n'est pas de cette forme. Ce programme est donc mal typé.

5.2 Définitions

Les types associés aux expressions sont décrits dans la figure 5.2.

Pour maintenir les contextes de typage, un environnement Γ associe un type à un ensemble de variables.

Plus précisément, un environnement Γ est une liste de couples (variable, type).

Par exemple, $(p, \text{INT}*) \in \Gamma$ permet de typer (sous Γ) l'expression p en $\text{INT}*$, $*p$ en INT et $p +_p 4$ en $\text{INT}*$.

Le type des fonctions semble faire apparaître un n -uplet (t_1, \dots, t_n) mais ce n'est qu'une notation : il n'y a pas de n -uplets de première classe, ils sont toujours présents dans un type fonctionnel.

1. Deux variables peuvent avoir le même nom dans deux fonctions différentes, par exemple. Dans ce cas il n'y a aucune contrainte particulière entre ces deux variables. L'analyse de typage se fait toujours dans un contexte précis.

Définition 5.1 (Typage d'une expression). *On note de la manière suivante le fait qu'une expression e (telle que définie dans la figure 4.1) ait pour type t dans le contexte Γ .*

$$\Gamma \vdash e : t$$

Définition 5.2 (Typage d'une instruction). *Les instructions n'ont en revanche pas de type. Mais il est tout de même nécessaire de vérifier que toutes les sous-expressions apparaissant dans une instruction sont cohérentes ensemble.*

On note de la manière suivante le fait que sous l'environnement Γ l'instruction i est bien typée :

$$\Gamma \vdash i$$

Définition 5.3 (Typage d'une phrase). *De par leur nature séquentielle, les phrases qui composent un programme altèrent l'environnement de typage. Par exemple, la déclaration d'une variable globale ajoute une valeur dans l'environnement.*

On note

$$\Gamma \vdash p \rightarrow \Gamma'$$

si le typage de la phrase p transforme l'environnement Γ en Γ' .

On étend cette notation aux suites de phrases :

$$\left\{ \begin{array}{l} \Gamma \vdash [] \rightarrow^* [] \\ \Gamma \vdash p :: ps \rightarrow^* p :: ps \text{ si } \exists \Gamma'', \left\{ \begin{array}{l} \Gamma \vdash p \rightarrow \Gamma'' \\ \Gamma'' \vdash ps \rightarrow^* ps \end{array} \right. \end{array} \right.$$

Définition 5.4 (Typage d'un programme). *Un programme est bien typé si on peut typer sa suite de phrases en partant d'un environnement vide. C'est à dire s'il existe un environnement final Γ_f tel que*

$$[] \vdash P \rightarrow^* P$$

Cela est indépendant de tout environnement ; on note alors :

$$\vdash P$$

5.3 Expressions

Littéraux

Le typage des littéraux numériques ne dépend pas de l'environnement de typage : ce sont toujours des entiers ou des flottants.

$$\frac{}{\Gamma \vdash i : \text{INT}} \text{ (CST-INT)} \qquad \frac{}{\Gamma \vdash f : \text{FLOAT}} \text{ (CST-FLOAT)}$$

Le pointeur nul, quant à lui, est compatible avec tous les types pointeur.

$$\frac{}{\Gamma \vdash \text{NULL} : t*} \text{ (CST-NULL)}$$

Left-values

Rappelons que l'environnement de typage Γ contient le type des variables accessibles du programme. Le cas où la left-value à typer est une variable est donc direct : il suffit de retrouver son type dans l'environnement.

$$\frac{x : t \in \Gamma}{\Gamma \vdash x : t} \text{ (LV-VAR)}$$

Dans le cas d'un déréférencement, on commence par typer la left-value déréférencée. Si elle a un type pointeur, la valeur déréférencée est du type pointé.

$$\frac{\Gamma \vdash lv : t*}{\Gamma \vdash *lv : t} \text{ (LV-DEREF)}$$

Pour une left-value indexée (l'accès à tableau), on s'assure que l'indice soit entier, et que la left-value a un type tableau : le type de l'élément est encore une fois le type de base du type tableau (t pour $t[]$).

$$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash lv : t[]}{\Gamma \vdash lv[e] : t} \text{ (LV-INDEX)}$$

$$\frac{(s, f, t_f) \in S \quad \Gamma \vdash lv : \text{struct } s}{\Gamma \vdash lv.f : t_f} \text{ (LV-FIELD)}$$

Opérateurs

Un certain nombre d'opérations est possible sur le type INT.

$$\frac{\boxplus \in \{+, -, \times, /, \&, |, ^, \&\&, ||, \ll, \gg\} \quad \Gamma \vdash e_1 : \text{INT} \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-INT)}$$

De même sur FLOAT.

$$\frac{\boxplus \in \{+, -, \times, /, .\} \quad \Gamma \vdash e_1 : \text{FLOAT} \quad \Gamma \vdash e_2 : \text{FLOAT}}{\Gamma \vdash e_1 \boxplus e_2 : \text{FLOAT}} \text{ (OP-FLOAT)}$$

Les opérateurs de comparaison peuvent s'appliquer à deux opérandes qui sont d'un type qui supporte l'égalité. Ceci est représenté par un jugement $\text{Eq}(t)$ qui est vrai pour les types INT, FLOAT et pointeurs. Les opérateurs $=$ et \neq renvoient alors un INT.

$$\frac{t \in \{\text{INT}, \text{FLOAT}\}}{\text{Eq}(t)} \text{ (EQ-NUM)} \quad \frac{\text{Eq}(t)}{\text{Eq}(t*)} \text{ (EQ-PTR)} \quad \frac{\text{Eq}(t)}{\text{Eq}(t[]) } \text{ (EQ-ARRAY)}$$

$$\frac{\boxplus \in \{=, \neq\} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \quad \text{Eq}(t)}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-EQ)}$$

Les comparaisons sont plus restrictives, et ne s'appliquent qu'aux types primitifs (on ne peut pas comparer deux pointeurs, ou deux tableaux).

$$\frac{\boxplus \in \{=, \neq, \leq, \geq, <, >\} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \quad t \in \{\text{INT}, \text{FLOAT}\}}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-COMPARABLE)}$$

Les opérateurs unaires "+" et "-" appliquent aux INT, et leurs équivalents "+." et "-." aux FLOAT.

$$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash +e : \text{INT}} \text{ (UNOP-PLUS-INT)} \quad \frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash +.e : \text{FLOAT}} \text{ (UNOP-PLUS-FLOAT)}$$

$$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash -e : \text{INT}} \text{ (UNOP-MINUS-INT)} \quad \frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash -.e : \text{FLOAT}} \text{ (UNOP-MINUS-FLOAT)}$$

Les opérateurs de négation unaires, en revanche, ne s'appliquent qu'aux entiers.

$$\frac{\Box \in \{\sim, !\} \quad \Gamma \vdash e : \text{INT}}{\Gamma \vdash \Box e : \text{INT}} \text{ (UNOP-NOT)}$$

L'arithmétique de pointeurs préserve le type des pointeurs.

$$\frac{\Box \in \{+_p, -_p\} \quad \Gamma \vdash e_1 : t^* \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \Box e_2 : t^*} \text{ (PTR-ARITH)}$$

Autres expressions

Prendre l'adresse d'une left-value rend un type pointeur sur le type de celle-ci.

$$\frac{\Gamma \vdash lv : t}{\Gamma \vdash \&lv : t^*} \text{ (ADDR)}$$

Pour typer une affectation, on vérifie que la left-value (à gauche) et l'expression (à droite) ont le même type. C'est alors le type résultat de l'expression d'affectation.

$$\frac{\Gamma \vdash lv : t \quad \Gamma \vdash e : t}{\Gamma \vdash lv \leftarrow e : t} \text{ (SET)}$$

Un littéral tableau a pour type $t[]$ où t est le type de chacun de ses éléments.

$$\frac{\forall i \in [1; n], \Gamma \vdash e_i : t}{\Gamma \vdash \{e_1; \dots; e_n\} : t[]} \text{ (ARRAY)}$$

Pour qu'un littéral de structure soit bien typé, il faut que chacun de ses noms de champs corresponde à un même nom de type structure, avec le bon type pour chaque champ.

$$\frac{\forall i \in [1; n], \Gamma \vdash e_i : t_i \quad \forall i \in [1; n], (s, f_i, t_i) \in S}{\Gamma \vdash \{f_1 : e_1; \dots; f_n : e_n\} : \text{struct } s} \text{ (STRUCT)}$$

Pour typer un appel de fonction, on s'assure que la fonction a bien un type fonctionnel. On type alors chacun des arguments avec le type attendu. Le résultat est du type de retour de la fonction.

$$\frac{\Gamma \vdash e : (t_1, \dots, t_n) \rightarrow t \quad \forall i \in [1; n], \Gamma \vdash e_i : t_i}{\Gamma \vdash e(e_1, \dots, e_n) : t} \text{ (CALL)}$$

5.4 Instructions

La séquence est simple à traiter : l'instruction vide est toujours bien typée, et la suite de deux instructions est bien typée si celles-ci le sont également.

$$\frac{}{\Gamma \vdash \text{PASS}} \text{ (PASS)} \qquad \frac{\Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash i_1; i_2} \text{ (SEQ)}$$

Une instruction constituée d'une expression est bien typée si celle-ci peut être typée dans ce même contexte.

$$\frac{\Gamma \vdash e : t}{\Gamma \vdash e} \text{ (EXP)}$$

Les constructions de contrôle sont bien typées si leurs sous-instructions sont bien typées, et si la condition est d'un type entier.

$$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash \text{IF}(e)\{i_1\}\text{ELSE}\{i_2\}} \text{ (IF)} \qquad \frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i}{\Gamma \vdash \text{WHILE}(e)\{i\}} \text{ (WHILE)}$$

5.5 Fonctions

Le typage des fonctions fait intervenir une variable virtuelle \underline{R} . Cela revient à typer l'instruction $\text{RETURN}(e)$ comme $\underline{R} \leftarrow e$.

$$\frac{\Gamma \vdash \underline{R} \leftarrow e}{\Gamma \vdash \text{RETURN}(e)} \text{ (RETURN)}$$

Pour typer une définition de fonction, on commence par créer un nouvel environnement de typage Γ' obtenu par la suite d'opérations suivantes :

- on enlève (s'il existe) le couple $\underline{R} : t_f$ correspondant à la valeur de retour de la fonction appelante
- on ajoute les types des arguments $a_i : t_i$
- on ajoute les types des variables locales $l_i : t'_i$
- on ajoute le type de la valeur de retour de la fonction appelée, $\underline{R} : t$

Il reste alors à vérifier que les initialiseurs e_i ont le bon type t'_i et que le corps de la fonction est bien typé sous Γ' . Le type de la fonction est alors $(t_1, \dots, t_n) \rightarrow t$.

$$\frac{\begin{array}{c} \Gamma' = (\Gamma - \underline{R}), a_1 : t_1, \dots, a_n : t_n, l_1 : t'_1, \dots, l_n : t'_n, \underline{R} : t \\ \forall i \in [1; p], \Gamma \vdash e_i : t'_i \quad \Gamma' \vdash i \end{array}}{\Gamma \vdash \text{fun}(a_1, \dots, a_n)((l_1, e_1), \dots, (l_p, e_p))\{i\} : (t_1, \dots, t_n) \rightarrow t} \text{ (FUN)}$$

5.6 Phrases

L'évaluation d'une expression est le cas le plus simple. En effet, il y a juste à vérifier que celle-ci est bien typable (avec ce type) dans l'environnement de départ. L'environnement n'est pas modifié.

$$\frac{\Gamma \vdash e : t}{\Gamma \vdash e \rightarrow \Gamma} \text{ (PH-EXP)}$$

La déclaration d'une variable globale commence de la même manière, mais on enrichit l'environnement de cette nouvelle valeur.

$$\frac{\Gamma \vdash e : t \quad \Gamma' = (x, t), \Gamma}{\Gamma \vdash p \rightarrow \Gamma'} \text{ (PH-VAR)}$$

$$\frac{S' = (x_1, t_1, s), \dots, (x_n, t_n, s)}{\Gamma \vdash \text{struct } s\{x_1 : t_1; \dots; x_n : t_n\} \rightarrow \Gamma} \text{ (PH-STRUCT)}$$

5.7 Sûreté du typage

Comme dit l'adage :

"Well-typed programs don't go wrong." (Robin Milner)

C'est à dire qu'un programme bien typé possède des propriétés de sûreté.

Progrès : l'évaluation d'un terme bien typé ne reste pas bloquée ; il y a toujours une règle qui s'applique.

Préservation : l'évaluation d'un terme bien typé produit un terme bien typé.

Puisqu'il s'agit de propriétés reliant la syntaxe à la sémantique, deux types d'environnement sont utilisés en même temps. D'une part, les environnement de typage Γ pour la syntaxe, et d'autre part les environnements mémoire m pour la sémantique. Il est nécessaire de définir une relation de compatibilité entre ces deux monde, pour exprimer par exemple qu'un état mémoire contient dans la variable x un entier.

Définition 5.5 (Compatibilité mémoire). *Soient Γ un environnement de typage et $m = x_1 \mapsto v_1, \dots, x_n \mapsto v_n$ un état mémoire. On dit que m est compatible avec Γ si*

$$\exists(t_1, \dots, t_n), \forall i \in [1; n], \begin{cases} \Gamma \vdash x_i : t_i \\ \Gamma \vdash v_i : t_i \end{cases}$$

On note alors $\Gamma \vdash_{mem} m$.

Théorème 5.1 (Progrès). *Supposons que $\Gamma \vdash e : t$. Soit m un état mémoire tel que $\Gamma \vdash_{mem} m$. Alors l'un des cas suivant est vrai :*

- *e est une valeur : $\exists v, e = v$*
- *il existe e' et m' tels que $\langle e, m \rangle \rightarrow \langle e', m' \rangle$*
- *une erreur d'accès tableau ou pointeur se produit*

Démonstration. On procède par induction sur la dérivation de $\Gamma \vdash e : t$, et plus précisément sur la dernière règle utilisée.

- LV-INDEX
PHI-INDEX s'applique, ou alors une erreur d'indice se produit.
- LV-FIELD
PHI-STRUCT s'applique.
- OP-INT
EXP-BINOP s'applique.
Les cas OP-FLOAT, OP-EQ, et OP-COMPARABLE sont similaires.
- UNOP-PLUS-INT
EXP-UNOP s'applique.
Les cas UNOP-MINUS-INT, UNOP-PLUS-FLOAT, UNOP-MINUS-FLOAT, et UNOP-NOT sont similaires.
- PTR-ARITH
EXP-BINOP s'applique ou une erreur se produit
- ADDR
EXP-ADDR-OF s'applique.

- SET
EXP-LV s'applique.
- ARRAY EXP-ARRAY s'applique.
- STRUCT EXP-STRUCT s'applique.
- CALL EXP-CALL s'applique.

□

Démonstration avec une sémantique par réduction. On procède par induction sur la dérivation de $\Gamma \vdash e : t$, et plus précisément sur la dernière règle utilisée.

CST-INT : e est alors de la forme n , qui est une valeur.

CST-FLOAT : e est alors de la forme f , qui est une valeur.

CST-NULL : e est alors égal à NULL, qui est une valeur.

LV-VAR : Puisque $(x, t) \in \Gamma$ $\Gamma \vdash_{mem} m$, il existe $(x \mapsto v) \in m$. La règle d'évaluation PHI-VAR s'applique donc.

LV-DEREF : On applique l'hypothèse de récurrence à lv .

- si $lv = \varphi$, alors la règle PHI-DEREF s'applique : et $\langle *lv, m \rangle \rightarrow \langle \hat{*}\varphi, m \rangle$
- $\exists(lv', m'), \langle lv, m \rangle \rightarrow \langle lv', m' \rangle$. On peut alors utiliser la règle CTX avec $C = * \bullet$, ce qui donne : $\langle *lv, m \rangle \rightarrow \langle *lv', m' \rangle$.

LV-INDEX :

LV-FIELD :

OP-INT : Cela implique que $e = e_1 \boxplus e_2$. Par le lemme 5.1, on en déduit que $\Gamma \vdash e_1 : \text{INT}$ et $\Gamma \vdash e_2 : \text{INT}$.

Appliquons l'hypothèse de récurrence sur e_1 . Trois cas peuvent se produire.

- $e_1 = v_1$. On a alors $\langle e_1, m \rangle \rightarrow \langle v_1, m' \rangle$ avec $m' = m$.
On applique l'hypothèse de récurrence à e_2 .
 - $e_2 = v_2$: alors $\langle e_2, m' \rangle \rightarrow \langle v_2, m'' \rangle$ avec $m'' = m$. On peut alors appliquer EXP-BINOP.

- $\exists(e'_2, m''), \langle e_2, m' \rangle \rightarrow \langle e'_2, m'' \rangle$.
En appliquant CTX avec $C = v_1 \boxplus \bullet$, on en déduit $\langle v_1 \boxplus e_2, m' \rangle \rightarrow \langle v_1 \boxplus e'_2, m'' \rangle$
soit $\langle e, m \rangle \rightarrow \langle v_1 \boxplus e'_2, m'' \rangle$.
- $\langle e_2, m' \rangle \rightarrow \Omega$. On pose alors $\langle e, m \rangle \rightarrow \Omega$.
- $\exists(e', m'), \langle e_1, m \rangle \rightarrow \langle e'_1, m' \rangle$. En appliquant CTX avec $C = \bullet \boxplus e_2$, on obtient
 $\langle e_1 \boxplus e_2, m \rangle \rightarrow \langle e'_1 \boxplus e_2, m' \rangle$, ou $\langle e, m \rangle \rightarrow \langle e'_1 \boxplus e_2, m' \rangle$.
- $\langle e_1, m \rangle \rightarrow \Omega$. On pose alors $\langle e, m \rangle \rightarrow \Omega$.

OP-FLOAT : Ce cas est similaire à OP-INT.

OP-EQ : Ce cas est similaire à OP-INT.

OP-COMPARABLE : Ce cas est similaire à OP-INT.

UNOP-PLUS-INT : Alors $e = + e_1$. En appliquant l'hypothèse d'induction sur e_1 :

- soit $e_1 = v_1$. Alors en appliquant EXP-UNOP, $\langle + v_1, m \rangle \rightarrow \langle \hat{+} v_1, m \rangle$, c'est à dire en posant $v = \hat{+} v_1$, $\langle e, m \rangle \rightarrow \langle v, m \rangle$.
- soit $\exists e'_1, m', \langle e_1, m \rangle \rightarrow \langle e'_1, m' \rangle$. Alors en appliquant CTX avec $C = + \bullet$, on obtient $\langle e, m \rangle \rightarrow \langle e'_1, m' \rangle$.
- soit $\langle e_1, m \rangle \rightarrow \Omega$. Alors on pose $\langle e, m \rangle \rightarrow \Omega$.

UNOP-PLUS-FLOAT : Ce cas est similaire à UNOP-PLUS-INT.

UNOP-MINUS-INT : Ce cas est similaire à UNOP-PLUS-INT.

UNOP-MINUS-FLOAT : Ce cas est similaire à UNOP-PLUS-INT.

UNOP-NOT :

PTR-ARITH :

ADDR :

SET :

ARRAY :

STRUCT :

CALL :

FUN :

□

Les règles précédentes ont la particularité suivante : pour chaque forme syntaxique, il n'y a souvent qu'une règle qui peut s'appliquer. Cela permet de déduire quelle règle il faut appliquer pour vérifier (ou inférer) le type d'une expression.

Lemme 5.1 (Inversion). *À partir d'un jugement de typage, on peut en déduire des informations sur les types de ses sous-expressions.*

• *Références mémoire :*

- si $\Gamma \vdash x : t$, $x : t \in \Gamma$
- si $\Gamma \vdash * \varphi : t$, alors $\Gamma \vdash \varphi : t*$
- si $\Gamma \vdash \varphi[] : t$, alors $\Gamma \vdash \varphi : t[]$
- si $\Gamma \vdash \varphi.f : t$, alors $\Gamma \vdash \varphi : \{f : t; \dots\}$

• *Appel de fonction :* si $\Gamma \vdash e(e_1, \dots, e_n) : t$, il existe (t_1, \dots, t_n) tels que

$$\begin{cases} \Gamma \vdash e : (t_1, \dots, t_n) \rightarrow t \\ \forall i \in [1; n], \Gamma \vdash e_i : t_i \end{cases}$$

• *Fonction :* si $\Gamma \vdash \text{fun}(a_1, \dots, a_n)((l_1, e_1), \dots, (l_p, e_p))\{i\} : t$, alors il existe (t_1, \dots, t_n) et t' tels que $t' = (t_1, \dots, t_n) \rightarrow t$.

• *Constantes*

- si $\Gamma \vdash n : t$, alors $t = \text{INT}$
- si $\Gamma \vdash f : t$, alors $t = \text{FLOAT}$
- si $\Gamma \vdash \text{NULL} : t$, alors $\exists t', t = t'*$

Démonstration. Pour chaque jugement, on considère les règles qui peuvent amener à cette conclusion.

• *Références mémoire :*

- $\Gamma \vdash x : t$

La seule règle de cette forme est LV-VAR. Puisque sa prémisse est vraie, on en conclut que $x : t \in \Gamma$.

- $\Gamma \vdash * \varphi : t$

De même, seule la règle LV-DEREF convient. On en conclut que $\Gamma \vdash \varphi : t*$.

- $\Gamma \vdash \varphi[] : t$
Idem avec LV-INDEX.
- $\Gamma \vdash \varphi.f : t$
 $\Gamma \vdash \varphi : \{f : t; \dots\}$
- Appel de fonction : pour en arriver à $\Gamma \vdash e(e_1, \dots, e_n) : t$, seule la règle CALL s'applique, ce qui permet de conclure.
- Fonction : la seule règle possible pour conclure une dérivation de

$$\Gamma \vdash \text{fun}(a_1, \dots, a_n)((l_1, e_1), \dots, (l_p, e_p))\{i\} : t$$

est FUN.

□

Il est aussi possible de réaliser l'opération inverse : à partir du type d'une valeur, on peut déterminer sa forme syntaxique. C'est bien sûr uniquement possible pour les valeurs, pas pour n'importe quelle expression (par exemple l'expression x (variable) peut avoir n'importe quel type t dans le contexte $\Gamma = x : t$).

Lemme 5.2 (Formes canoniques). *Il est possible de déterminer la forme syntaxique d'une valeur étant donné son type.*

- si $\Gamma \vdash v : \text{INT}$, $v = n$.
- si $\Gamma \vdash v : (t_1, \dots, t_n) \rightarrow t$, $v = f$.
- si $\Gamma \vdash v : t^*$, $v = \varphi$.
- si $\Gamma \vdash v : t[]$, $v = [v_1, \dots, v_n]$.
- si $\Gamma \vdash v : (t_1, \dots, t_n) \rightarrow t$, $v = \text{fun}(a_1, \dots, a_n)((l_1, e_1), \dots, (l_p, e_p))$.

Lemme 5.3 (Permutation). *L'ordre dans lequel les variables apparaissent dans un environnement n'influe pas sur la relation de typage.*

Pour toute permutation σ de $[1; n]$, on note $\sigma(x_1 : t_1, \dots, x_n : t_n) = x_{\sigma(1)} : t_{\sigma(1)}, \dots, x_{\sigma(n)} : t_{\sigma(n)}$.

Alors : si $\Gamma \vdash e : t$ et $\Gamma' = \sigma(\Gamma)$, alors $\Gamma' \vdash e : t$.

Lemme 5.4 (Affaiblissement). *De même que l'ordre n'influe pas le typage, on peut aussi ajouter des associations supplémentaires dans l'environnement sans modifier les typages dans cet environnement.*

Si $\Gamma \vdash e : t$ et $x \notin \text{dom}(\Gamma)$, alors $\Gamma, x : t' \vdash e : t$.

Lemme 5.5 (Substitution). *Si dans une expression e il apparaît une variable x de type t' , le typage est préservé lorsqu'on remplace ses occurrences par une expression e' de même type.*

Si $\Gamma, x : t' \vdash e : t$ et $\Gamma \vdash e' : t'$, alors $\Gamma \vdash e[x/e'] : t$.

Ces lemmes permettent de prouver le théorème suivant :

Théorème 5.2 (Préservation). *Si une expression est typable et que son évaluation produit une valeur, alors cette valeur est du même type que l'expression.*

Si $\Gamma \vdash e : t$ et $e \rightarrow v$

alors $\Gamma \vdash v : t$.

TODO

- ordre des sections
- noter les flottants, fonctions et champs différemment
- versions $\Gamma \vdash i$ des propriétés

QUALIFICATEURS DE TYPE

Dans le chapitre 5, nous avons vu comment ajouter un système de types forts statiques à un langage impératif. Ici, nous étendons ce système afin de lui ajouter des *qualificateurs de type* qui décrivent l'origine des données. Ils permettent de restreindre certaines opérations sensibles à des expressions dont la valeur est sûre.

6.1 Provenance des pointeurs

6.1.1 Éditions et ajouts

La sémantique d'évaluation du langage n'est pas modifiée. En revanche, on modifie légèrement le système de types (figure 6.1) afin d'ajouter à chaque pointeur un *qualificateur* qui représente qui contrôle sa valeur.

Les deux qualificateurs possibles sont :

- **KERNEL** : il s'applique aux pointeurs contrôlés par le noyau. Par exemple, prendre l'adresse d'un objet donne un pointeur noyau.
- **USER** : il s'applique aux pointeurs qui proviennent de l'espace utilisateur. Ces pointeurs proviennent toujours d'interfaces particulières, comme les appels système ou les paramètres de la fonction `ioctl`.

Règle de sûreté du déréférencement

$$\frac{\Gamma \vdash e : \tau \text{ KERNEL}^*}{\Gamma \vdash *e : \tau} \text{ (LV-DEREF-KERNEL)}$$

Règle d'évaluation

Qualificateurs	$q ::= \text{KERNEL}$ USER	Donnée noyau (sûre) Donnée utilisateur (non sûre)
Types	$t ::= \dots$ t^* tq^*	Pointeur Pointeur qualifié
Instructions		$i ::= \dots$ $\text{TAIN}(x)$

FIGURE 6.1 – Changements liés aux qualificateurs de types

$$\frac{}{\langle \text{TAIN}(x), m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{ (TAIN)}$$

Règle de taintage

$$\frac{\Gamma \vdash x : t \text{ USER } *}{\Gamma \vdash \text{TAIN}(x)} \text{ (TAIN)}$$

6.1.2 Propriété d'isolation mémoire

Le déréférencement d'un pointeur dont la valeur est contrôlée par l'utilisateur ne peut se faire qu'à travers une fonction qui vérifie la sûreté de celui-ci.

6.2 Analyse de terminaison des chaînes C

Dans ce chapitre, nous présentons une autre extension au système de types du chapitre 5, similaire à celle du chapitre 6. Il s'agit cette fois-ci de détecter les pointeurs sur caractères ($\text{char } *$) qui sont terminés par un caractère NUL et donc une chaîne C correcte. La bibliothèque C propose quantité de fonctions manipulant ces chaînes et appeler une fonction comme `strcpy` sur un pointeur quelconque est un problème de sécurité que nous cherchons à détecter.

6.2.1 But

Le langage C ne fournit pas directement de type "chaîne de caractère". C'est au programmeur de les gérer via des pointeurs sur caractère ($\text{char } *$).

En théorie le programmeur est libre de choisir une représentation : des chaînes préfixées par la longueur, une structure contenant la taille et un pointeur vers les données, ou encore une chaîne avec un terminateur comme 0.

Néanmoins c'est ce dernier style qui est le plus idiomatique : par exemple, les littéraux de chaîne ("comme ceci") ajoutent un octet nul à la fin. De plus, le standard décrit dans la bibliothèque d'exécution de nombreuses fonctions destinées à les manipuler — c'est le fichier `<string.h>` ([ISO99] section 7.21).

Ainsi la fonction `strcpy` a pour prototype :

```
char *strcpy(char *dest, const char *src);
```

Elle réalise la copie de la chaîne pointée par `src` à l'endroit pointé par `dest`. Pour détecter la fin de la chaîne, cette fonction parcourt la mémoire jusqu'à trouver un caractère nul. Une implémentation naïve pourrait être :

```
char *strcpy(char *dest, const char *src)
{
    int i;
    for(i=0;src[i]!=0;i++) {
        dest[i] = src[i];
    }
    return dest;
}
```

La copie n'est arrêtée que lorsqu'un 0 est lu. Autrement dit, si quelqu'un contrôle la valeur pointée par `src`, il pourra écraser autant de données qu'il le désire. On est dans le cas d'école du débordement de tampon sur la pile tel que décrit dans [One96]. Considérons la fonction suivante :

```
void f(char *src)
{
    char buf[100];
    strcpy(buf, src);
}
```

Si le pointeur `src` pointe sur une chaîne de longueur supérieure à 100 (ou une zone mémoire qui n'est pas une chaîne et ne contient pas de 0), les valeurs placées sur la pile juste avant `buf` (à une adresse supérieure) seront écrasées. Avec les conventions d'appel habituelles, il s'agit de l'adresse de retour de la fonction. Un attaquant pourra donc détourner le flot d'exécution du programme.

Pour éviter ces cas de fonctions vulnérables, on peut introduire une distinction entre les pointeurs `char *` classiques (représentant l'adresse d'un caractère par exemple) et les pointeurs sur une chaîne terminée par un caractère nul.

Dans certaines bases de code (la plus célèbre étant celle de Microsoft), une convention syntaxique est utilisée : les pointeurs vers des chaînes terminées par 0 ont un nom qui commence par `sz`, comme `"szTitle"`. C'est pourquoi nous appellerons ce qualificateur de type `sz`.

6.2.2 Approche

Cette propriété est un peu différente de la séparation entre espace utilisateur et espace noyau modélisée dans le chapitre 6 : autant un pointeur reste contrôlé par l'utilisateur (ou sûr) toute sa vie, autant le fait d'être terminé par un octet nul dépend de l'ensemble de l'état mémoire. Il y a deux problèmes principaux à considérer.

D'une part, l'*aliasing* rend l'analyse difficile : si `p` et `q` pointent tous les deux vers une même zone mémoire, le fait de modifier l'un peut modifier l'autre. D'autre part, ce n'est pas parce qu'une fonction maintient l'invariant de terminaison, qu'elle le maintient à chaque instruction.

On peut résoudre en partie le problème d'*aliasing* en étant très conservateur, c'est à dire en sous-approximant l'ensemble des chaînes du programme (on traitera une chaîne légitime comme une chaîne non terminée, interdisant par excès de zèle les fonctions comme `strcpy`).

Le second problème est plus délicat puisqu'il casse l'hypothèse habituelle que chaque variable conserve le même type au long de sa vie. Plusieurs techniques sont possibles pour contourner ce problème : la première est d'être encore une fois conservateur et d'interdire ces constructions (on ne pourrait alors analyser que les programmes ne manipulant les chaînes qu'à travers les fonctions de la bibliothèque standard). Une autre est d'insérer des annotations permettant de s'affranchir localement du système de types. Enfin, il est possible d'utiliser un système de types où les variables ont en plus d'un type, un automate d'états possible dépendant de la position dans le programme : c'est le concept de *typestates*[SY86].

6.2.3 Annotation de `string.h`

Une première étape est d'annoter l'ensemble des fonctions manipulant les chaînes de caractères.

Fonctions de copie

`memcpy`

```
void *memcpy(void *dest, const void *src, size_t n);
```

memmove

```
void *memmove(void *dest, const void *src, size_t n);
```

strcpy

```
char *strcpy(char *dest, const char *src);
```

strncpy

```
char *strncpy(char *dest, const char *src, size_t n);
```

Fonctions de concaténation**strcat**

```
char *strcat(char *dest, const char *src);
```

strncat

```
char *strncat(char *dest, const char *src, size_t n);
```

Fonctions de comparaison**memcmp**

```
int memcmp(const void *s1, const void *s2, size_t n);
```

strcmp

```
int strcmp(const char *s1, const char *s2);
```

strncmp

```
int strncmp(const char *s1, const char *s2, size_t n);
```

strcoll

```
int strcoll(const char *s1, const char *s2);
```

strxfrm

```
size_t strxfrm(char *dest, const char *src, size_t n);
```

Fonctions de recherche**memchr**

```
void *memchr(const void *s, int c, size_t n);
```

strchr

```
char *strchr(const char *s, int c);
```

strcspn

```
size_t strcspn(const char *s, const char *reject);
```

strpbrk

```
char *strpbrk(const char *s, const char *accept);
```

strrchr

```
char *strrchr(const char *s, int c);
```

strspn

```
size_t strspn(const char *s, const char *accept);
```

strstr

```
char *strstr(const char *haystack, const char *needle);
```

strtok

```
char *strtok(char *str, const char *delim);
```

Fonctions diverses**memset**

```
void *memset(void *s, int c, size_t n);
```

strerror

```
char *strerror(int errnum);
```


strlen

```
size_t strlen(const char *s);
```

6.2.4 Typage des primitives

6.2.5 Extensions au système de types

6.2.6 Résultats

Troisième partie

Expérimentation

IMPLANTATION

Dans ce chapitre, nous décrivons la mise en œuvre des analyses statiques précédentes. Nous commençons par un tour d’horizon des représentations intermédiaires possibles, avant de décrire celle retenue : Newspeak. La chaîne de compilation est explicitée, partant de C pour aller au langage impératif décrit dans le chapitre 4. Enfin, nous donnons les détails d’un algorithme d’inférence de types à la Hindley-Milner, reposant sur l’unification et le partage de références.

7.1 Langages intermédiaires

Le langage C [KR88, ISO99] a été conçu pour être une sorte d’assembleur portable, permettant décrire du code indépendamment de l’architecture sur laquelle il sera compilé. Historiquement, c’est il a permis de créer Unix, et ainsi de nombreux logiciels bas niveau sont écrits en C. En particulier, il existe des compilateurs de C vers les différents langages machine pour à peu près toutes les architectures.

Lors de l’écriture d’un compilateur, on a besoin d’un langage intermédiaire qui

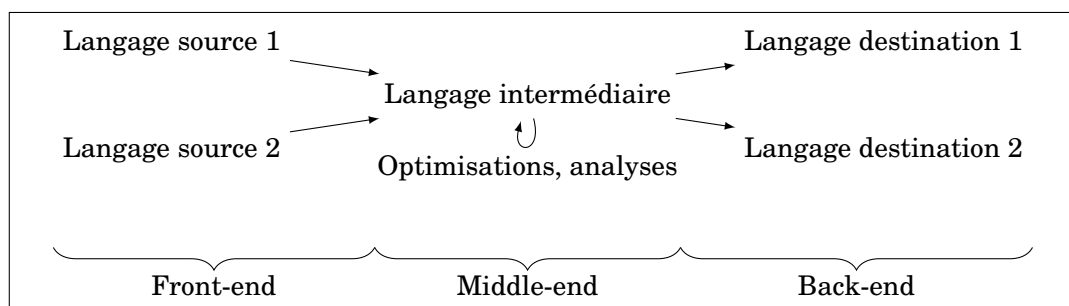


FIGURE 7.1 – Décomposition d’un compilateur : front-ends, middle-end, back-ends

fasse l'intermédiaire entre *front-end* et *back-end* (figure 7.1). Depuis ce langage on doit pouvoir exprimer des transformations intermédiaires sur cette représentation (analyses sémantiques, optimisations, etc), mais aussi compiler ce langage vers un langage machine.

L'idée de prendre C comme langage intermédiaire est très séduisante, mais malheureusement sa sémantique est trop complexe et trop peu spécifiée. Il est donc judicieux d'utiliser un langage plus simple à cet effet. Dans de nombreux projets, des sous-ensembles de C ont été définis pour aller dans ce sens.

Langages

Les premiers candidats sont bien entendu les représentations intermédiaires utilisées dans les compilateurs C. Elles ont l'avantage d'accepter en plus du C standard, les diverses extensions (GNU, Microsoft, Plan9) utilisées par la plupart des logiciels. En particulier, le noyau Linux repose fortement sur les extensions GNU.

GCC utilise une représentation interne nommée GIMPLE[Mer03]. Il s'agit d'une structure d'arbre écrite en C, reposant sur de nombreuses macros afin de cacher les détails d'implémentation pouvant varier entre deux versions de GCC. Cette représentation étant réputée difficile à manipuler, le projet MELT[Sta11] permet de générer une passe de compilation à partir d'un dialecte de Lisp.

LLVM [LA04] est un compilateur développé par la communauté puis sponsorisé Apple. À la différence de GCC, sa base de code est écrite en C++. Il utilise une représentation intermédiaire qui peut être manipulée soit sous forme d'une structure de données C++, soit d'un fichier de code-octet compact, soit sous forme textuelle.

Objective Caml [Cam1] utilise pour sa génération de code une représentation interne nommée Cmm, disponible dans les sources du compilateur sous le chemin `asmcomp/cmm.mli` (il s'agit donc d'une structure de données OCaml). Ce langage a l'avantage d'être très restreint, mais malheureusement il n'existe pas directement de traducteur permettant de compiler C vers Cmm.

C- [PJN097] [Cam7], dont le nom est inspiré du précédent, est un projet qui visait à unifier les langages intermédiaires utilisés par les compilateurs. L'idée est que si un front-end peut émettre du C- (sous forme de texte), il est possible d'obtenir du code machine efficace. Le compilateur Haskell GHC utilise une représentation intermédiaire très similaire à C-.

Comme le problème de construire une représentation intermédiaire adaptée à une analyse statique n'est pas nouveau, plusieurs projets ont déjà essayé d'y apporter une

solution. Puisque qu'ils sont développés en parallèle des compilateurs, le support des extensions est en général moins important dans ces langages.

CIL [NMRW02] [6] est une représentation en OCaml d'un programme C, développée depuis 2002. Grâce à un mécanisme de greffons, elle permet de prototyper rapidement des analyses statiques de programmes.

Newspeak [HL08] est un langage intermédiaire développé par EADS Innovation Works, et qui est spécialisé dans l'analyse de valeurs par interprétation abstraite. Il sera décrit plus en détails dans la section 7.2.

Compcert est un projet qui vise à produire un compilateur certifié pour C. C'est à dire que le fait que les transformations conservent la sémantique est prouvé. Il utilise de nombreux langages intermédiaires, dont CIL. Pour le front-end, le langage se nomme Clight[BDL06]. Les passes de middle-end, quant à elles, utilisent Cminor[AB07].

7.2 Newspeak

(wip)

7.3 Chaîne de compilation

La compilation vers C est faite en trois étapes (figure 7.2) : prétraitement du code source, compilation de C prétraité vers NEWSPEAK, puis compilation de NEWSPEAK vers ce langage.

7.3.1 Prétraitement

C2NEWSPEAK travaillant uniquement sur du code prétraité (dans directives de préprocesseur), la première étape consiste donc à faire passer le code par CPP : les macros sont développées, les constantes remplacées par leurs valeurs, les commentaires supprimés, les fichiers d'en-tête inclus, etc.

7.3.2 Compilation (levée des ambiguïtés)

Cette passe est réalisée par l'utilitaire C2NEWSPEAK. L'essentiel de la compilation consiste à mettre à plat les définition de types, et à simplifier le flot de contrôle. C en effet propose de nombreuses constructions ambiguës ou redondantes.

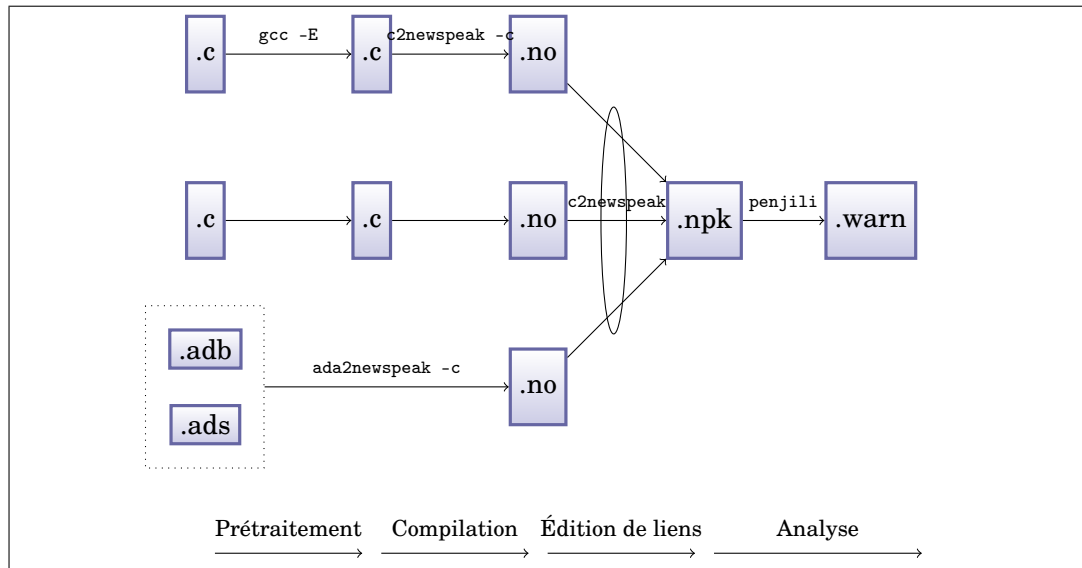


FIGURE 7.2 – Compilation depuis Newspeak

Au contraire, NEWSPEAK propose un nombre réduit de constructions. Rappelons que le but de ce langage est de faciliter l'analyse statique : des constructions orthogonales permettent donc d'éviter la duplication de règles sémantique, ou de code lors de l'implémentation d'un analyseur.

Par exemple, plutôt que de fournir une boucle *while*, une boucle *do/while* et une boucle *for*, NEWSPEAK fournit une unique boucle `WHILE(1){}`. La sortie de boucle est compilée vers un `GOTO`, qui est toujours un saut vers l'avant (similaire à un "break" généralisé).

La sémantique de NEWSPEAK et la traduction de C vers NEWSPEAK sont décrites dans [HL08]. En ce qui concerne l'élimination des sauts vers l'arrière, on peut se référer à [EH94].

7.3.3 Annotations

NEWSPEAK a de nombreux avantages, mais pour une analyse par typage il est trop bas niveau. Par exemple, dans le code suivant

```

struct s {
    int a;
    int b;
};

int main(void)

```



```

{
    struct s x;
    int y[10];
    x.b = 1;
    y[1] = 1;
    return 0;
}

```

(wip)

7.3.4 Implantation de l'algorithme de typage

Commençons par étudier le cas du lambda-calcul simplement typé (figure 7.3). Prenons l'exemple de la fonction suivante¹ :

$$f = \lambda x. \lambda y. \text{plus}(\text{plus}(\text{fst } x)(\text{snd } x))y$$

On voit que puisque `fst` et `snd` sont appliqués à `x`, ce doit être un tuple. En outre on additionne ces deux composantes ensemble, donc elles doivent être de type `INT` (et le résultat aussi). Par le même argument, `y` doit aussi être de type `INT`. En conclusion, `x` est de type `INT × INT` et `y` de type `INT`, donc `f` est de type `INT × INT → INT → INT`.

Mais comment faire pour implanter cette analyse ? En fait le système de types de la figure 7.3 a une propriété particulièrement intéressante : chaque forme syntaxique (variable, abstraction, etc) est en conclusion exactement d'une règle de typage. Cela permet de toujours savoir quelle règle il faut appliquer.

Partant du terme de conclusion (f), on peut donc en déduire un squelette d'arbre d'inférence (figure 7.4)²

Une fois à cette étape, on peut donner un nom à chaque type inconnu : τ_1, τ_2, \dots . L'utilisation qui en est faite permet de générer un ensemble de contraintes d'unification. Par exemple, pour chaque application de la règle (APP) :

$$\frac{\Gamma \vdash \dots : \tau_3 \quad \Gamma \vdash \dots : \tau_1}{\Gamma \vdash \dots : \tau_2} \text{ (APP)}$$

on doit déduire que $\tau_3 = \tau_1 \rightarrow \tau_2$.

-
1. On suppose que `plus` est une fonction de l'environnement global qui a pour type `INT → INT → INT`.
 2. Par souci de clarté, les prémisses des applications de (VAR) ne sont pas notées.

Termes	$t ::= x$	Variable	
	$\lambda x. t$	Abstraction	
	t	Application	
	n	Entier	
	f	Flottant	
	(t, t)	Couple	
	$\text{fst } t$	Projection gauche	
	$\text{snd } t$	Projection droite	
Types	$\tau ::= \text{INT}$	Entier	
	FLOAT	Flottant	
	$\tau \rightarrow \tau$	Fonction	
	$\tau \times \tau$	Produit	
Contextes	$\Gamma ::= \varepsilon$	Contexte vide	
	$\Gamma, x : \tau$	Extension	
Règles	$\frac{}{\Gamma \vdash n : \text{INT}} \text{ (INT)}$	$\frac{}{\Gamma \vdash f : \text{FLOAT}} \text{ (FLOAT)}$	$\frac{x : \tau \in \Gamma}{\Gamma \vdash x : \tau} \text{ (VAR)}$
	$\frac{\Gamma \vdash f : \tau_1 \rightarrow \tau_2 \quad \Gamma \vdash x : \tau_1}{\Gamma \vdash f x : \tau_2} \text{ (APP)}$	$\frac{\Gamma, x : \tau_1 \vdash y : \tau_2}{\Gamma \vdash \lambda x. y : \tau_1 \rightarrow \tau_2} \text{ (ABS)}$	
	$\frac{\Gamma \vdash x : \tau_1 \times \tau_2}{\Gamma \vdash \text{fst } x : \tau_1} \text{ (PROJ-G)}$	$\frac{\Gamma \vdash x : \tau_1 \times \tau_2}{\Gamma \vdash \text{snd } x : \tau_2} \text{ (PROJ-D)}$	
	$\frac{\Gamma \vdash x : \tau_1 \quad \Gamma \vdash y : \tau_2}{\Gamma \vdash (x, y) : \tau_1 \times \tau_2} \text{ (TUP)}$		

FIGURE 7.3 – Lambda calcul simplement typé avec entiers, flottants et couples

$$\begin{array}{c}
\frac{}{\Gamma_2 \vdash x} \text{ (VAR)} \\
\frac{}{\Gamma_2 \vdash \text{fst } x} \text{ (PROJ-G)} \\
\frac{}{\Gamma_2 \vdash \text{plus}} \text{ (VAR)} \quad \vdots \\
\frac{}{\Gamma_2 \vdash \text{plus}(\text{fst } x)} \text{ (APP)} \\
\vdots \quad \frac{}{\Gamma_2 \vdash x} \text{ (VAR)} \\
\vdots \quad \frac{}{\Gamma_2 \vdash \text{snd } x} \text{ (PROJ-D)} \\
\frac{}{\Gamma_2 \vdash \text{plus}(\text{fst } x)(\text{snd } x)} \text{ (APP)} \\
\vdots \quad \frac{}{\Gamma_2 \vdash \text{plus}} \text{ (VAR)} \\
\frac{}{\Gamma_2 \vdash \text{plus}(\text{plus}(\text{fst } x)(\text{snd } x))} \text{ (APP)} \\
\vdots \quad \frac{}{\Gamma_2 \vdash y} \text{ (VAR)} \\
\frac{}{\Gamma_2 \vdash \text{plus}(\text{plus}(\text{fst } x)(\text{snd } x))y} \text{ (APP)} \\
\frac{}{\Gamma_1 \vdash \lambda y. \text{plus}(\text{plus}(\text{fst } x)(\text{snd } x))y} \text{ (ABS)} \\
\frac{}{\Gamma_0 \vdash \lambda x. \lambda y. \text{plus}(\text{plus}(\text{fst } x)(\text{snd } x))y} \text{ (ABS)}
\end{array}$$

$\Gamma_1 = \Gamma_0, x \quad \Gamma_2 = \Gamma_0, x, y$

FIGURE 7.4 – Arbre d'inférence : règles à utiliser

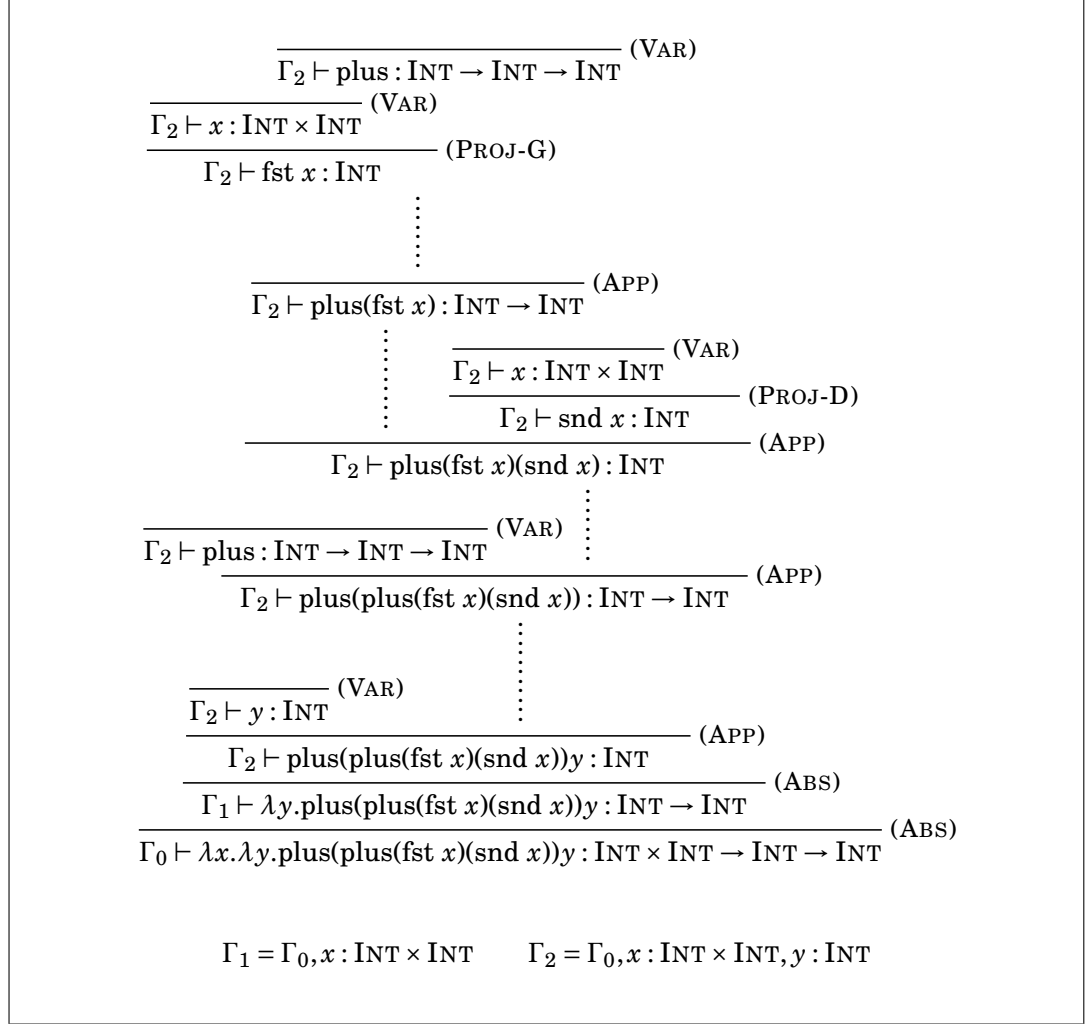


FIGURE 7.5 – Arbre d'inférence complet

Ce signe $=$ est à prendre comme une contrainte d'égalité : partant d'un ensemble de contraintes de la forme "type avec inconnue = type avec inconnue", on veut obtenir une substitution "inconnue \rightarrow type concret".

Pour résoudre ces contraintes, on commence par les simplifier : si $\tau_a \rightarrow \tau_b = \tau_c \rightarrow \tau_d$, alors $\tau_a = \tau_c$ et $\tau_b = \tau_d$. De même si $\tau_a \times \tau_b = \tau_c \times \tau_d$. Au contraire, si $\tau_a \rightarrow \tau_b = \tau_c \times \tau_d$, il est impossible d'unifier les types et il faut abandonner l'inférence de types. D'autres cas sont impossibles, par exemple $\text{INT} = \tau_1 \rightarrow \tau_2$ ou $\text{INT} = \text{FLOAT}$.

Une fois ces simplifications réalisées, les contraintes restantes sont d'une des formes suivantes :

- $\tau_i = \tau_i$. Il n'y a rien à faire, cette contrainte peut être supprimée.
- $\tau_i = \tau_j$ avec $i \neq j$: toutes les occurrences de τ_j dans les autres contraintes peuvent être remplacées par τ_i .
- $\tau_i = x$ (ou $x = \tau_i$) où x est un type concret : idem.

Une fois toutes les substitutions effectuées, on obtient un arbre de typage correct (figure 7.5, donc un programme totalement inféré.

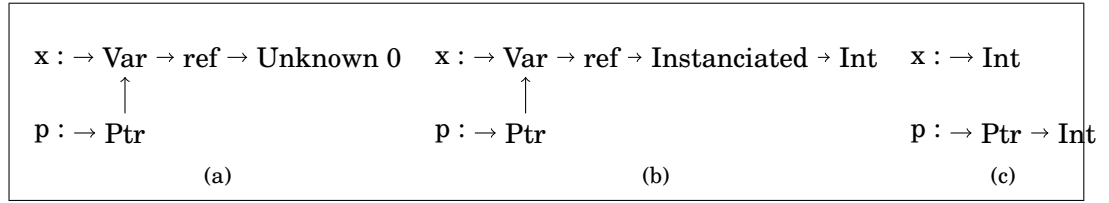


FIGURE 7.6 – Unification par partage

```
int x;
int *p = &x;
x = 0;
```

FIGURE 7.7 – Compilation d'un programme C - avant

Plutôt que de modifier toutes les occurrences d'un type τ_i , on va affecter à τ_i la valeur du nouveau type.

L'implémentation de cet algorithme utilise le partage et les références (figure 7.6).

D'abord 7.6a, ensuite 7.6b, et enfin 7.6c.

Prenons l'exemple de la figure 7.7 et typons-le "à la main". On commence par oublier toutes les étiquettes de type présentes dans le programme. Celui-ci devient alors :

```
var x, p;
p = &x;
x = 0;
```

La première ligne introduit deux variables. On peut noter leurs types respectifs (inconnus pour le moment) t_1 et t_2 . La première affectation $p = \&x$ implique que les deux côtés du signe "=" ont le même type. À gauche, le type est t_2 , et à droite $\text{Ptr}(t_1)$. On applique le même raisonnement à la seconde affectation : à gauche, le type est t_1 et à droite Int. On en déduit que le type de x est Int et celui de p est $\text{Ptr}(\text{Int})$.

```
type var_type =
| Unknown of int
| Instanciated of ml_type

and const_type =
| Int_type
| Float_type

and ml_type =
```

```

| Var_type of var_type ref
| Const_type of const_type
| Pair_type of ml_type * ml_type
| Fun_type of ml_type * ml_type

```

Pour implanter cet algorithme, on représente les types de données du programmes à typer par une valeur de type `ml_type`. En plus des constantes de types comme `int` ou `float`, et des constructeurs de type comme `pair` et `fun`, le constructeur `Var` permet d'exprimer les variables de types (inconnues ou non).

Celles-ci sont numérotées par un `int`, on suppose avoir à disposition deux fonctions manipulant un compteur global d'inconnues.

```

module Counter : sig
  val reset_unknowns : unit -> unit
  val new_unknown : unit -> int
end

```

De plus, on a un module gérant les environnements de typage. Il pourra être implanté avec des listes d'association ou des tables de hachage, par exemple. Sa signature est :

```

module Env : sig
  type t

  (* Construction *)
  val empty : t
  val extend : ml_ident -> ml_type -> t -> t

  (* Interrogation *)
  val get : ml_ident -> t -> ml_type option
end

```

Reprenons l'exemple précédent. Partant d'un environnement vide (`Env.empty`), on commence par l'étendre de deux variables. Comme on n'a aucune information, il fait allouer des nouveaux noms d'inconnues (qui correspondent à t_1 et t_2) :

```

let t1 = Var_type (Unknown (new_unknown ())) in
let t2 = Var_type (Unknown (new_unknown ())) in
let env =
  Env.extend "p" t2
  (Env.extend "x" t1
   Env.empty
  ) in

```

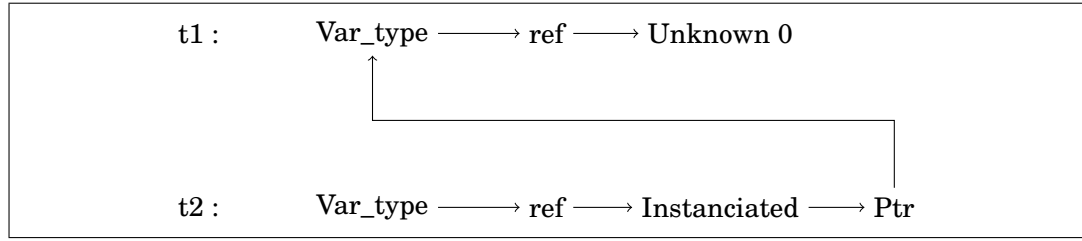


FIGURE 7.8 – Unification : partage

La première instruction indique que les deux côtés de l'affectation doivent avoir le même type.

```

let lhs1 = Lv_var "p"
and rhs1 = AddrOf (Exp_var "x") in
let t_lhs1 = typeof lhs1 env
and t_rhs1 = typeof rhs1 env in
unify t_lhs1 t_rhs1;

```

Ici il se passe plusieurs choses intéressantes. D'une part nous faisons appel à une fonction externe `typeof` qui retourne le type d'une expression sous un environnement (dans une implantation complète il s'agirait d'un appel récursif). Dans ce cas, `typeof lhs1 env` est identique à `Env.get lhs1 env` et `typeof rhs1 env` à `Ptr_type t1`. L'autre aspect intéressant est la dernière ligne : la fonction `unify` va modifier en place les représentations des types afin de les rendre égales. L'implantation de `unify` sera décrite plus tard. Dans ce cas précis, elle va faire pointer la référence dans `t2` vers `t1` (figure 7.8).

Enfin, la seconde affectation se déroule à peu près de la même manière.

```

let lhs2 = Lv_deref (Lv_var "p")
and rhs2 = Exp_int 0 in
let t_lhs2 = typeof lhs2 env
and t_rhs2 = typeof rhs2 env in
unify t_lhs2 t_rhs2;

```

Ici `typeof lhs2 env` est identique à `Ptr_type (Env.get "p" env)` et `typeof rhs2 env` à `Const_type Int_type`. Et dans ce cas, l'unification doit se faire entre `t1` et `Const_type Int_type` : cela mute la référence derrière `t1` (figure 7.9).

L'essence de l'algorithme d'inférence se situe donc dans 2 fonctions. D'une part, `unify` qui réalise l'unification des types grâce à au partage des références. D'autre part, la `typeof` qui encode les règles de typage elles-mêmes et les applique à l'aide de `unify`.

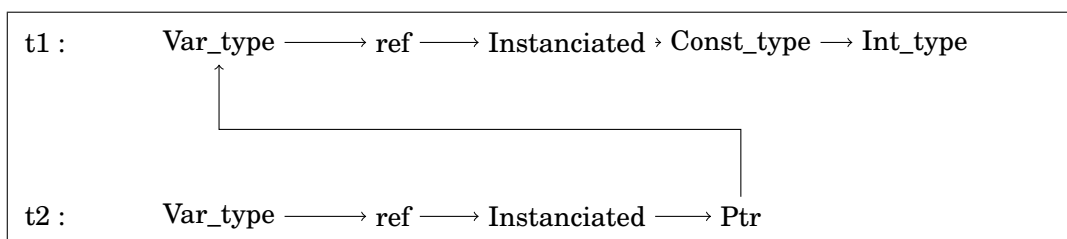


FIGURE 7.9 – Unification par mutation de références

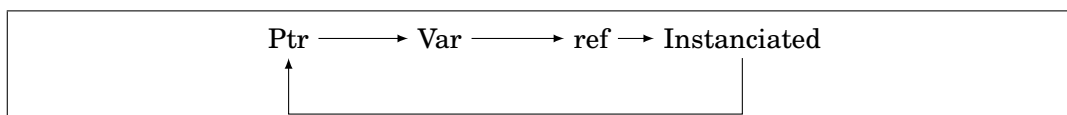


FIGURE 7.10 – Cycle dans le graphe de types

7.3.5 Algorithme d'unification

Voici une implantation de la fonction `unify`.

Celle-ci prend en entrée deux types t_1 et t_2 . À l'issue de l'exécution de `unify`, ces deux types doivent pouvoir être considérés comme égaux. Si ce n'est pas possible, une erreur sera levée.

La première étape est de réduire ces deux types, c'est à dire à transformer les constructions `Var (ref (Instancié t))` en `t`.

Ensuite, cela dépend des formes qu'ont les types réduits :

- si les deux types sont inconnus (de la forme `Var (ref (Instanciated t))`), on fait pointer l'une des deux références vers le premier type. Notons que cela crée un type de la forme `Var (ref (Instanciated (Var (ref (Unknown n)))))` qui sera réduit lors d'une prochaine étape d'unification.
- si un type est inconnu et pas l'autre, il faut de la même manière affecter la référence. Mais en faisant ça inconditionnellement, cela peut poser problème : par exemple en tentant d'unifier `a` avec `Ptr(a)` on pourrait créer un cycle dans le graphe (figure 7.10). Pour éviter cette situation, il suffit de s'assurer que le type inconnu n'est pas présent dans le type à affecter.
- si les deux types sont des types de base (comme `INT` ou `FLOAT`) égaux, on ne fait rien.
- si les deux types sont des constructeurs de type, il faut que les constructeurs soient égaux. On unifie en outre leurs arguments deux à deux.
- dans les autres cas, l'algorithme échoue.
- TODO sous typage pour les structures

- implem du polymorphisme
- implem du sous-typage
- généralisation depuis le toy language

```

Decl
( "x"
, Newspeak.Scalar (Newspeak.Int (Newspeak.Signed, 32))
, ()
, [ Decl
    ( "p"
    , Newspeak.Scalar Newspeak.Ptr
    , ()
    , [ Set
        ( Local "p"
        , ( AddrOf (Local "x")
          , ()
          )
        , Newspeak.Scalar Newspeak.Ptr
        )
      ; Set
        ( Local "x"
        , ( Const (CInt Nat.zero)
          , ()
          )
        , Newspeak.Scalar (Newspeak.Int (Newspeak.Signed, 32))
        )
      ]
    )
  ]
)

```

FIGURE 7.11 – Compilation d'un programme C - après

Le programme C (figure 7.7) est compilé ainsi en Tyspeak (figure 7.11).

ÉTUDE DE CAS : UN PILOTE DE CARTE GRAPHIQUE

8.1 Description du problème

(wip)

Un système d'exploitation moderne comme GNU/Linux est séparé en deux niveaux de privilèges : le noyau, qui gère directement le matériel, et les applications de l'utilisateur, qui communiquent avec le noyau par l'interface restreinte des *appels système*.

Pour assurer l'isolation, ces deux parties n'ont pas accès aux mêmes zones mémoire (cf. figure 2.5).

Si le code utilisateur tente d'accéder à la mémoire du noyau, une erreur sera déclenchée. En revanche, si cette écriture est faite au sein de l'implantation d'un appel système, il n'y aura pas d'erreur puisque le noyau a accès à toute la mémoire : l'isolation aura donc été brisée.

Pour celui qui implante un appel système, il faut donc empêcher qu'un pointeur passé en paramètre référence le noyau. Autrement dit, il est indispensable de vérifier dynamiquement que la zone dans laquelle pointe le paramètre est accessible par l'appelant[Har88].

Si au contraire un tel pointeur est déréférencé sans vérification (avec * ou une fonction comme memcpy), le code s'exécutera correctement mais en rendant le système vulnérable, comme le montre la figure 8.1.

Pour éviter cela, le noyau fournit un ensemble de fonctions qui permettent de vérifier dynamiquement la valeur d'un pointeur avant de le déréférencer. Par exemple, dans la figure précédente, la ligne 8 aurait dû être remplacée par :

```
copy_from_user(&value, value_ptr, sizeof(value));
```

cf annexe A

FIGURE 8.1 – Bug freedesktop.org #29340. Le paramètre `data` provient de l'espace utilisateur via un appel système. Un appelant malveillant peut se servir de cette fonction pour lire la mémoire du noyau à travers le message d'erreur.

L'analyse présentée ici permet de vérifier automatiquement et statiquement que les pointeurs qui proviennent de l'espace utilisateur ne sont déréférencés qu'à travers une de ces fonctions sûres.

8.2 Principes de l'analyse

Le problème est modélisé de la façon suivante : on associe à chaque variable `x` un type de données `t`, ce que l'on note `x:t`. En plus des types présents dans le langage C, on ajoute une distinction supplémentaire pour les pointeurs. D'une part, les pointeurs "noyau" (de type `t *`) sont créés en prenant l'adresse d'un objet présent dans le code source. D'autre part, les pointeurs "utilisateurs" (leur type est noté `t user*`) proviennent des interfaces avec l'espace utilisateur.

Il est sûr de déréférencer un pointeur noyau, mais pas un pointeur utilisateur. L'opérateur `*` prend donc un `t *` en entrée et produit un `t`.

Pour faire la vérification de type sur le code du programme, on a besoin de quelques règles. Tout d'abord, les types suivent le flot de données. C'est-à-dire que si on trouve dans le code `a = b`, `a` et `b` doivent avoir un type compatible. Ensuite, le qualificateur `user` est récursif : si on a un pointeur utilisateur sur une structure, tous les champs pointeurs de la structure sont également utilisateur. Enfin, le déréférencement s'applique aux pointeurs noyau seulement : si le code contient l'expression `*x`, alors il existe un type `t` tel que `x:t*` et `*x:t`.

Appliquons ces règles à l'exemple de la figure 8.1 : on suppose que l'interface avec l'espace utilisateur a été correctement annotée. Cela permet de déduire que `data:void user*`. En appliquant la première règle à la ligne 6, on en déduit que `info:struct drm_radeon_info user*` (comme en C, on peut toujours convertir de et vers un pointeur sur `void`).

Pour déduire le type de `value_ptr` dans la ligne 7, c'est la deuxième règle qu'il faut appliquer : le champ `value` de la structure est de type `uint32_t *` mais on y accède à travers un pointeur utilisateur, donc `value_ptr:uint32_t user*`.

À la ligne 8, on peut appliquer la troisième règle : à cause du déréférencement, on en déduit que `value_ptr:t *`, ce qui est une contradiction puisque d'après les lignes précédentes, `value_ptr:uint32_t user*`.

Si la ligne 3 était remplacée par l'appel à `copy_from_user`, il n'y aurait pas d'erreur de typage car cette fonction peut accepter les arguments (`uint32_t *`, `uint32_t user*`, `size_t`).

8.3 Implantation

Une implantation est en cours. Le code source est d'abord prétraité par `gcc -E` puis converti en Newspeak [HL08], un langage destiné à l'analyse statique. Ce traducteur peut prendre en entrée tout le langage C, y compris de nombreuses extensions GNU utilisées dans le noyau. En particulier, l'exemple de la figure 8.1 peut être analysé.

À partir de cette représentation du programme et d'un ensemble d'annotations globales, on propage les types dans les sous-expressions jusqu'aux feuilles.

Si aucune contradiction n'est trouvée, c'est que le code respecte la propriété d'isolation. Sinon, cela peut signifier que le code n'est pas correct, ou bien que le système de types n'est pas assez expressif pour le code en question.

Le prototype, disponible sur [♣⁸], fera l'objet d'une démonstration.

8.4 Conclusion

Nous avons montré que le problème de la manipulation de pointeurs non sûrs peut être traité avec une technique de typage. Elle est proche des analyses menées dans CQual [FFA99] ou Sparse [♣⁵].

Plusieurs limitations sont inhérentes à cette approche : notamment, la présence d'unions ou de *casts* entre entiers et pointeurs fait échouer l'analyse.

Le principe de cette technique (associer des types aux valeurs puis restreindre les opérations sur certains types) peut être repris. Par exemple, si on définit un type "numéro de bloc" comme étant un nouvel alias de `int`, on peut considérer que multiplier deux telles valeurs est une erreur.

CHAPITRE 9

CONCLUSION

9.1 Limitations

9.2 Perspectives



CODE DU MODULE NOYAU

```
/* from drivers/gpu/drm/radeon/radeon_kms.c */
int radeon_info_ioctl(struct drm_device *dev, void *data, struct drm_file *filp)
{
    struct radeon_device *rdev = dev->dev_private;
    struct drm_radeon_info *info;
    struct radeon_mode_info *minfo = &rdev->mode_info;
    uint32_t *value_ptr;
    uint32_t value;
    struct drm_crtc *crtc;
    int i, found;

    info = data;
    value_ptr = (uint32_t *)((unsigned long)info->value);
    value = *value_ptr;
    switch (info->request) {
    case RADEON_INFO_DEVICE_ID:
        value = dev->pci_device;
        break;
    case RADEON_INFO_NUM_GB_PIPES:
        value = rdev->num_gb_pipes;
        break;
    case RADEON_INFO_NUM_Z_PIPES:
        value = rdev->num_z_pipes;
        break;
    case RADEON_INFO_ACCEL_WORKING:
        /* xf86-video-ati 6.13.0 relies on this being false for evergreen */

```

```

    if ((rdev->family >= CHIP_CEDAR) && (rdev->family <= CHIP_HEMLOCK))
        value = false;
    else
        value = rdev->accel_working;
    break;
case RADEON_INFO_CRTC_FROM_ID:
    for (i = 0, found = 0; i < rdev->num_crtc; i++) {
        crtc = (struct drm_crtc *)minfo->crtcs[i];
        if (crtc && crtc->base.id == value) {
            struct radeon_crtc *radeon_crtc = to_radeon_crtc(crtc);
            value = radeon_crtc->crtc_id;
            found = 1;
            break;
        }
    }
    if (!found) {
        DRM_DEBUG_KMS("unknown crtc id %d\n", value);
        return -EINVAL;
    }
    break;
case RADEON_INFO_ACCEL_WORKING2:
    value = rdev->accel_working;
    break;
case RADEON_INFO_TILING_CONFIG:
    if (rdev->family >= CHIP_CEDAR)
        value = rdev->config.evergreen.tile_config;
    else if (rdev->family >= CHIP_RV770)
        value = rdev->config.rv770.tile_config;
    else if (rdev->family >= CHIP_R600)
        value = rdev->config.r600.tile_config;
    else {
        DRM_DEBUG_KMS("tiling config is r6xx+ only!\n");
        return -EINVAL;
    }
case RADEON_INFO_WANT_HYPERZ:
    mutex_lock(&dev->struct_mutex);
    if (rdev->hyperz_filp)
        value = 0;
    else {
        rdev->hyperz_filp = filp;
        value = 1;
    }

```

```

    }
    mutex_unlock(&dev->struct_mutex);
    break;
default:
    DRM_DEBUG_KMS("Invalid request %d\n", info->request);
    return -EINVAL;
}
if (DRM_COPY_TO_USER(value_ptr, &value, sizeof(uint32_t))) {
    DRM_ERROR("copy_to_user\n");
    return -EFAULT;
}
return 0;
}

/* from drivers/gpu/drm/radeon/radeon_kms.c */
struct drm_ioctl_desc radeon_ioctls_kms[] = {
    /* KMS */
    DRM_IOCTL_DEF(DRM_RADEON_INFO, radeon_info_ioctl, DRM_AUTH|DRM_UNLOCKED)
};

/* from drivers/gpu/drm/radeon/radeon_drv.c */

static struct drm_driver kms_driver = {
    .driver_features =
        DRIVER_USE_AGP | DRIVER_USE_MTRR | DRIVER_PCI_DMA | DRIVER_SG |
        DRIVER_HAVE_IRQ | DRIVER_HAVE_DMA | DRIVER_IRQ_SHARED | DRIVER_GEM,
    .dev_priv_size = 0,
    .ioctls = radeon_ioctls_kms,
    .name = "radeon",
    .desc = "ATI Radeon",
    .date = "20080528",
    .major = 2,
    .minor = 6,
    .patchlevel = 0,
};

/* from drivers/gpu/drm/drm_drv.c */
int drm_init(struct drm_driver *driver)
{
    DRM_DEBUG("\n");
    INIT_LIST_HEAD(&driver->device_list);

```

```
if (driver->driver_features & DRIVER_USE_PLATFORM_DEVICE)
    return drm_platform_init(driver);
else
    return drm_pci_init(driver);
}
```



RÈGLES D'ÉVALUATION

	$\frac{\langle e, m \rangle \rightarrow \langle e', m' \rangle}{\langle C[e], m \rangle \rightarrow \langle C[e'], m' \rangle} \text{ (CTX)}$
	$C ::= \bullet$
	$ C \boxplus e$
	$ v \boxplus C$
	$ \boxminus C$
	$ * C$
	$ \varphi.C$
	$ \varphi[C]$
	$ C[e]$
	$ C \leftarrow e$
	$ \varphi \leftarrow C$
	$ \{f_1 : v_1; \dots; f_i : C; \dots; f_n : e_n\}$
	$ [v_1; \dots; C; \dots; e_n]$
	$ C(e_1, \dots, e_n)$
	$ f(v_1, \dots, C, \dots, e_n)$
	$ C; i$
	$ \text{IF}(C)\{i_1\}\text{ELSE}\{i_2\}$
	$ \text{RETURN}(C)$
Contextes	

FIGURE B.1 – Règles d'évaluation - contextes

$\frac{a = \text{Lookup}(x, m)}{\langle x, m \rangle \rightarrow \langle a, m \rangle} \text{ (PHI-VAR)}$	$\frac{}{\langle * \varphi, m \rangle \rightarrow \langle \widehat{*} \varphi, m \rangle} \text{ (PHI-DEREF)}$
$\frac{}{\langle lv.f, m \rangle \rightarrow \langle lv.\widehat{f}, m \rangle} \text{ (PHI-STRUCT)}$	$\frac{}{\langle \varphi[n], m \rangle \rightarrow \langle \varphi[\widehat{n}], m \rangle} \text{ (PHI-ARRAY)}$
$\frac{}{\langle c, m \rangle \rightarrow \langle \widehat{c}, m \rangle} \text{ (EXP-CST)}$	$\frac{}{\langle f, m \rangle \rightarrow \langle \widehat{f}, m \rangle} \text{ (EXP-FUN)}$
$\frac{}{\langle \varphi, m \rangle \rightarrow \langle m[\varphi], m \rangle} \text{ (EXP-LV)}$	

FIGURE B.2 – Règles d'évaluation - constantes et left-values

$\frac{}{\langle \boxminus v, m \rangle \rightarrow \langle \widehat{\boxminus} v, m \rangle} \text{ (EXP-UNOP)}$	$\frac{\boxplus \notin \{/, \%\}}{\langle v_1 \boxplus v_2, m \rangle \rightarrow \langle v_1 \widehat{\boxplus} v_2, m \rangle} \text{ (EXP-BINOP)}$
$\frac{\boxdiv \in \{/, \%\} \quad v_2 \neq \widehat{0}}{\langle v_1 \boxdiv v_2, m \rangle \rightarrow \langle v_1 \widehat{\boxdiv} v_2, m \rangle} \text{ (EXP-DIV)}$	$\frac{\boxdiv \in \{/, \%\}}{\langle v_1 \boxdiv 0, m \rangle \rightarrow \Omega_{div}} \text{ (EXP-DIV-ZERO)}$
$\frac{a = \text{Lookup}(x, m)}{\langle \&x, m \rangle \rightarrow \langle a, m \rangle} \text{ (EXP-ADDR OF)}$	$\frac{}{\langle \varphi \leftarrow v, m \rangle \rightarrow \langle v, m[\varphi \leftarrow v] \rangle} \text{ (EXP-SET)}$

FIGURE B.3 – Règles d'évaluation - opérations

$$\begin{array}{c}
\frac{}{\langle \{f_1 : v_1; \dots; f_n : v_n\}, m \rangle \rightarrow \langle \{f_1 : v_1; \dots; f_n : v_n\}, m \rangle} \text{(EXP-STRUCT)} \\
\\
\frac{}{\langle [v_1, \dots, v_n], m \rangle \rightarrow \langle [v_1, \dots, v_n], m \rangle} \text{(EXP-ARRAY)} \\
\\
\begin{array}{l}
f = \text{fun}(a_1, \dots, a_n)((l'_1, e'_1), \dots, (l'_p, e'_p))\{i\} \\
m_1 = \text{Push}(m_0, ((a_1, v_1), \dots, (a_n, v_n))) \\
\langle \left(\begin{array}{c} e'_1 \\ \vdots \\ e'_p \end{array} \right), m_1 \rangle \rightarrow \langle \left(\begin{array}{c} v'_1 \\ \vdots \\ v'_p \end{array} \right), m_2 \rangle \quad m_3 = \text{Extend}(m_2, ((l_1, v_1), \dots, (l_n, v_n))) \\
\langle i, m_3 \rangle \rightarrow \langle \text{RETURN}(v), m_4 \rangle \quad m_5 = \text{Pop}(m_4) \quad m_6 = \text{Cleanup}(m_5) \\
\hline
\langle f(v_1, \dots, v_n), m_0 \rangle \rightarrow \langle v, m_6 \rangle
\end{array} \text{(EXP-CALL)}
\end{array}$$

FIGURE B.4 – Règles d'évaluation - expressions composées

$$\begin{array}{c}
\frac{\langle i, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle (i; i'), m \rangle \rightarrow \langle i', m' \rangle} \text{(SEQ)} \quad \frac{}{\langle (\text{PASS}; i), m \rangle \rightarrow \langle i, m \rangle} \text{(PASS)} \\
\\
\frac{}{\langle v, m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{(EXP)} \quad \frac{}{\langle \text{IF}(0)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_f, m \rangle} \text{(IF-FALSE)} \\
\\
\frac{v \neq 0}{\langle \text{IF}(v)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_t, m \rangle} \text{(IF-TRUE)} \\
\\
\frac{}{\langle \text{WHILE}(e)\{i\}, m \rangle \rightarrow \langle \text{IF}(e)\{i; \text{WHILE}(e)\{i\}\}, m \rangle} \text{(WHILE)} \\
\\
\frac{}{\langle \text{RETURN}(v); i, m \rangle \rightarrow \langle \text{RETURN}(e), m \rangle} \text{(RETURN)} \quad \frac{}{m \vdash \text{struct } s\{\dots\} \rightarrow m} \text{(PH-STRUCT)} \\
\\
\frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{m \vdash e \rightarrow m'} \text{(PH-EXP)} \quad \frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{(s, g) \vdash x = e \rightarrow (s, (x, v) :: g)} \text{(PH-VAR)} \\
\\
\frac{}{m \vdash [] \rightarrow^* m} \text{(PH*-NIL)} \quad \frac{m \vdash p \rightarrow m' \quad m' \vdash ps \rightarrow^* m''}{m \vdash p :: ps \rightarrow^* m''} \text{(PH*-CONS)}
\end{array}$$

FIGURE B.5 – Règles d'évaluation

RÈGLES DE TYPAGE

$\frac{}{\Gamma \vdash i : \text{INT}} \text{ (CST-INT)}$	$\frac{}{\Gamma \vdash f : \text{FLOAT}} \text{ (CST-FLOAT)}$	$\frac{}{\Gamma \vdash \text{NULL} : t*} \text{ (CST-NULL)}$
$\frac{x : t \in \Gamma}{\Gamma \vdash x : t} \text{ (LV-VAR)}$	$\frac{\Gamma \vdash lv : t*}{\Gamma \vdash *lv : t} \text{ (LV-DEREF)}$	$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash lv : t[]}{\Gamma \vdash lv[e] : t} \text{ (LV-INDEX)}$
$\frac{(s, f, t_f) \in S \quad \Gamma \vdash lv : \text{struct } s}{\Gamma \vdash lv.f : t_f} \text{ (LV-FIELD)}$		

FIGURE C.1 – Règles de typage - constantes et variables

$\frac{\boxplus \in \{+, -, \times, /, \&, , ^, \&\&, , \ll, \gg\} \quad \Gamma \vdash e_1 : \text{INT} \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-INT)}$		
$\frac{\boxplus \in \{+., -., \times., /.\} \quad \Gamma \vdash e_1 : \text{FLOAT} \quad \Gamma \vdash e_2 : \text{FLOAT}}{\Gamma \vdash e_1 \boxplus e_2 : \text{FLOAT}} \text{ (OP-FLOAT)}$		
$\frac{t \in \{\text{INT}, \text{FLOAT}\}}{\text{EQ}(t)} \text{ (EQ-NUM)}$	$\frac{\text{EQ}(t)}{\text{EQ}(t*)} \text{ (EQ-PTR)}$	$\frac{\text{EQ}(t)}{\text{EQ}(t[]) } \text{ (EQ-ARRAY)}$
$\frac{\boxplus \in \{=, \neq\} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \quad \text{EQ}(t)}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-EQ)}$		
$\frac{\boxplus \in \{=, \neq, \leq, \geq, <, >\} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \quad t \in \{\text{INT}, \text{FLOAT}\}}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}} \text{ (OP-COMPARABLE)}$		
$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash +e : \text{INT}} \text{ (UNOP-PLUS-INT)}$		$\frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash +.e : \text{FLOAT}} \text{ (UNOP-PLUS-FLOAT)}$
$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash -e : \text{INT}} \text{ (UNOP-MINUS-INT)}$		$\frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash -.e : \text{FLOAT}} \text{ (UNOP-MINUS-FLOAT)}$
$\frac{\boxminus \in \{\sim, !\} \quad \Gamma \vdash e : \text{INT}}{\Gamma \vdash \boxminus e : \text{INT}} \text{ (UNOP-NOT)}$		
$\frac{\boxtimes \in \{+_p, -_p\} \quad \Gamma \vdash e_1 : t* \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxtimes e_2 : t*} \text{ (PTR-ARITH)}$		

FIGURE C.2 – Règles de typage - opérateurs

$\frac{\Gamma \vdash lv : t}{\Gamma \vdash \&lv : t*} \text{ (ADDR)}$	$\frac{\Gamma \vdash lv : t \quad \Gamma \vdash e : t}{\Gamma \vdash lv \leftarrow e : t} \text{ (SET)}$	$\frac{\forall i \in [1;n], \Gamma \vdash e_i : t}{\Gamma \vdash \{e_1; \dots; e_n\} : t[]} \text{ (ARRAY)}$
$\frac{\forall i \in [1;n], \Gamma \vdash e_i : t_i \quad \forall i \in [1;n], (s, f_i, t_i) \in S}{\Gamma \vdash \{f_1 : e_1; \dots; f_n : e_n\} : \text{struct } s} \text{ (STRUCT)}$		
$\frac{\Gamma \vdash e : (t_1, \dots, t_n) \rightarrow t \quad \forall i \in [1;n], \Gamma \vdash e_i : t_i}{\Gamma \vdash e(e_1, \dots, e_n) : t} \text{ (CALL)}$	$\frac{}{\Gamma \vdash \text{PASS}} \text{ (PASS)}$	
$\frac{\Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash i_1; i_2} \text{ (SEQ)}$	$\frac{\Gamma \vdash e : t}{\Gamma \vdash e} \text{ (EXP)}$	$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash \text{IF}(e)\{i_1\}\text{ELSE}\{i_2\}} \text{ (IF)}$
$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i}{\Gamma \vdash \text{WHILE}(e)\{i\}} \text{ (WHILE)}$	$\frac{\Gamma \vdash \underline{R} \leftarrow e}{\Gamma \vdash \text{RETURN}(e)} \text{ (RETURN)}$	
$\frac{\Gamma' = (\Gamma - \underline{R}), a_1 : t_1, \dots, a_n : t_n, l_1 : t'_1, \dots, l_n : t'_n, \underline{R} : t \quad \forall i \in [1;p], \Gamma \vdash e_i : t'_i \quad \Gamma' \vdash i}{\Gamma \vdash \text{fun}(a_1, \dots, a_n)((l_1, e_1), \dots, (l_p, e_p))\{i\} : (t_1, \dots, t_n) \rightarrow t} \text{ (FUN)}$		$\frac{\Gamma \vdash e : t}{\Gamma \vdash e \rightarrow \Gamma} \text{ (PH-EXP)}$
$\frac{\Gamma \vdash e : t \quad \Gamma' = (x, t), \Gamma}{\Gamma \vdash p \rightarrow \Gamma'} \text{ (PH-VAR)}$	$\frac{S' = (x_1, t_1, s), \dots, (x_n, t_n, s)}{\Gamma \vdash \text{struct } s\{x_1 : t_1; \dots; x_n : t_n\} \rightarrow \Gamma} \text{ (PH-STRUCT)}$	

FIGURE C.3 – Règles de typage

TABLE DES FIGURES

2.1	Cadres de pile	9
2.2	Les différents <i>rings</i>	11
2.3	Implantation de la mémoire virtuelle	12
2.4	Mécanisme de mémoire virtuelle.	12
2.5	Espace d'adressage d'un processus	13
2.6	Appel de <code>gettimeofday</code>	16
2.7	Zones mémoire	16
2.8	Implantation de l'appel système <code>gettimeofday</code>	17
4.1	Syntaxe - expressions	24
4.2	Syntaxe - instructions	25
4.3	Syntaxe - opérateurs	25
4.4	Composantes d'un état mémoire	27
4.5	Opérations de pile	28
4.6	Contextes d'exécution	32
4.7	Substitution dans les contextes d'évaluation	33
4.8	Évaluation des left-values.	34
4.9	L'appel d'une fonction. La taille de la pile croît de gauche à droite, et les réductions se font de haut en bas.	37
4.10	Nettoyage d'un cadre de pile	38
5.1	Programmes bien et mal formés	45
5.2	Types et environnements de typage	46
6.1	Changements liés aux qualificateurs de types	60
7.1	Décomposition d'un compilateur : front-ends, middle-end, back-ends	69
7.2	Compilation depuis Newspeak	72
7.3	Lambda calcul simplement typé avec entiers, flottants et couples	74
7.4	Arbre d'inférence : règles à utiliser	75
7.5	Arbre d'inférence complet	76
7.6	Unification par partage	78
7.7	Compilation d'un programme C - avant	78
7.8	Unification : partage	80
7.9	Unification par mutation de références	81

7.10	Cycle dans le graphe de types	81
7.11	Compilation d'un programme C - après	83
8.1	Bug freedesktop.org #29340	86
B.1	Règles d'évaluation - contextes	96
B.2	Règles d'évaluation - constantes et left-values	97
B.3	Règles d'évaluation - opérations	97
B.4	Règles d'évaluation - expressions composées	98
B.5	Règles d'évaluation	98
C.1	Règles de typage - constantes et variables	99
C.2	Règles de typage - opérateurs	100
C.3	Règles de typage	101

RÉFÉRENCES WEB

- [🌐¹] The Objective Caml system, documentation and user's manual – release 3.12
<http://caml.inria.fr/pub/docs/manual-ocaml/>
- [🌐²] Haskell Programming Language – Official Website
<http://www.haskell.org/>
- [🌐³] Python Programming Language – Official Website
<http://www.python.org/>
- [🌐⁴] Perl Programming Language – Official Website
<http://www.perl.org/>
- [🌐⁵] Sparse - a Semantic Parser for C
https://sparse.wiki.kernel.org/index.php/Main_Page
- [🌐⁶] CIL - C Intermediate Language
<http://kerneis.github.com/cil/>
- [🌐⁷] The C - - language
<http://www.cminusminus.org/>
- [🌐⁸] Penjili project
<http://www.penjili.org/>

BIBLIOGRAPHIE

- [AB07] Andrew W. Appel and Sandrine Blazy. Separation logic for small-step Cminor (extended version). Research report 6138, INRIA, 2007. 29 pages. 71
- [BBC⁺10] Al Bessey, Ken Block, Ben Chelf, Andy Chou, Bryan Fulton, Seth Hallem, Charles Henri-Gros, Asya Kamsky, Scott McPeak, and Dawson Engler. A few billion lines of code later : using static analysis to find bugs in the real world. *Commun. ACM*, 53(2) :66–75, February 2010. 20
- [BC05] Daniel P. Bovet and Marco Cesati. *Understanding the Linux Kernel, Third Edition*. O'Reilly Media, third edition edition, November 2005. 12
- [BDH⁺09] Julien Brunel, Damien Doligez, René Rydhof Hansen, Julia L. Lawall, and Gilles Muller. A foundation for flow-based program matching using temporal logic and model checking. In *The 36th Annual ACM SIGPLAN - SIGACT Symposium on Principles of Programming Languages*, pages 114–126, Savannah, GA, USA, January 2009. 19
- [BDL06] Sandrine Blazy, Zaynah Dargaye, and Xavier Leroy. Formal verification of a C compiler front-end. In *FM 2006 : Int. Symp. on Formal Methods*, volume 4085 of *Lecture Notes in Computer Science*, pages 460–475. Springer, 2006. 71
- [CC77] Patrick Cousot and Radhia Cousot. Abstract interpretation : a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *POPL '77 : Proceedings of the 4th ACM SIGACT-SIGPLAN symposium on Principles of Programming Languages*, pages 238–252, New York, NY, USA, 1977. ACM. 19
- [CC92] P. Cousot and R. Cousot. Abstract interpretation and application to logic programs. *Journal of Logic Programming*, 13(2–3) :103–179, 1992. (The editor of Journal of Logic Programming has mistakenly published the unreadable galley proof. For a correct version of this paper, see <http://www.di.ens.fr/~cousot>.) 19
- [CCF⁺05] Patrick Cousot, Radhia Cousot, Jérôme Feret, Laurent Mauborgne, Antoine Miné, David Monniaux, and Xavier Rival. The ASTREE analyzer. In Shmuel Sagiv, editor, *ESOP*, volume 3444 of *Lecture Notes in Computer Science*, pages 21–30. Springer, 2005. 20

- [CCF⁺09] Patrick Cousot, Radhia Cousot, Jérôme Feret, Laurent Mauborgne, Antoine Miné, and Xavier Rival. Why does astrée scale up? *Formal Methods in System Design*, 35(3) :229–264, 2009. 20
- [CMP10] Dumitru Ceară, Laurent Mounier, and Marie-Laure Potet. Taint dependency sequences : A characterization of insecure execution paths based on input-sensitive cause sequences. In *ICST Workshops*, 2010. 20
- [DRS00] Nurit Dor, Michael Rodeh, and Mooly Sagiv. CSSV : Towards a realistic tool for statically detecting all buffer overflows in C, 2000. 20
- [EH94] Ana Erosa and Laurie J. Hendren. Taming control flow : A structured approach to eliminating goto statements. In *In Proceedings of 1994 IEEE International Conference on Computer Languages*, pages 229–240. IEEE Computer Society Press, 1994. 72
- [FFA99] Jeffrey S. Foster, Manuel Fähndrich, and Alexander Aiken. A theory of type qualifiers. In *Programming language design and implementation, PLDI '99*, pages 192–203, 1999. 19, 87
- [FJKA06] Jeffrey S. Foster, Robert Johnson, John Kodumal, and Alex Aiken. Flow-insensitive type qualifiers. *ACM Trans. Program. Lang. Syst.*, 28 :1035–1087, November 2006. 19
- [FTA02] Jeffrey S. Foster, Tachio Terauchi, and Alex Aiken. Flow-sensitive type qualifiers. In *PLDI '02 : Proceedings of the ACM SIGPLAN 2002 Conference on Programming language design and implementation*, volume 37, pages 1–12, New York, NY, USA, May 2002. ACM Press. 19
- [Gor04] Mel Gorman. *Understanding the Linux Virtual Memory Manager*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 2004. 12
- [Gra92] Philippe Granger. Improving the results of static analyses programs by local decreasing iteration. In *Proceedings of the 12th Conference on Foundations of Software Technology and Theoretical Computer Science*, pages 68–79, London, UK, UK, 1992. Springer-Verlag. 20
- [Har88] Norm Hardy. The confused deputy (or why capabilities might have been invented). *ACM Operating Systems Review*, 22(4) :36–38, October 1988. 17, 85
- [HL08] Charles Hymans and Olivier Levillain. Newspeak, Doubleplussimple Minilang for Goodthinkful Static Analysis of C. Technical Note 2008-IW-SE-00010-1, EADS IW/SE, 2008. 71, 72, 87
- [Int] Intel Corporation. *Intel[®] 64 and IA-32 Architectures Software Developer's Manual*. 7, 13

- [ISO99] ISO. The ANSI C standard (C99). Technical Report WG14 N1124, ISO/IEC, 1999. 61, 69
- [JW04] Robert Johnson and David Wagner. Finding user/kernel pointer bugs with type inference. In *USENIX Security Symposium*, pages 119–134, 2004. 19
- [KcS07] Oleg Kiselyov and Chung chieh Shan. Lightweight static capabilities. *Electr. Notes Theor. Comput. Sci.*, 174(7) :79–104, 2007. 20
- [KR88] Brian W. Kernighan and Dennis M. Ritchie. *The C Programming Language Second Edition*. Prentice-Hall, Inc., 1988. 69
- [LA04] Chris Lattner and Vikram Adve. LLVM : A Compilation Framework for Lifelong Program Analysis & Transformation. In *Proceedings of the 2004 International Symposium on Code Generation and Optimization (CGO'04)*, Palo Alto, California, Mar 2004. 70
- [LZ06] Peng Li and Steve Zdancewic. Encoding information flow in Haskell. In *Proceedings of the 19th IEEE Workshop on Computer Security Foundations (CSFW '06)*, Washington, DC, USA, 2006. IEEE Computer Society. 20
- [Mau04] Laurent Mauborgne. ASTRÉE : Verification of absence of run-time error. In René Jacquart, editor, *Building the information Society (18th IFIP World Computer Congress)*, pages 384–392. The International Federation for Information Processing, Kluwer Academic Publishers, Aug 2004. 20
- [Mer03] J. Merrill. GENERIC and GIMPLE : a new tree representation for entire functions. In *GCC developers summit 2003*, pages 171–180, 2003. 70
- [NMRW02] George C. Necula, Scott McPeak, Shree Prakash Rahul, and Westley Weimer. Cil : Intermediate language and tools for analysis and transformation of c programs. In *Proceedings of the 11th International Conference on Compiler Construction, CC '02*, pages 213–228, London, UK, UK, 2002. Springer-Verlag. 71
- [oEE08] Institute of Electrical and Electronics Engineers. IEEE Standard for Floating-Point Arithmetic. Technical report, Microprocessor Standards Committee of the IEEE Computer Society, 3 Park Avenue, New York, NY 10016-5997, USA, August 2008. 26
- [One96] Aleph One. Smashing the stack for fun and profit. *Phrack*, 1996. 7, 61
- [Pie02] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002. 20

- [PJNO97] Simon L. Peyton Jones, Thomas Nordin, and Dino Oliva. C- : A portable assembly language. In Chris Clack, Kevin Hammond, and Antony J. T. Davie, editors, *Implementation of Functional Languages*, volume 1467 of *Lecture Notes in Computer Science*, pages 1–19. Springer, 1997. 70
- [PTS⁺11] Nicolas Palix, Gaël Thomas, Suman Saha, Christophe Calvès, Julia Lawall, and Gilles Muller. Faults in Linux : Ten years later. In *Sixteenth International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS 2011)*, Newport Beach, CA, USA, March 2011. 19
- [SAB10] Edward J. Schwartz, Thanassis Avgerinos, and David Brumley. All you ever wanted to know about dynamic taint analysis and forward symbolic execution (but might have been afraid to ask). In *Proceedings of the IEEE Symposium on Security and Privacy*, 2010. 20
- [SM03] Andrei Sabelfeld and Andrew C. Myers. Language-based information-flow security. *IEEE Journal on Selected Areas in Communications*, 21 :2003, 2003. 20
- [Spe05] Brad Spengler. grsecurity 2.1.0 and kernel vulnerabilities. *Linux Weekly News*, 2005. 19
- [Sta11] Basile Starynkevitch. Melt - a translated domain specific language embedded in the gcc compiler. In Olivier Danvy and Chung chieh Shan, editors, *DSL*, volume 66 of *EPTCS*, pages 118–142, 2011. 70
- [STFW01] Umesh Shankar, Kunal Talwar, Jeffrey S. Foster, and David Wagner. Detecting format string vulnerabilities with type qualifiers. In *SSYM'01 : Proceedings of the 10th conference on USENIX Security Symposium*, page 16, Berkeley, CA, USA, 2001. USENIX Association. 19
- [SY86] R E Strom and S Yemini. Typestate : A programming language concept for enhancing software reliability. *IEEE Trans. Softw. Eng.*, 12(1) :157–171, January 1986. 62
- [Tan07] Andrew S. Tanenbaum. *Modern Operating Systems*. Prentice Hall Press, Upper Saddle River, NJ, USA, 3rd edition, 2007. 5
- [VB04] Arnaud Venet and Guillaume Brat. Precise and efficient static array bound checking for large embedded c programs. In *Proceedings of the ACM SIGPLAN 2004 conference on Programming language design and implementation, PLDI '04*, pages 231–242, New York, NY, USA, 2004. ACM. 20