

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

Version du 2 juillet 2013
Soutenue le xx-yy-2013 devant un jury composé de

Président

Gilles Muller Université Pierre et Marie Curie

Rapporteurs

Sandrine Blazy IRISA

Pierre Jouvelot Mines ParisTech

Directeurs

Emmanuel Chailloux Université Pierre et Marie Curie

Sarah Zennou EADS Innovation Works

Examineurs

Olivier Levillain ANSSI

Vincent Simonet Google

TABLE DES MATIÈRES

Table des matières	i
1 Introduction	1
1.1 Rôle d'un système d'exploitation	2
1.2 Séparation entre noyau et espace utilisateur	3
1.3 Systèmes de types	4
1.4 Langages d'implantation	5
1.5 Langages d'analyse	6
1.6 Plan de la thèse	8
 I Méthodes formelles pour la sécurité	 9
2 Systèmes d'exploitation	13
2.1 Architecture physique	13
2.2 Tâches et niveaux de privilèges	14
2.3 Appels système	15
2.4 Le <i>Confused Deputy Problem</i>	16
 3 État de l'art	 19
3.1 Taxonomie	19
3.2 Méthodes syntaxiques	20
3.3 Interprétation abstraite	20
3.4 Typage	22
3.5 Analyse de code système	23
3.6 Langages sûrs	24
3.7 Logique de Hoare	24
3.8 Assistants de preuve	25
3.9 Conclusion	26
 Conclusion de la partie I	 29
 II Un langage pour l'analyse de code système : SAFESPEAK	 31
4 Syntaxe et sémantique	35
4.1 Notations	35
4.2 Fonctionnalités	40
4.3 Syntaxe	41
4.4 Mémoire et valeurs	41
4.5 Opérations sur les valeurs	45
4.6 Opérations sur les états mémoire	46
4.7 Accesseurs	48

4.8	Contextes d'évaluation	51
4.9	Left-values	52
4.10	Expressions	53
4.11	Instructions	55
4.12	Erreurs	58
4.13	Phrases	58
4.14	Exécution	58
5	Typage	61
5.1	Principe	61
5.2	Environnements et notations	62
5.3	Expressions	63
5.4	Instructions	66
5.5	Fonctions	67
5.6	Phrases	67
5.7	Sûreté du typage	68
5.8	Typage des valeurs	68
5.9	Progrès et préservation	71
6	Analyse de provenance de pointeurs	77
6.1	Extensions noyau pour SAFESPEAK	78
6.2	Extensions sémantiques	79
6.3	Insuffisance des types simples	80
6.4	Extensions du système de types	81
6.5	Sûreté du typage	83
	Conclusion de la partie II	87
	III Expérimentation	89
7	Implantation	93
7.1	Newspeak	93
7.2	Chaîne de compilation	93
7.3	Architecture de ptrtype	98
7.4	Inférence de types	99
7.5	Unification	100
7.6	Exemple	102
8	Étude de cas : un pilote de carte graphique	107
8.1	Linux	107
8.2	GNU C	107
8.3	Configuration	108
8.4	Appels systèmes sous Linux	108
8.5	Bug	110
8.6	Détails	111

IV Conclusion	113
9 Conclusion	115
9.1 Contributions	115
9.2 Limitations et travaux futurs	115
A Module Radeon KMS	123
B Règles d'évaluation	127
C Règles de typage	133
D Preuves	137
D.1 Composition de lentilles	137
D.2 Progrès	138
D.3 Préservation	142
D.4 Progrès pour les types qualifiés	145
D.5 Préservation pour les types qualifiés	146
E Analyse de bitmasks	147
E.1 Ajouts	147
Liste des figures	149
Liste des définitions	151
Liste des théorèmes et lemmes	151
Références web	153
Bibliographie	155

INTRODUCTION

Communication, audiovisuel, transports, médecine : tout ces domaines se sont transformés dans les dernières décennies, en particulier grâce à la révolution numérique. En effet le plus petit appareil électrique contient maintenant des composants matériels programmables.

En 2013, on pense bien sûr aux téléphones portables dont la fonctionnalité et la complexité les rapprochent des ordinateurs de bureau. Par exemple, le système d'exploitation Android de Google est basé sur le noyau Linux, destiné à la base aux micro-ordinateurs.

Le noyau d'un système d'exploitation est chargé de faire l'intermédiaire entre le matériel (processeur, mémoire, périphériques, ...) et les applications exécutées sur celui-ci (par exemple un navigateur web, une calculatrice ou un carnet d'adresses). Il doit aussi garantir la sécurité de celles-ci.

En tant qu'intermédiaire de confiance, le noyau a un certain nombre de responsabilités et est le seul à avoir accès à certaines informations sensibles. Il est capital de s'assurer qu'il est bien le seul à pouvoir y accéder. En particulier, il faut pouvoir vérifier que les requêtes faites par l'utilisateur au noyau ne peuvent pas volontairement ou involontairement détourner ce dernier et lui faire fuir des informations confidentielles.

Le problème est que comme tous les logiciels, les noyaux de système d'exploitation sont écrits par des humains, qui ne sont pas parfaits. Les activités de relecture et de débogage ont beau prendre la majeure partie du temps de développement, il est facile de laisser passer des défauts de programmation.

Ces erreurs, ou *bugs*, peuvent avoir des conséquences dramatiques sur le plan matériel ou humain. Dans le cadre des produits d'EADS, il est donc capital de vérifier la correction des logiciels embarqués. De plus, le coût de déployer des correctifs est assez important : il est d'autant plus important d'être correct dès la première version.

Une technique efficace est de réaliser des tests, c'est-à-dire exécuter le programme sous un environnement contrôlé. On peut alors détecter des comportements non désirés. Mais même avec des tests les plus exhaustifs il n'est pas possible de couvrir tous les cas d'utilisation.

Une autre approche est d'analyser le code source du programme avant de l'exécuter et de refuser de lancer les programmes qui contiennent des erreurs, quitte à limiter la possibilité d'en écrire certains. C'est l'analyse statique de programmes.

Une des techniques d'analyse statique les plus répandues est le typage statique, qui consiste à associer à chaque morceau de programme, une étiquette décrivant quel genre de valeur sera produite par son évaluation. Par exemple, si n est le nom d'une variable entière, alors $n + 2$ produira toujours une valeur entière.

Pour garantir l'isolation d'un noyau de système d'exploitation, un des points cruciaux est de restreindre la manière dont sont traitées les informations provenant des programmes utilisateur. Le but de cet thèse est de montrer que le typage statique peut être utilisé pour détecter et interdire ces manipulations dangereuses.

1.1 Rôle d'un système d'exploitation

Un ordinateur est constitué de nombreux composants matériels : microprocesseur, mémoire, et divers périphériques. Et 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 aspects.

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 ; et chaque programme aura à sa disposition une partie de la mémoire principale, ou du disque dur.

Si plusieurs programmes 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 n'utilisent pas la carte réseau en même temps. Ce sont des rôles du système d'exploitation.

Ainsi, au lieu d'accéder directement au matériel via des instructions de bas niveau, les programmes communiquent avec le noyau, qui centralise donc les appels au matériel, et abstrait certaines opérations.

Par exemple, comparons ce qui se passe concrètement lors de la copie de données depuis un cédérom ou une clef USB.

- Dans le cas du cédérom, il faut 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.
- Avec une clef, il faut 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 la même intention (copier de la mémoire depuis un périphérique amovible), ne sont pas effectuées en extension 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 est responsable de la gestion de la mémoire, des périphériques et des processus. Les détails d'implantation ne sont pas présentés à l'utilisateur ; à la place il manipule des abstractions, comme la notion de fichier. Pour une explication détaillée du concept de système d'exploitation ainsi que des cas d'étude, on pourra se référer à [Tan07].

Mode du processeur	Privilège (code)	Privilège (données)	Accès possible
Noyau	Noyau	Noyau	☑
Noyau	Noyau	Utilisateur	☑
Noyau	Utilisateur	Noyau	☑
Noyau	Utilisateur	Utilisateur	☑
Utilisateur	Noyau	Noyau	☐
Utilisateur	Noyau	Utilisateur	☐
Utilisateur	Utilisateur	Noyau	☐
Utilisateur	Utilisateur	Utilisateur	☑

FIGURE 1.1: Règles d'exécution de code et d'accès à la mémoire.

1.2 Séparation entre noyau et espace utilisateur

Puisque le noyau est garant du bon fonctionnement du système, il ne doit pas pouvoir être manipulé directement par l'utilisateur ou les programmes exécutés. Ainsi, il est nécessaire de mettre en place des protections entre les espaces noyau et utilisateur.

Au niveau matériel, on utilise la notion de *niveaux de privilèges* pour déterminer s'il est possible d'exécuter une instruction.

D'une part, le processeur contient un niveau de privilège intrinsèque. D'autre part, chaque zone mémoire contenant du code ou des données possède également un niveau de privilège minimum nécessaire. L'exécution d'une instruction est alors possible si et seulement si le niveau de privilège du processeur est supérieur à celui de l'instruction et des opérandes mémoires qui y sont présentes¹.

Par exemple, supposons qu'un programme utilisateur contienne l'instruction `mov %eax, 0x54c3b2f7`. Celle-ci consiste à déplacer le contenu du registre EAX vers l'adresse mémoire 0x54c3b2f7. Puisque cette adresse est dans l'espace utilisateur (la raison est décrite dans le chapitre 2), aucune erreur de protection mémoire n'est déclenchée.

Ainsi, pour une instruction manipulant des données en mémoire, les accès possibles sont décrits dans la figure 1.1. En cas d'impossibilité (signalée par un ☐), une erreur se produit et l'exécution s'arrête.

En plus de cette vérification, certains types d'instructions sont explicitement réservés au mode le plus privilégié : par exemple les lectures ou écritures sur des ports matériels, ou celles qui permettent de définir les niveaux de privilèges des différentes zones mémoire.

Les programmes utilisateur ne pouvant pas accéder à ces instructions de bas niveau, ils sont très limités dans ce qu'ils peuvent faire. Puisque l'interaction avec le matériel (comme l'écriture sur un disque dur) se fait uniquement via des instructions privilégiées, ils sont limités à l'utilisation du processeur et de la mémoire, permettant uniquement de réaliser des calculs.

Pour utiliser le matériel ou accéder à des abstractions de haut niveau (comme créer un nouveau processus), ils doivent donc passer par l'intermédiaire du noyau. La communication entre le noyau et les programmes utilisateur est constituée par le mécanisme des *appels système*.

Lors d'un appel système, une fonction du noyau est invoquée (en mode noyau) avec des paramètres provenant de l'utilisateur. Il faut donc être particulièrement précautionneux dans

1. Ici "supérieur" est synonyme de "plus privilégié". Dans l'implantation d'Intel présentée dans le chapitre 2, les niveaux sont numérotés de 0 à 3 où le niveau 0 est le plus privilégié.

le traitement de ces données.

Par exemple, considérons un appel système de lecture depuis un disque : on passe au noyau les arguments (d, o, n, a) où d est le nom du disque, o (pour *offset*) l'adresse sur le disque où commencer la lecture, n le nombre d'octets à lire et a l'adresse en mémoire où commencer à stocker les résultats.

Dans le cas d'utilisation prévu, le noyau va copier la mémoire lue dans a . Le processeur est en mode noyau, en train d'exécuter une instruction du noyau manipulant des données utilisateur. D'après la figure 1.1 aucune erreur ne se produit.

Mais même si ce cas ne produit pas d'erreur à l'exécution, il est tout de même erroné. En effet, si on passe à l'appel système une adresse a faisant partie de l'espace noyau, que se passe-t-il ?

L'exécution est presque identique : au moment de la copie on est en mode noyau, en train d'exécuter une instruction du noyau manipulant des données noyau. Encore une fois il n'y a pas d'erreur à l'exécution.

On peut donc écrire n'importe où en mémoire. De même, une fonction d'écriture sur un disque (et lisant en mémoire) permettrait de lire de la mémoire du noyau. À partir de ces primitives, on peut accéder aux autres processus exécutés, ou détourner l'exécution vers du code arbitraire. L'isolation est totalement brisée à cause de ces appels système.

La cause de ceci est qu'on a accédé à la mémoire en testant les privilèges du noyau au lieu de tester les privilèges de celui qui a fait la requête (l'utilisateur). Ce problème est connu sous le nom de *confused deputy problem* [Har88].

La bonne manière d'implanter un appel système est donc d'interdire le déréférencement direct des pointeurs dont la valeur peut être contrôlée par l'utilisateur. Dans le cas du passage par adresse d'un argument, il aurait fallu vérifier à l'exécution que celui-ci a bien les mêmes privilèges que l'appelant.

Il est facile d'oublier d'ajouter cette vérification, puisque le cas "normal" fonctionne. Avec ce genre d'exemple on voit comment les bugs peuvent arriver si fréquemment et pourquoi il est aussi capital de les détecter avant l'exécution.

1.3 Systèmes de types

La plupart des langages de programmation incorporent la notion de type, dont un des buts est d'empêcher de manipuler des données incompatibles entre elles.

En mémoire, les seules données qu'un ordinateur manipule sont des nombres. Selon les opérations effectuées, ils seront interprétés comme des entiers, des adresses mémoires, ou des caractères. Pourtant il est clair que certaines opérations n'ont pas de sens : par exemple, multiplier un nombre par une adresse, ou déréférencer le résultat d'une division sont des comportements qu'on voudrait pouvoir empêcher.

En un mot, le but du typage est de classer les objets et de restreindre les opérations possibles selon la classe d'un objet : en somme, "ne pas ajouter des pommes et des oranges". Le modèle qui permet cette classification est appelé *système de types* et est en général constitué d'un ensemble de *règles de typage*, comme "un entier plus un entier égale un entier".

Typage dynamique Dans ce cas, chaque valeur manipulée par le programme est décorée d'une étiquette définissant comment interpréter la valeur en question. Les règles de typage sont alors réalisées à l'exécution. Par exemple, l'opérateur "+" vérifie que ces deux opérandes ont une étiquette "entier", et construit alors une valeur obtenue en faisant l'addition des deux valeurs, avec une étiquette "entier". Par exemple, le langage Python [P+3] utilise cette stratégie.

Typage statique Dans ce cas on fait les vérifications à la compilation. Pour vérifier ceci, on donne à chaque fonction un contrat comme “si deux entiers sont passés, et que la fonction renvoie une valeur, alors cette valeur sera un entier”. Cet ensemble de contrats peut être vérifié statiquement par le compilateur, à l’aide d’un système de types statique.

Par exemple, on peut dire que l’opérateur “+” a pour type $(\text{INT}, \text{INT}) \rightarrow \text{INT}$. Cela veut dire que si on lui passe deux entiers (INT, INT) , alors la valeur obtenue est également un entier.

Typage fort ou faible Contrairement à la distinction claire entre typage statique ou dynamique, la séparation entre typage fort et faible est moins nette. En poussant à l’extrême, les systèmes de types forts permettent d’éliminer totalement la nécessité de réaliser des tests de typage. Mais souvent ce n’est pas le cas, car il peut y avoir des constructions au sein du langage qui permettent de contourner le système de types, comme un opérateur de transtypage.

Polymorphisme Parfois, il est trop restrictif de donner un unique type à une fonction. Si on considère une fonction ajoutant un élément à une liste, ou une autre triant un tableau en place, leur type doit-il faire intervenir le type des éléments manipulés ?

En première approximation, on peut imaginer fournir une version du code par type de données à manipuler. C’est la solution retenue par les premières versions du langage Pascal, ce qui rendait très difficile l’écriture de bibliothèques [Ker81]. On parle alors de monomorphisme.

Une autre manière de procéder est d’autoriser plusieurs fonctions à avoir le même nom, mais avec des types d’arguments différents. Par exemple, on peut définir séparément l’addition entre deux entiers, entre deux flottants, ou entre un entier et un flottant. Selon les informations connues à la compilation, la bonne version sera choisie. C’est ainsi que fonctionnent les opérateurs en C++. On parle de polymorphisme *ad hoc*, ou de surcharge.

La dernière possibilité est le polymorphisme paramétrique, qui consiste à utiliser le même code quelque soit le type des arguments. Dans ce cas, on utilise une seule fonction pour traiter une liste d’entiers ou une liste de flottants, par exemple. Au lieu d’associer à chaque fonction un type, dans certains cas on lui associe un type paramétré, instanciable en un type concret. Dans le cas des fonctions de traitement de liste, l’idée est que lorsqu’on ne touche pas aux éléments, alors le traitement est valable quelque soit leur type. Cette technique a été décrite en premier dans [Mil78].

Pour un tour d’horizon de différents systèmes types de statiques, avec en particulier du polymorphisme, on pourra se référer à [Pie02].

1.4 Langages d’implantation

Puisque notre but est d’analyser du code provenant de systèmes d’exploitation, il est nécessaire de s’intéresser aux langages de programmation dans lesquels ils sont écrits.

Assembleur Le noyau d’un système d’exploitation nécessite d’accéder au matériel, donc il est naturellement bas niveau. Pour accéder aux fonctionnalités spécifiques de chaque processeur, il est donc nécessaire d’en implanter une partie dans le langage d’assemblage natif de chaque architecture. Historiquement, les premiers systèmes d’exploitations étaient entièrement écrits en assembleur. Cela est de plus en plus rare, sauf dans les cas où les ressources sont trop limitées pour exécuter du code compilé, comme dans les systèmes embarqués ou temps-réel.

C Le système Unix, développé à partir de 1969, a tout d’abord été développé en assembleur sur un mini-ordinateur PDP-7, puis a été porté sur d’autres architectures matérielles. Pour aider ce portage, il a été nécessaire de créer un “assembleur portable”, le langage C [KR88, ISO99]. Son but est de fournir des abstractions au dessus du langage d’assemblage. Les structures de contrôle (`if`, `while`, `for`) permettent d’utiliser la programmation structurée, c’est-à-dire en limitant l’utilisation de l’instruction `goto`. Les types de données sont également abstraits de la machine : ainsi, `int` désigne un entier machine, indépendamment de sa taille concrète. Son système de types est assez rudimentaire : toutes les formes de transtypage sont acceptées, certaines conversions sont insérées automatiquement par le compilateur, et la plupart des abstractions fournies par le langage sont perméables. Le noyau Linux est écrit dans un dialecte du langage C. Le noyau du système Mac OSX d’Apple est également un dérivé d’Unix, et est donc aussi écrit dans ce langage.

C++ Le manque de modularité de C a conduit à la création d’un nouveau langage plus riche autour de celui-ci. Dans les années 1980, C++ est apparu comme un C amélioré, avec de nombreuses fonctionnalités supplémentaires comme un modèle objet avec des classes, du polymorphisme *ad hoc* et un système de *templates* puissant permettant d’implanter une forme de polymorphisme universel. Microsoft Windows, initialement écrit en C, s’est enrichi de plus en plus de code C++, en faisant un des principaux systèmes d’exploitation écrits dans ce langage.

Ada Le langage C est bas niveau, mais son système de types n’est pas suffisamment évolué pour le rendre plus sûr que le langage d’assemblage vers lequel il est compilé. Dans les années 1970, le département de la défense américain (*DoD*) a créé un comité dans le but de créer un nouveau langage pour ses systèmes temps-réels embarqués. Le résultat est Ada, un langage impératif avec un système de types statique plus contraignant mais plus sûr que celui de C. Par exemple, deux types entiers différents (comme les entiers de 16 bits et ceux de 32 bits) ne sont pas compatibles. Il faut ajouter une conversion explicite dans le code, qui se traduira par un test à l’exécution et éventuellement par une erreur `Constraint_Error`. Bien qu’il ne se soit pas beaucoup développé dans les applications grand public, Ada garde toujours une place de choix dans l’implantation de systèmes embarqués, notamment dans les applications spatiales et de défense.

Autres langages Il est bien sûr possible d’implanter un système d’exploitation dans n’importe quel langage, si celui-ci peut être compilé vers du langage machine et qu’il permet facilement d’inclure des fragments d’assembleur pour les parties de bas niveau. On peut par exemple trouver des systèmes d’exploitation écrits en C# [HLA⁺05], Haskell [HJLT05] ou Ocaml [MMR⁺13].

1.5 Langages d’analyse

Les langages décrits précédemment sont faits pour être facilement écrits par des programmeurs humains. En général ils sont ambigus ou ont des comportements implicites. Pour analyser du code source, il est plus pratique d’avoir une représentation intermédiaire plus simple afin d’avoir moins de traitements dupliqués. Dans de nombreux projets, des sous-ensembles de C ont été définis pour aller dans ce sens.

Middle-ends

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'implantation interne pouvant varier entre deux versions. Cette représentation étant réputée difficile à manipuler, le projet MELT[Sta11] permet de générer un greffon de compilateur 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.

Cmm est une représentation interne utilisée pour la génération de code lors de la compilation d'Objective Caml [♣¹], et 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- - [PJNO97] [♣⁷], 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, par exemple, utilise une représentation intermédiaire très similaire à C- -.

Langages intermédiaires ad hoc

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.

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].

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 7.1.

La solution retenue : SAFESPEAK

Dans la suite de cette thèse nous nous appuierons sur SAFESPEAK, un langage intermédiaire minimal, pour exprimer les règles de typage de la manière la plus simple possible.

1.6 Plan de la thèse

Cette thèse comporte trois parties.

La partie I présente le contexte de ces travaux. Le fonctionnement général d'un système d'exploitation y est détaillé, et les problèmes de manipulation de pointeurs contrôlés par l'utilisateur y sont introduits. On fait ensuite un tour d'horizon des techniques existantes permettant de traiter ce problème. Cette partie se conclut par la proposition au cœur de cette thèse.

La partie II décrit notre solution : SAFESPEAK, un langage impératif. On y décrit sa syntaxe, sa sémantique ainsi qu'un système de types statiques. On l'étend ensuite pour capturer les problèmes d'adressage mémoire présents dans les systèmes d'exploitation en ajoutant des pointeurs contrôlés par l'utilisateur. Le système de types est également étendu. Pour chacune de ces variantes, on établit la propriété de sûreté de typage reliant la sémantique d'exécution aux types statiques.

La partie III documente la démarche expérimentale associée à ces travaux. L'implantation du système de types est décrite, ainsi que la manière de transformer automatiquement du code C en SAFESPEAK. Un cas d'étude est déroulé portant sur un bug ayant touché le noyau Linux. Il est démontré que le système de type capture précisément ce genre d'erreur de programmation. Enfin, les possibilités d'extension tant théoriques qu'expérimentales sont présentées.

Première partie

Méthodes formelles pour la sécurité

Après avoir décrit le contexte général de ces travaux, nous on décrivons les enjeux.

Le chapitre 2 explore plus en détail le fonctionnement d'un système d'exploitation, y compris la séparation du code en plusieurs niveaux de privilèges. L'architecture Intel 32 bits est prise comme support. En particulier, le mécanisme des appels système est décrit et on montre qu'une implantation naïve de la communication entre espaces utilisateur et noyau casse toute isolation.

Le chapitre 3 consiste en un tour d'horizon des techniques existantes en analyses de programmes. Ces analyses se centrent autour des problèmes liés à la vérification de code système ou embarqué, y compris le problème de manipulation mémoire évoqué dans le chapitre 2.

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.

2.1 Architecture physique

Un système informatique est principalement constitué d'un processeur (ou CPU pour *Central Processing Unit*), de mémoire principale (ou RAM pour *Random Access Memory*), et de divers périphériques.

Le processeur est constitué de divers registres internes qui permettent d'encoder l'état dans lequel il se trouve : quelle est l'instruction courante (registre EIP), quelle est la hauteur de la pile système (registre ESP), etc. Son fonctionnement peut être vu de la manière la plus simple qui soit comme le cycle :

- charger depuis la mémoire la prochaine instruction
- (*optionnel*) charger depuis la mémoire les données référencées par l'instruction
- effectuer l'instruction
- (*optionnel*) stocker en la mémoire les données modifiées

Les instructions sont constituées d'un mnémonique (indiquant quelle opération faire) et d'un ensemble d'opérandes. La signification des opérandes dépend du mnémonique, mais en général elles permettent de désigner les sources et la destination. Les opérandes peuvent avoir plusieurs formes : une valeur immédiate (\$4), un nom de registre (%eax) ou une référence à la mémoire (directement : addr ou indirectement : (%ecx)¹).

Parmi les mnémoniques courants on peut citer les 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.

1. Cela consiste à interpréter le contenu du registre ECX comme une adresse mémoire.

- `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`.

Certaines de ces instructions font référence à la pile par le biais du registre ESP. Cette zone mémoire n'est pas gérée de manière particulière. Elle permet de la pile des appels de fonction en cours grâce à la manière dont `jmp` et `ret` fonctionnent. Elle sert aussi à stocker les variables locales des fonctions.

À l'aide de ces quelques instructions on peut implanter des algorithmes impératifs. Mais pour faire quelque chose de visible, comme afficher à l'écran ou envoyer un paquet sur le réseau, cela ne suffit pas : il faut parler au reste du matériel.

Pour ceci, il y a deux techniques principales. D'une part, certains périphériques sont dits *memory-mapped* : ils sont associés à un espace mémoire particulier, qui ne permet pas de stocker des informations mais de lire ou d'écrire des données dans le périphérique. Par exemple, écrire à l'adresse `0xB8000` permet d'écrire des caractères à l'écran. L'autre système principal est l'utilisation des ports d'entrée/sortie. Cela correspond à des instructions spéciales `in %ax, port` et `out port, %ax` où `port` est un numéro qui correspond à un périphérique particulier. Par exemple, en écrivant dans le port `0x60` on peut contrôler l'état des indicateurs lumineux du clavier PS/2.

2.2 Tâches et niveaux de privilèges

Alternance des tâches

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.

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*) : 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*.

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 fuiter ; et s'il peut y écrire, il peut en détourner l'exécution.

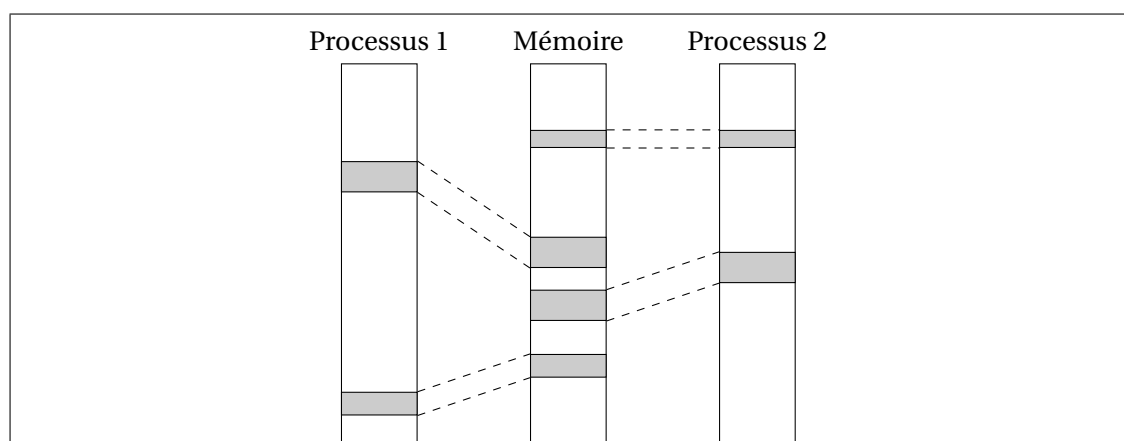


FIGURE 2.1: Mécanisme de mémoire virtuelle.

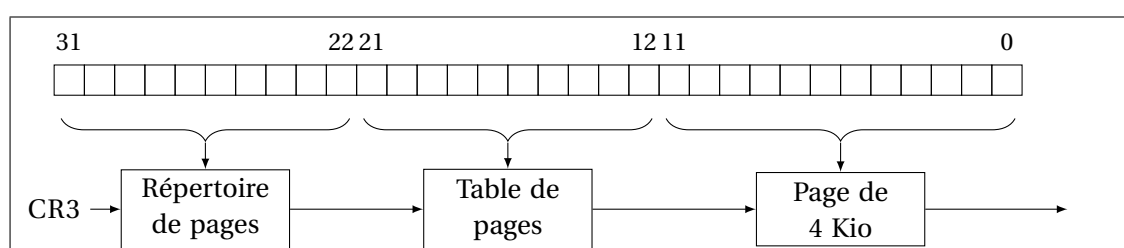


FIGURE 2.2: Implantation de la mémoire virtuelle

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 (figure 2.1).

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 kioctets (figure 2.2).

2.3 Appels système

Avec une telle isolation, tout le code qui est exécuté en *ring* 3 a une expressivité limitée. Il ne peut pas contenir d'instructions privilégiées comme `in` ou `out`, ni faire référence à des périphériques mappés en mémoire. C'est en effet au noyau d'accéder au matériel, et pas au code utilisateur.

Il est donc nécessaire d'appeler une routine du noyau depuis le code utilisateur. C'est le but des appels système. Cela consiste à coupler une fonction du *ring* 3 à une fonction du *ring* 0 : en appelant la fonction non privilégiée, le flot d'exécution se retrouve dans le noyau avec les bons privilèges.

Bien sûr, il n'est pas possible de faire directement un `call` puisque cela consisterait à faire un saut vers une zone plus privilégiée. Il y a plusieurs manières d'implanter ce mécanisme. Nous décrivons ici la technique historique à l'aide d'interruptions.

Le processeur peut répondre à des interruptions, qui sont des événements extérieurs. Cela permet d'écrire du code asynchrone. Par exemple, une fois qu'un long transfert mémoire est terminé, une interruption est reçue. D'autres interruptions dites logicielles peuvent arriver lorsqu'une erreur se produit. Par exemple, diviser par zéro provoque l'interruption 0, et

tenter d'exécuter une instruction privilégiée provoque l'interruption 14. On peut aussi provoquer une interruption par une instruction `int` dédiée.

Une table globale définit, pour chaque numéro d'interruption : quelle est la routine à appeler pour la traiter, avec quel niveau de privilège, ainsi que le niveau de privilège requis pour pouvoir déclencher celle-ci avec l'instruction `int`.

Il est donc possible de créer une interruption purement logicielle (on utilise en général le numéro 128 soit 0x80), déclenchable en *ring* 3 et traitée en *ring* 0. Les registres sont préservés, donc on peut les utiliser pour passer un numéro d'appel système (par exemple 3 pour `read()` et 5 pour `open()`) et leurs arguments.

2.4 Le *Confused Deputy Problem*

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 8.4.

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 2.3).

```
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 2.3: Appel de `gettimeofday`

Considérons une implémentation naïve de cet appel système qui écrirait directement à l'adresse pointée. Si le pointeur fourni est dans l'espace d'adressage du processus, on est dans le cas d'utilisation normal et l'écriture est donc possible.

```
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 2.4: Implantation de l'appel système gettimeofday

Si l'utilisateur passe un pointeur dont la valeur correspond à la mémoire réservée au noyau, que se passe-t-il? Comme le déréférencement est fait dans le code du 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 : s'il pointe en espace noyau, il faut indiquer une erreur plutôt que d'écrire (figure ??). Sinon, il peut toujours y avoir une 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 2.4).

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

Dans ce chapitre, nous présentons un tour d'horizon des techniques existantes permettant d'analyser des programmes. En particulier, on s'intéresse à la propriété de sécurité décrite dans le chapitre 2, mais on ne se limite pas à celle-ci : il est également intéressant de considérer des analyses développées pour d'autres propriétés, celles-ci pouvant potentiellement s'adapter.

L'analyse statique de programmes est un sujet de recherche actif depuis l'apparition de la science informatique. On commence par en présenter une classification, puis on montrera des exemples pertinents permettant d'analyser du code système ou embarqué.

3.1 Taxonomie

Techniques statiques et dynamiques L'analyse peut être faite au moment de la compilation, ou au moment de l'exécution. En général on peut obtenir des informations plus précises de manière dynamique, mais cela ne prend en compte que les parties du programme qui seront vraiment exécutées. Un autre problème des techniques dynamiques est qu'il est souvent nécessaire d'instrumenter l'environnement d'exécution (ce qui — dans le cas où cela est possible — peut se traduire par un impact en performances). L'approche statique, en revanche, nécessite de construire à l'arrêt une carte mentale du programme, ce qui n'est pas toujours possible dans certains langages.

Les techniques dynamiques sont néanmoins les plus répandues, puisqu'elles sont plus simples à mettre en œuvre et permettent de trouver des erreurs pendant le processus de développement. De plus, on peut considérer qu'un programme avec une forte couverture par les tests a de grandes chances d'être correct pour toutes les entrées. C'est par exemple ce qui se passe dans l'avionique civile.

Mais pour s'assurer de la correction d'un programme, on ne peut pas s'appuyer uniquement sur les tests — ou de manière générale sur des analyses dynamiques — car il est impossible d'étudier l'ensemble complet de tous les comportements possibles. Par exemple, si un bug se présente lors d'une interaction entre deux composants qui n'a pas été testée, il passera inaperçu. Pour cette raison, la plupart des analyses présentées ici sont statiques.

Cohérence et complétude Le but d'une analyse statique est de catégoriser les programmes selon s'ils satisfont ou non un ensemble de propriétés fixées à l'avance. Malheureusement, cela n'est que rarement possible car l'ensemble des valeurs possibles lors de l'exécution d'un programme n'est pas un ensemble calculable (théorème de Rice [Ric53]). Autrement dit, il

ne peut exister une procédure de décision prenant un programme et le déclarant correct ou incorrect.

Il n'est donc pas possible d'écrire un analyseur statique parfait, ne se trompant jamais. Toute technique statique va donc de se retrouver dans au moins un des cas suivants :

- un programme valide (pour une propriété donnée) est rejeté : on parle de *faux positif*.
- un programme invalide n'est pas détecté : on parle de *faux négatif*.

En général, et dans notre cas, on préfère s'assurer que les programmes acceptés possèdent la propriété recherchée, quitte à en rejeter certains. C'est l'approche que nous retiendrons. Tolérer les faux négatifs n'est cependant pas toujours une mauvaise idée. Par exemple, si le but est de trouver des constructions dangereuses dans les programmes, on peut signaler certains cas qui empiriquement valent d'être vérifiés manuellement.

Par ailleurs la plupart des techniques ne concernent que les programmes qui terminent. On étudie donc la correction, ou les propriétés des termes convergents. Prouver automatiquement que l'exécution ne boucle pas est une propriété toute autre qui n'est pas ici considérée.

3.2 Méthodes syntaxiques

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.

Ces techniques sont utiles parce qu'elles permettent de plonger rapidement dans le code, en identifiant par exemple des appels à des fonctions dangereuses. En revanche, cela n'est utile que lorsque les propriétés que l'on recherche sont très locales. Elles offrent également peu de garantie (ce n'est pas parce qu'une construction dangereuse n'est pas présente que le code est correct).

3.3 Interprétation abstraite

L'interprétation abstraite est une technique d'analyse générique qui permet de simuler statiquement tous les comportements d'un programme [CC77, CC92]. Un exemple d'application est de calculer les bornes de variation des variables pour s'assurer qu'aucun débordement de tableau n'est possible [AH07].

L'idée est d'associer à chaque ensemble concret de valeurs, une représentation abstraite. Sur celle-ci, on peut définir des opérations indépendantes de la valeur exacte des données, mais préservant l'abstraction (figure 3.1). Par exemple, les règles comme $- \times - = +$ définissent le domaine abstrait des signes. Les domaines ont une structure de treillis, c'est-à-dire qu'ils possèdent les notions d'ordre partiel et d'union de valeurs. En calculant les valeurs extrémales d'une variable, on obtient le domaine des intervalles.

De tels domaines ne capturent aucune relation entre variables. Ils sont dits non relationnels. Lorsque plusieurs variables sont analysées en même temps, utiliser de tels domaines consiste à considérer un produit cartésien d'ensembles abstraits (figure 3.2(a)).

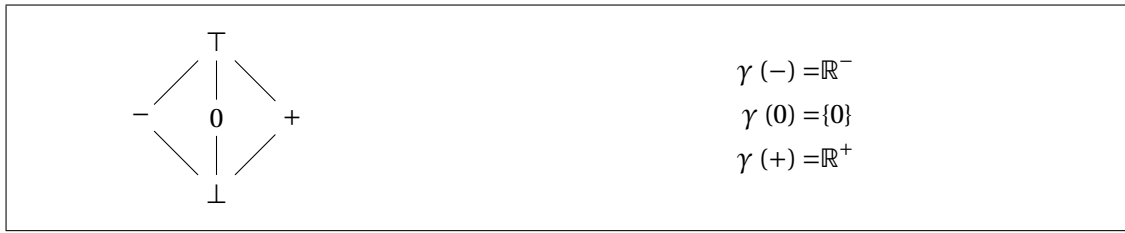


FIGURE 3.1: Domaine des signes

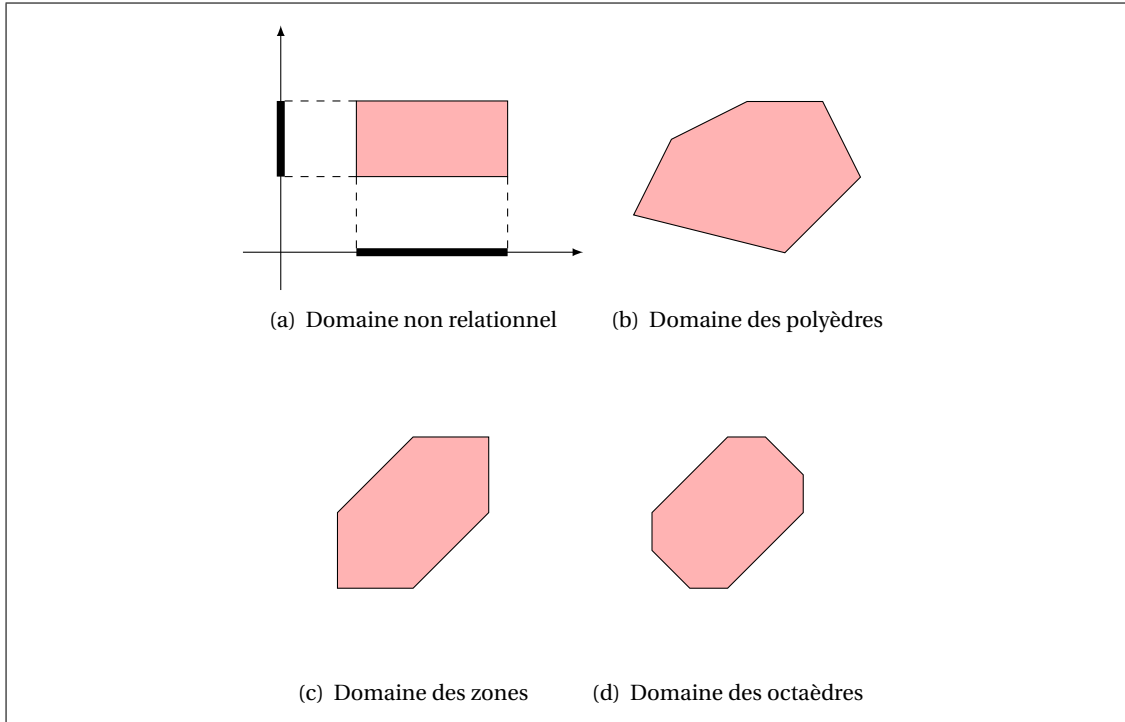


FIGURE 3.2: Quelques domaines abstraits

Des domaines abstraits plus précis permettent de retenir celles-ci. Pour ce faire, il faut modéliser l'ensemble des valeurs des variables comme un tout. Parmi les domaines relationnels courants on peut citer : le domaine des polyèdres, permettant de retenir tous les invariants affines entre fonctions (figure 3.2(b)) ; le domaine des zones, permettant de représenter des relations affines de la forme $v_i - v_j \leq c$ (figure 3.2(c)) ; ou encore le domaine des octogones qui est un compromis entre les polyèdres et les zones. Il permet de représenter les relations $\pm v_i \pm v_j \leq c$ (figure 3.2(d)).

En plus des domaines numériques, il est nécessaire d'employer des domaines spécialisés dans la modélisation de la mémoire. Cela est nécessaire pour pouvoir prendre en compte les pointeurs. Par exemple, on peut représenter un pointeur par un ensemble de variables possiblement pointées et une valeur abstraite représentant le décalage (*offset*) du pointeur par rapport au début de la zone mémoire. Cette valeur peut elle-même être abstraite par un domaine numérique.

Au delà des domaines eux-mêmes, l'analyse se fait sous forme d'un calcul de point fixe. La manière la plus simple est d'utiliser un algorithme de *liste de travail*, décrit par exemple dans [SRH95]. Les raffinements en revanche sont nombreux.

Dès [CC77] il est remarqué que la terminaison de l'analyse n'est assurée que si le treillis des valeurs abstraites est de hauteur finie, ou qu'un opérateur d'élargissement (*widening*) ∇ est employé. L'idée est qu'une fois qu'on a calculé quelques termes d'une suite croissante,

on peut réaliser une projection de celle ci. Par exemple, dans le domaine des intervalles, $[0; 2] \vee [0; 3] = [0; +\infty[$. On atteint alors un point fixe mais qui est plus grand que celui qu'on aurait obtenu sans cette accélération : on perd en précision. Pour en gagner, on peut redescendre sur le treillis des points fixe avec une suite d'itérations décroissantes [Gra92]. Dans l'itération de point fixe, il est possible d'obtenir les résultats de manière plus efficace en choisissant un ordre particulier dans les calculs des sous-itérations[GGTZ07].

En termes d'ingénierie logicielle, implanter un analyseur statique est un défi en soi. En plus des domaines abstraits, d'un itérateur, il faut traduire le code source à analyser dans un langage, et traduire les résultats de l'analyse en un ensemble d'alarmes à présenter à l'utilisateur.

Cette technique est très puissante : si un interpréteur abstrait *sound* (c'est-à-dire réalisant une surapproximation) analyse un programme et ne renvoie pas d'erreur, alors on a prouvé que le programme est correct (par rapport aux propriétés que vérifient les domaines abstraits). Cela a été appliqué avec succès avec les analyseurs Astrée [Mau04, CCF⁺05, CCF⁺09] chez Airbus ou CGS [VB04] à la NASA par exemple.

Cependant, ces analyses sont difficiles à mettre en œuvre. Avec des domaines abstraits classiques comme ceux présentés ci-dessus, les premières analyses peuvent remonter un nombre prohibitif de fausses alarmes. Pour "aider" l'analyse, il faut soit annoter le code soit développer des domaines abstraits *ad hoc* au programme à analyser.

Il existe également des analyseurs statiques par interprétation abstraite qui ne sont pas *sound*, c'est-à-dire qui peuvent manquer des comportements erronés. Leur approche est plus d'aider le programmeur à détecter certains types de bugs pendant le développement. On citer l'exemple de Coverity [BBC⁺10], qui publie régulièrement des rapports de qualité sur certains logiciels open source. Néanmoins, de part leur aspect non *sound*, les analyses réalisées ne peuvent pas être assimilées à de la vérification formelle en tant que telle.

3.4 Typage

Le typage, introduit dans la section 1.3, peut aussi être utilisé pour la vérification de programmes. On peut le voir comme une manière de catégoriser les types de données manipulés par la machine, mais également à plus au niveau comme une manière d'articuler les différents composants d'un programme. Mais on peut aussi programmer avec les types, c'est-à-dire utiliser le compilateur (dans le cas statique) ou l'environnement d'exécution (dans le cas dynamique) pour vérifier des propriétés écrites par le programmeur.

Systèmes ad hoc Les systèmes de types les plus simples expriment des contrats essentiellement liés à la sûreté d'exécution, pour ne pas utiliser des valeurs de types incompatibles entre eux. Mais il est possible d'étendre le langage avec des annotations plus riches : par exemple en vérifiant statiquement que des listes ne sont pas vides[KcS07], ou dans le domaine de la sécurité, d'empêcher des fuites d'information [LZ06].

Qualificateurs de types 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. Celles-ci sont ignorées par le compilateur, mais un outil d'analyse statique ad-hoc nommé Sparse [S05] peut être utilisé pour détecter les cas les plus simples d'erreurs. Il demande aussi au programmeur d'ajouter de nombreuses annotations dans le programme.

Cette solution se rapproche de la solution décrite dans ce manuscrit. Ce système d'annotations sur les types a été formalisé sous le nom de *qualificateurs de types* [FJKA06] : chaque

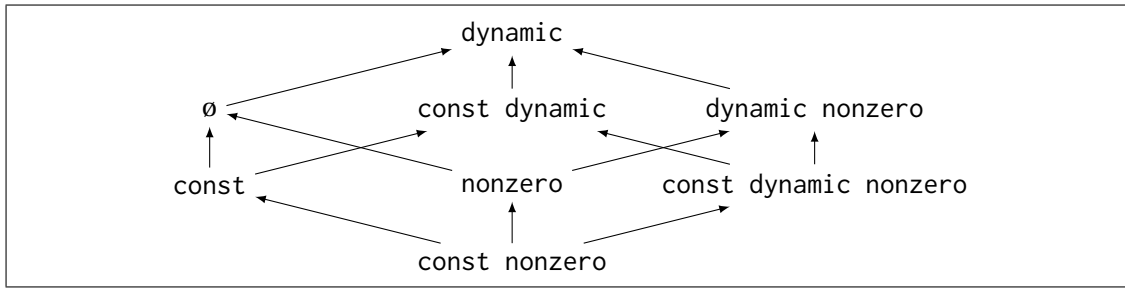


FIGURE 3.3: Treillis de qualificateurs

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.

Plus précisément, les jugements de typage de la forme $\Gamma \vdash e : t$ sont remplacés par des jugements de typage qualifiés $\Gamma \vdash e : t \ q$. Les qualificateurs q permettent d'exprimer plusieurs jugements. Par exemple, pointeurs on peut étudier le fait qu'une variable soit constante ou pas, que sa valeur soit connue à la compilation, ou encore qu'elle puisse être nulle ou pas. La spécificité de ce système est que les qualificateurs sont ordonnés, du plus spécifique au moins spécifique, et que l'on forme alors un treillis à partir de ces informations. Partant des deux caractéristiques précédentes, on forme le treillis de la figure 3.3. Le qualificateur `const` désigne les données dont la valeur change au cours de l'exécution ; `dynamic` celles qui peuvent être connues à la compilation ; et `nonzero` celles qui ne peuvent jamais être nulles. Le cube sur lequel se trouvent correspond à une relation d'ordre, du plus spécifique (en bas) au plus général (en haut).

Cette relation d'ordre \leq entre qualificateurs, induit une relation de sous-typage \sqsubseteq entre les types qualifiés : si $q \leq q'$, alors $t \ q \sqsubseteq t \ q'$.

Ces analyses ont été implantée dans l'outil CQual. Ce système peut servir à inférer les annotations `const` [FFA99], à l'analyse de souillure pour les chaînes de format [STFW01] (pouvant poser des problèmes de sécurité [New00]) et des propriétés dépendantes du flot de contrôle, comme des invariants sur les verrous [FTA02], à rapprocher du concept de *types-tates* [SY86]. Il a également été appliqué à la classe de vulnérabilités sur les pointeurs utilisateurs dont il est ici l'objet [JW04].

Cette approche étant assez proche de la nôtre, la conclusion de la partie II (page 87) sera dédiée à un comparatif entre les deux.

Ces techniques de typage sont séduisantes parce qu'elles sont en général simples à mettre en place : à l'aide d'un ensemble de règles, on attribue un type à chaque expression. Si le typage se termine sans erreur, alors on est assuré de la correction du programme (par rapport aux propriétés capturées par le système de types).

Le typage statique peut également être implanté de manière efficace. La plupart des systèmes de types peuvent être vérifiés en pratique dans un temps linéaire en la taille du programme considéré. L'inférence, en revanche, peut parfois "exploser" en temps quadratique.

3.5 Analyse de code système

Les logiciels système demandent des garanties de sécurité et de fiabilité. Pour les assurer, de nombreuses analyses *ad-hoc* ciblent les noyaux de systèmes d'exploitation.

Ajouter un système de types forts à C est l'idée centrale de CCured [NCH⁺05]. Dans les cas où il n'est pas possible de conclure, des vérifications à l'exécution sont ajoutées. Cependant, cela nécessite une instrumentation dynamique qui se paye en performances.

Saturn [ABD⁺07] est un système pour analyser du code système écrit en C. Il traite le problème des pointeurs utilisateur en utilisant une analyse de forme “pointe-sur” [BA08]. Comme l’interprétation abstraite, son but est d’être très précis, au détriment d’un temps de calcul important dans certains cas.

3.6 Langages sûrs

Une autre approche est de concevoir un langage à la fois bas niveau et sûr, permettant d’exprimer des programmes proches de la machine tout en interdisant les constructions dangereuses.

Le langage Cyclone [JMG⁺02] est conçu comme un C “sûr”. Afin d’apporter plus de sûreté au modèle mémoire de C, des tests dynamiques sont ajoutés, par exemple aux endroits où des conversions implicites peuvent poser problème. Le langage se distingue par le fait qu’il possède plusieurs types de pointeurs : des pointeurs classiques (`int *`), des pointeurs “jamais nuls” (`int @`; un test à l’exécution est alors inséré), et des “pointeurs lourds” (`int ?`; qui contiennent des informations sur la zone mémoire pointée). L’arithmétique des pointeurs n’est autorisée que sur ces derniers, rendant impossibles les débordements de tableaux (ceux-ci étant détectés au pire à l’exécution). Le problème des *dangling pointers* est résolu en utilisant un système de régions [GMJ⁺02], inspiré des travaux de Jouvelot, Talpin et Tofte [TJ92, TT93, TT94]. Cela permet d’interdire statiquement les constructions où l’on déréfère un pointeur faisant référence à une région de mémoire qui n’est plus allouée.

Le langage Rust [R⁹] développé par Mozilla prend une approche similaire en distinguant plusieurs types de pointeurs pour gérer la mémoire de manière plus fine. Les *managed pointers* (notés `@int`) utilisent un ramasse-miette pour libérer la mémoire allouée lorsqu’ils ne sont plus accessibles. Les *owning pointers* (notés `~int`) décrivent une relation 1 à 1 entre deux objets, comme les `std::unique_ptr` de C++ : la mémoire est libérée lorsque le pointeur l’est. Les *borrowed pointers* (notés `&int`) correspondent aux pointeurs obtenus en prenant l’adresse d’un objet, ou d’un champ d’un objet. Une analyse statique faite lors de la compilation s’assure que la durée de vie de ces pointeurs est plus longue que l’objet pointé, afin d’éviter les *dangling pointers* (pointeurs fous). Cette analyse est également basée sur les régions. Une fonction qui retourne l’adresse d’une variable locale sera donc rejetée par le compilateur. Enfin, le dernier type est celui des *raw pointers* (notés `*int`), pour lesquels le langage n’apporte aucune garantie (il faut d’ailleurs encapsuler chaque utilisation dans un bloc marqué explicitement `unsafe`). Ils sont équivalents aux pointeurs de C.

Ces techniques sont utiles pour créer des nouveaux programmes sûrs, mais on ne peut pas les appliquer pour étudier la correction de logiciels existants. Leurs systèmes de types apportent dans le langage différents types de pointeurs, permettant de manipuler finement la mémoire, à la manière des *smart pointers* de C++. De cette approche on retient surtout l’analyse de régions de Rust qui permet de manipuler de manière sûre les adresses des variables locales, et les pointeurs lourds de Cyclone qui apportent une sûreté à l’arithmétique de pointeurs, au prix d’un test dynamique.

3.7 Logique de Hoare

Une technique pour vérifier statiquement des propriétés sur la sémantique d’un programme a été formalisée par Robert Floyd [Flo67] et Tony Hoare [Hoa69].

Elle consiste à écrire les invariants qui sont maintenus à un point donné du programme. Ces propositions sont écrites dans une logique \mathcal{L} . Chaque instruction i est annotée d’une

pré-condition P et d'une post-condition Q , ce que l'on note $\{P\} i \{Q\}$. Cela signifie que si P est vérifiée et que l'exécution de i se termine, alors Q sera vérifiée.

En plus des règles de \mathcal{L} , des règles d'inférence traduisent la sémantique du programme ; par exemple la règle de composition est :

$$\frac{\{P\} i_1 \{Q\} \quad \{Q\} i_2 \{R\}}{\{P\} i_1; i_2 \{R\}} \text{ (HOARE-SEQ)}$$

Les pré-conditions peuvent être renforcées et les post-conditions relâchées :

$$\frac{\vdash_{\mathcal{L}} P \Rightarrow P' \quad \{P\} i \{Q\} \quad \vdash_{\mathcal{L}} Q' \Rightarrow Q}{\{P'\} i \{Q'\}} \text{ (HOARE-CONSEQUENCE)}$$

Il est alors possible d'annoter le programme avec ses invariants formalisés de manière explicite dans \mathcal{L} . Ceux-ci seront vérifiés à la compilation lorsque c'est possible, sinon à l'exécution.

La règle de conséquence permet de découpler les propriétés du programme lui-même : plusieurs niveaux d'annotations sont possibles, du moins précis au plus précis. En fait, il est même possible d'annoter chaque point de contrôle par l'ensemble d'annotations vide : $\{T\} i \{T\}$ est toujours vrai.

Augmenter graduellement les pré- et post-conditions est néanmoins assez difficile, puisqu'il peut être nécessaire de modifier l'ensemble des conditions à la fois. Cette difficulté est mentionnée dans [DRS00], où un système de programmation par contrats est utilisé pour vérifier la correction de routines de manipulation de chaînes en C.

Ce type d'annotations a été implémenté par exemple pour le langage Java dans le système JML [LBR99] ou pour le langage C# dans Spec# [BLS05].

3.8 Assistants de preuve

Avec un système de types classique, le fait qu'un terme soit bien typé amène quelques propriétés sur son exécution. Par exemple, le fait que seulement un ensemble réduit d'erreurs puisse arriver (comme la division par zéro).

En enrichissant le langage des types, on peut augmenter l'expressivité du typage. Par exemple, on peut former des types "entier pair", "vecteur de n entiers", ou encore "liste triée d'entiers".

Habituellement, les termes peuvent dépendre d'autres termes (par composition) ou de types (par des annotations). Les types peuvent également dépendre d'autres types (par composition de types : par exemple, un couple de a et de b a pour type $a * b$). Enrichir l'expressivité du typage revient essentiellement à introduire des termes dans les types, comme n dans l'exemple précédent du vecteurs de n entiers. C'est pourquoi on parle de types dépendants. Parmi les langages proposant ces types on peut citer Coq [The04], Agda [BDN09] ou Isabelle [NPW02].

Dans un langage classique, la plupart des types sont habitables, c'est à dire qu'il existe des termes ayant ces types. En revanche, avec les types dépendants ce n'est pas toujours vrai : par exemple "vecteur de -1 entiers" n'a pas d'habitants. Ainsi, pouvoir construire un terme d'un type donné est une information en soi.

On peut voir ce phénomène sous un autre angle : les termes sont à leur type ce que les preuves sont à leur théorème. C'est la correspondance de Curry-Howard. Prenons un

exemple : si on dispose d'un terme de type a , et d'une fonction de type $a \rightarrow b$, alors on peut en déduire (par simple application) un terme de type b . Vu sous forme logique : si on a une preuve de A , et une preuve de $A \Rightarrow B$, alors on peut en tirer une preuve de B . À l'aide de cette correspondance, on peut donc voir les systèmes de types dépendants comme des assistants de preuves généralistes où l'on peut faire vérifier par un compilateur des théorèmes mathématiques comme le théorème des 4 couleurs [Gon05].

Cette technique fait se rejoindre les mondes du typage et des annotations par triplets de Hoare : on enrichit le langage des types suffisamment pour encoder une logique, et on vérifie suffisamment les triplets statiquement pour éliminer tous les tests dynamiques. Cependant cette technique est très complexe à mettre en œuvre, puisqu'il faut encoder toutes les propriétés voulues dans un formalisme de très bas niveau (similaire à la théorie des ensembles).

3.9 Conclusion

Il existe de nombreuses techniques pour vérifier du code système ou embarqué. Il y a divers choix à faire entre l'expressivité, l'intégration de tests dynamiques, ou la facilité de mise en œuvre.

Pour résoudre les problèmes des pointeurs utilisateurs dans les noyaux, le typage statique est une solution assez pragmatique, puisqu'elle peut s'appliquer à des programmes existants, et est efficace à utiliser. Son expressivité limitée nous empêche de reposer entièrement sur elle pour garantir l'absence d'erreur dans les programmes système. C'est pourquoi dans notre solution, SAFESPEAK, nous approchons la sûreté de la manière suivante :

Tout d'abord, on utilise le typage pour manipuler les données de manière compatible : les types des opérations et fonctions sont vérifiés à la compilation.

Ensuite, les accès aux tableaux et aux pointeurs sont vérifiés dynamiquement. Dans le cas où une erreur est déclenchée, l'exécution s'arrête plutôt que de corrompre la mémoire. La pile est également nettoyée à chaque retour de fonction afin d'éviter les pointeurs fous (*dangling pointers*).

Enfin, les pointeurs provenant de l'espace utilisateur sont repérés statiquement afin que leur déréférencement se fasse au sein de fonctions sûres. Cela permet de préserver l'isolation entre le noyau et l'espace utilisateur.

CONCLUSION DE LA PARTIE I

Nous avons montré que l'écriture de noyaux de systèmes d'exploitation nécessite de manipuler des données provenant d'une zone non sûre, l'espace utilisateur. Parmi ces données, il arrive de récupérer des pointeurs qui servent à passer des données par référence à l'appellant, dans certains appels systèmes. Si on déréférence ces pointeurs sans vérifier qu'ils pointent bien vers une zone mémoire également contrôlée par l'appellant, on risque de lire ou d'écrire dans des zones mémoires réservées au noyau seul.

Nous proposons une technique de typage pour détecter ces cas dangereux. Pour ce faire, il faudra tout d'abord définir un langage impératif bien typé que nous appellerons SAFESPEAK. Ensuite, nous distinguons dans le système de types les pointeurs dont la valeur est contrôlée par l'utilisateur.

Cela revient à introduire une séparation entre le type des données et sa représentation : c'est-à-dire définir un type abstrait. Les pointeurs utilisateur sont en fait des pointeurs classiques déguisés, pour lesquels on interdit l'opérateur de déréférencement.

Cette technique est en fait générique : on peut également l'appliquer à certains types d'entiers. En C, il est commun d'utiliser des `int` pour tout et n'importe quoi : pour des entiers bien sûr (au sens de \mathbb{Z}), mais aussi comme identificateurs pour lesquels les opérations usuelles comme l'addition n'ont pas de sens. Par exemple, sous Linux, l'opération d'ouverture de fichier renvoie un entier, dit *descripteur de fichier*, qui identifie ce fichier pour ce processus. Il est passé à toutes les fonctions de lecture et d'écriture. Un autre idiome est d'utiliser des entiers comme des ensembles de bits (on parle de *bitmask*). Un prototype d'analyse statique sur ces *bitmasks* est présenté dans l'annexe ??.

En n'offrant pas cette distinction, le langage C permet d'écrire du code qui peut s'exécuter mais dont la sémantique n'est en quelque sorte pas bien fondée. Par exemple, en additionnant deux descripteurs de fichiers, ou en déréférençant un pointeur contrôlé par l'utilisateur. En effet, le système de types de C est trop primitif pour pouvoir garantir une véritable isolation entre deux types de même représentation : il n'y a pas de types abstraits. Certes, `typedef` permet d'introduire un nouveau nom pour un type, mais ce n'est qu'un raccourci syntaxique. Le compilateur ne peut en effet pas considérer un programme sans avoir la définition quasi-complète des types qui y apparaissent. La seule exception concerne les pointeurs sur structures : si on ne fait que les passer, il n'est pas nécessaire de connaître la taille ni la disposition de la structure, donc il est possible de ne pas connaître ces informations. Cette technique, connue sous le nom de *pointeurs opaques*, n'est pas applicable aux autres types.

En ajoutant une couche de typage, on interdit ces opérations à la compilation. Cela permet deux choses : pour le code déjà écrit, de détecter et corriger les manipulations dangereuses ; et pour le nouveau code, de s'assurer qu'il est correct.

Le but de cette thèse est donc de définir un langage intermédiaire proche de C, mais bien typable, et de l'adjoindre d'un système de types tel que les programmes bien typés manipulent les pointeurs utilisateur sans causer de problèmes de sécurité.

Deuxième partie

Un langage pour l'analyse de code système : SAFESPEAK

Dans cette partie, nous allons présenter un langage impératif modélisant le langage C. Le chapitre 4 décrit sa syntaxe, ainsi que sa sémantique. À ce point, de nombreux programmes sont acceptés mais qui provoquent des erreurs à l'exécution.

Afin de rejeter ces programmes incorrects, on définit ensuite dans le chapitre 5 une sémantique statique s'appuyant sur un système de types simples. Des propriétés de sûreté de typage sont ensuite établies, permettant de catégoriser l'ensemble des erreurs à l'exécution possibles.

Le chapitre 6 commence par étendre notre langage avec une nouvelle classe d'erreurs à l'exécution, modélisant les accès à la mémoire utilisateur catégorisé comme dangereux dans le chapitre 2. Une extension au système de types du chapitre 5 est ensuite établie, et on prouve que les programmes ainsi typés ne peuvent pas atteindre ces cas d'erreur.

Trois types d'erreurs à l'exécution sont possibles :

- les erreurs de typage (dynamique), lorsqu'on tente d'appliquer à une opération des valeurs incompatibles (additionner un entier et une fonction par exemple).
- les erreurs de sécurité, qui consistent en le déréférencement d'un pointeur dont la valeur est contrôlée par l'espace utilisateur. Celles-ci sont uniquement possibles en contexte noyau.
- les erreurs mémoire, qui résultent d'un débordement de tableau, du déréférencement d'un pointeur invalide ou d'arithmétique de pointeur invalide.

En résumé, l'introduction des types simples enlève la possibilité de rencontrer des erreurs de typage dynamique, et l'ajout des qualificateurs interdit les erreurs de sécurité.

Langage	Types	Erreurs possibles		
		Typage	Sécurité	Mémoire
SAFESPEAK	sans	☑	N/A	☑
SAFESPEAK	simples	☐	N/A	☑
SAFESPEAK noyau	simples	☐	☑	☑
SAFESPEAK noyau	qualifiés	☐	☐	☑

SYNTAXE ET SÉMANTIQUE

Dans ce chapitre nous présentons SAFESPEAK, un langage impératif inspiré de C. Sa syntaxe est tout d’abord décrite ; puis une sémantique opérationnelle est explicitée.

Ce langage servira de support aux systèmes de types décrits dans le chapitre 5 et enrichi dans le chapitre 6.

La traduction depuis C sera explicitée dans le chapitre 7.

4.1 Notations

Ensembles inductifs

Dans ce chapitre (et les chapitres suivants), on définit de nombreux ensembles inductifs. Plutôt que d’écrire la construction explicite par point fixe, on emploie une notation en grammaire.

Étudions l’exemple des listes chaînées composées d’éléments de \mathbb{N} .

Notons L cet ensemble ; si $[]$ est la liste vide et $n :: l$ la liste formée d’une “tête” $n \in \mathbb{N}$ et d’une “queue” $l \in L$. Toute liste est donc d’une des formes suivantes :

- $[]$
- $n_1 :: []$
- $n_1 :: n_2 :: []$
- etc.

On peut donc L de la manière inductive suivante :

$$L = \text{fix}(L')$$

$$L'(E) = \{[]\} \cup \{n :: l \mid n \in \mathbb{N}, l \in E\}$$

où

$$\text{fix}(f) = \bigcup_{n=0}^{\infty} f^n(\emptyset)$$

$$f^0(x) = x$$

$$\forall n > 0, f^n(x) = f^{n-1}(f(x))$$

(L'itération n de l'union correspond aux listes comprenant au plus n éléments)

Plutôt que d'écrire cette définition précise mais chargée, on écrira à la place une définition en compréhension :

Listes	$l ::= []$	Liste vide
	$ \quad n :: l$	Construction de liste

Chaque ensemble est identifié de manière unique par les noms de variables métasyntaxiques : n pour les entiers et l pour les listes ici. Si plusieurs métavariabes du même ensemble doivent apparaître, elles sont indicées. Par exemple, on peut définir des arbres binaires d'entiers de la manière suivante :

Arbres	bi-	$a ::= F$	Feuille
nares		$ \quad N(a_1, n, a_2)$	Nœud

Cette notation a aussi l'avantage de s'étendre facilement aux définitions mutuellement récursives.

Inférence

La sémantique opérationnelle consiste en la définition d'une relation de transition $\cdot \rightarrow \cdot$ entre états de l'interpréteur¹.

Cette relation est définie inductivement sur la syntaxe du programme. Plutôt que de présenter l'induction explicitement, elle est représentée par des jugements logiques et des règles d'inférences, de la forme :

$$\frac{P_1 \quad \dots \quad P_n}{C} \text{ (NOM)}$$

Les P_i sont les prémisses, et C la conclusion. Cette règle s'interprète de la manière suivante : si les P_i sont prouvées, alors C est prouvée.

Certaines règles n'ont pas de prémisses, ce sont des axiomes :

$$\frac{}{A} \text{ (Ax)}$$

Compte-tenu de la structure des règles, la preuve d'un jugement pourra donc être vue sous la forme d'un arbre :

$$\frac{\frac{\frac{}{A_1} \text{ (R3)} \quad \frac{}{A_2} \text{ (R4)}}{B_1} \text{ (R2)} \quad \frac{\frac{}{A_3} \text{ (R6)}}{B_2} \text{ (R5)}}{C} \text{ (R1)}$$

1. Dans le chapitre 5, la relation de typage $\cdot \vdash \cdot$ sera définie par la même technique.

Listes

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 $i \in [1; |u|]$, 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 :

$$\begin{aligned} v_1 &= x \\ \forall i \in [1; |u|], v_{i+1} &= u_i \end{aligned}$$

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

$$\begin{aligned} |w| &= |u| + |v| \\ \forall i \in [1; |u|], w_i &= u_i \\ \forall j \in [1; |v|], w_{|u|+j} &= v_j \end{aligned}$$

Lentilles

La notion d'accessor utilisée ici est directement inspirée des *lentilles* utilisées en programmation fonctionnelle, décrite dans [FGM⁺07] et [vL11].

Définition 4.1 (Lentille). *Étant donnés deux ensembles R et A , une lentille $\mathcal{L} \in \text{Lens}_{R,A}$ (ou accessor) est un moyen d'accéder en lecture ou en écriture à sous-valeur de type A au sein d'une valeur de type R (pour record). Elle est constituée des opérations suivantes :*

- une fonction de lecture $\text{get}_{\mathcal{L}} : R \rightarrow A$
- une fonction de mise à jour $\text{put}_{\mathcal{L}} : (A \times R) \rightarrow R$

telles que pour tous $a \in A, a' \in A, r \in R$:

$$\begin{aligned} \text{put}_{\mathcal{L}}(\text{get}_{\mathcal{L}}(r), r) &= r && (\text{GetPut}) \\ \text{get}_{\mathcal{L}}(\text{put}_{\mathcal{L}}(a, r)) &= a && (\text{PutGet}) \\ \text{put}_{\mathcal{L}}(a', \text{put}_{\mathcal{L}}(a, r)) &= \text{put}_{\mathcal{L}}(a', r) && (\text{PutPut}) \end{aligned}$$

On note $\mathcal{L} = \langle \text{get}_{\mathcal{L}} | \text{put}_{\mathcal{L}} \rangle$.

GetPut signifie que si on lit une valeur puis qu'on la réécrit, l'objet n'est pas modifié ; PutGet décrit l'opération inverse : si on écrit une valeur dans le champ, c'est la valeur qui sera lue ; enfin, PutPut évoque le fait que chaque écriture est totale : quand deux écritures se suivent, seule la seconde compte.

Une illustration se trouve dans la figure 4.1.

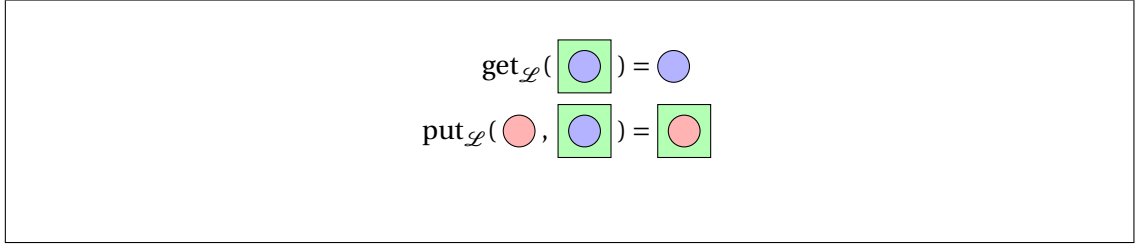


FIGURE 4.1: Fonctionnement d'une lentille

Exemple 4.1 (Lentilles de tête et de queue de liste). Soit E un ensemble. On considère $L(E)$, l'ensemble des listes d'éléments de E .

On définit les fonctions suivantes. Notons qu'elles ne sont pas définies sur la liste vide $[]$, qui pourra être traité comme un cas d'erreur.

$$\begin{aligned} \text{get}_T(t :: q) &= t \\ \text{put}_T(t', t :: q) &= t' :: q \\ \text{get}_Q(t :: q) &= q \\ \text{put}_Q(q', t :: q) &= t :: q' \end{aligned}$$

Alors $T = \langle \text{get}_T | \text{put}_T \rangle \in \text{Lens}_{L(E), E}$ et $Q = \langle \text{get}_Q | \text{put}_Q \rangle \in \text{Lens}_{L(E), L(E)}$.

On a par exemple :

$$\text{get}_T(1 :: 6 :: 1 :: 8 :: []) = 1$$

et :

$$\text{put}_Q(7, 3 :: 6 :: 1 :: 5 :: []) = 7 :: 6 :: 1 :: 5 :: [].$$

Définition 4.2 (Lentille indexée). Les objets de certains ensembles R sont composés de plusieurs sous-objets accessibles à travers un indice $i \in I$. Une lentille indexée est une fonction Δ qui associe à un indice i une lentille entre R et un de ses champs A_i :

$$\forall i \in I, \exists A_i, \Delta(i) \in \text{Lens}_{R, A_i}$$

On note alors :

$$\begin{aligned} r[i]_{\Delta} &\stackrel{\text{def}}{=} \text{get}_{\Delta(i)}(r) \\ r[i \leftarrow a]_{\Delta} &\stackrel{\text{def}}{=} \text{put}_{\Delta(i)}(a, r) \end{aligned}$$

Un exemple est illustré dans la figure 4.2.

Exemple 4.2 (Lentille “n^e élément d'un tuple”). Soient $n \in \mathbb{N}$, et n ensembles E_1, \dots, E_n .

Pour tout $i \in [1; n]$, on définit :

$$\begin{aligned} g_i((x_1, \dots, x_n)) &= x_i \\ p_i(y, (x_1, \dots, x_n)) &= (x_1, \dots, x_{i-1}, y, x_{i+1}, \dots, x_n) \end{aligned}$$

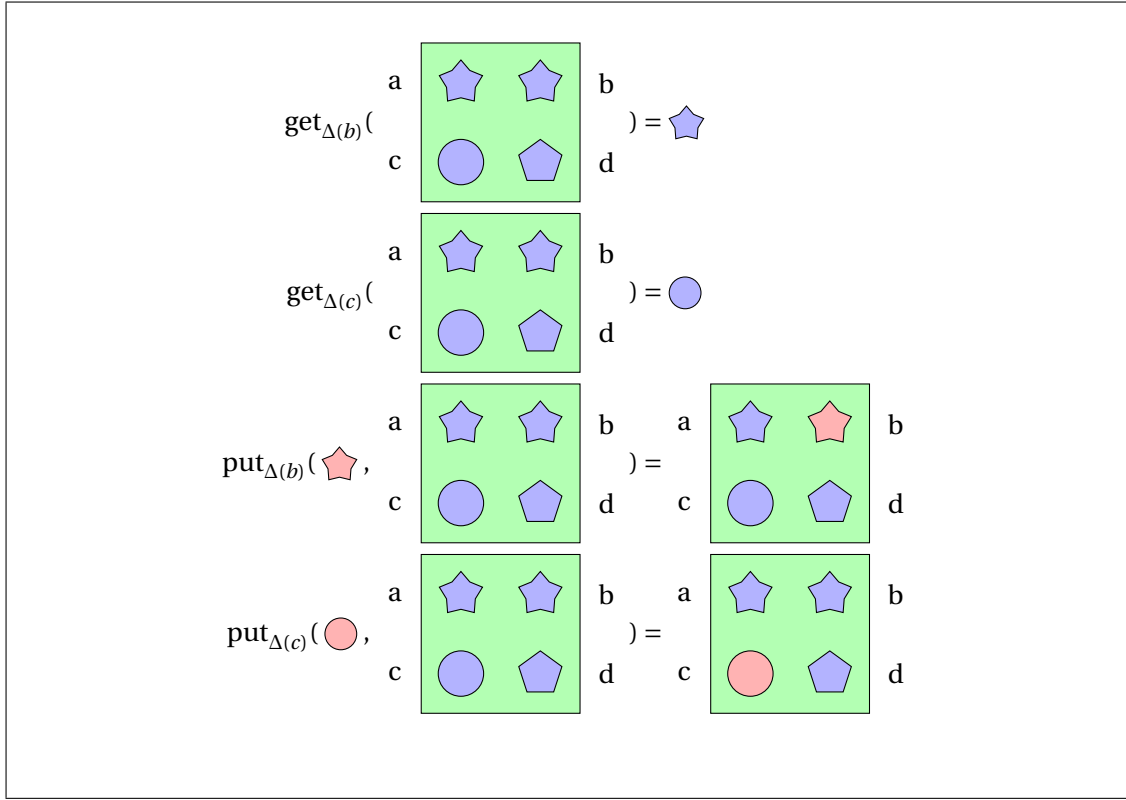


FIGURE 4.2: Fonctionnement d'une lentille indexée

Définissons $T(i) = \langle g_i | p_i \rangle$. Alors $T(i) \in \text{Lens}_{(E_1 \times \dots \times E_n), E_i}$.

Donc T est une lentille indexée, et on a par exemple :

$$(3, 1, 4, 1, 5)[2]_T = \text{get}_{T(2)}((3, 1, 4, 1, 5)) \\ = 1$$

$$(9, 2, 6, 5, 3)[3 \leftarrow 1]_T = \text{put}_{T(3)}(1, (9, 2, 6, 5, 3)) \\ = (9, 2, 1, 5, 3)$$

Définition 4.3 (Composition de lentilles). Soient $\mathcal{L}_1 \in \text{Lens}_{A,B}$ et $\mathcal{L}_2 \in \text{Lens}_{B,C}$.

La composition de \mathcal{L}_1 et \mathcal{L}_2 est la lentille $\mathcal{L} \in \text{Lens}_{A,C}$ définie de la manière suivante :

$$\text{get}_{\mathcal{L}}(r) = \text{get}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_1} r) \\ \text{put}_{\mathcal{L}}(a, r) = \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1} r), r)$$

On notera alors $\mathcal{L} = \mathcal{L}_1 \gg \mathcal{L}_2$.

Cette définition est illustrée dans la figure 4.3. Une preuve que la composition est une lentille est donnée en annexe D.1.

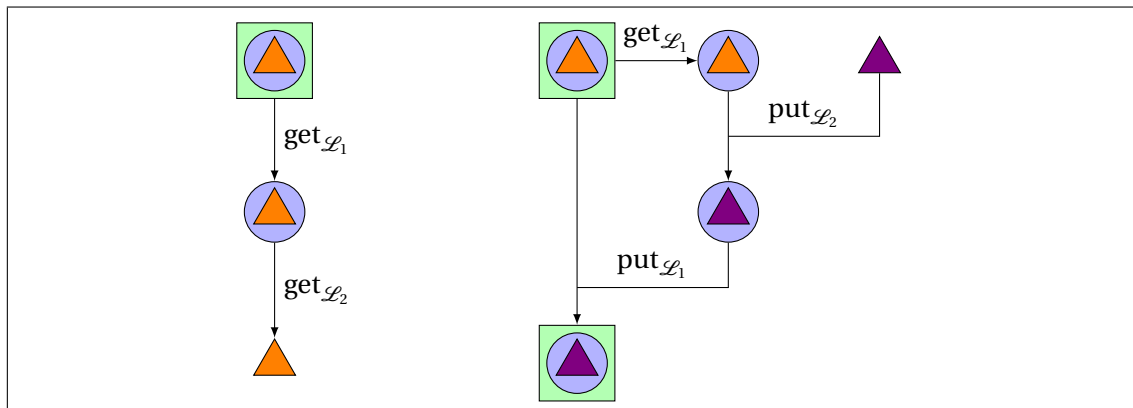


FIGURE 4.3: Composition de lentilles

4.2 Fonctionnalités

Fonctions et procédures : Un des problèmes classiques dans les langages impératifs est de distinguer les fonctions (qui retournent une valeur) et les procédures (qui n'en retournent pas). La solution choisie par C est de marquer les procédures comme retournant un “faux” type void. Mais c’est uniquement syntaxique : il n’est pas possible de manipuler cette valeur de retour de type void.

L’autre possibilité, souvent prise dans les langages fonctionnels, est de ne pas faire de distinction entre ces deux cas et d’interdire les procédures. Les fonctions ne retournant pas de valeur “intéressante” renvoient alors une valeur d’un type à un seul élément appelé $()^2$, et donc le type sera noté UNIT.

En C, puisqu’il n’y a pas de valeurs de type void, la notation void * a un sens particulier : elle désigne les pointeurs de type non défini, qui sont compatibles avec n’importe quel autre type de pointeur (c’est la seule forme — rudimentaire — de généricité qu’offre le langage). Ici, la valeur $()$ est une valeur comme les autres, et on peut construire un pointeur de type UNIT* qui n’aura pas de signification particulière : c’est seulement un pointeur vers une valeur de type UNIT.

Tableaux : ce sont des valeurs composées qui contiennent un certain nombre de valeurs d’un même type. Par exemple, 100 entiers. On accède à ces valeurs par un indice entier, qui dans le cas général n’est pas connu à la compilation. C’est une erreur (Ω_{array}) d’accéder à un tableau en dehors de ses bornes, c’est-à-dire en dehors de $[0; n - 1]$ pour un tableau à n éléments.

Les tableaux sont notés $[e_1; \dots; e_n]$, et le cas dégénéré ($n = 0$) est interdit.

Structures : comme les tableaux, ce sont des valeurs composées mais hétérogènes. Les différents éléments (appelés *champs*) sont désignés par noms l (pour *label*) et de manière statique (il n’y a pas de mécanisme pour faire référence à un nom dans le programme).

Les structures sont notées $\{l_1 : e_1; \dots; l_n : e_n\}$ et comme dans le cas des tableaux, le cas dégénéré ($n = 0$) est interdit.

Dans le programme, le nom de champ l est complété de la définition complète de la structure S . Celle-ci n’est pas utilisée dans l’évaluation et sera donc décrite au chapitre 5. Bien sûr, écrire la totalité de la structure dans le code du programme serait fastidieux. Nous

2. Cette notation évoque un n -uplet à 0 composante.

décrivons dans le chapitre 7 comment obtenir automatiquement ces annotations depuis un langage comme C qui utilise des noms de structures.

Types de données : très simples. Entiers machine, flottants, pointeurs et types composés (structures et tableaux) composés de ceux-ci.

Variables : elles sont mutables, et on peut passer des données par valeur ou par pointeur.

Flot de contrôle : il repose sur les constructions “if” et “while”. Les autres types de boucle (“for” et “do/while”) peuvent être construits avec ces opérateurs.

Fonctions : le code est organisé en fonctions “simples”, c’est-à-dire qui ne sont pas des fermetures. Même si le corps d’une fonction peut être inclus dans le corps d’une autre, il n’est pas possible d’accéder aux variables de la portée entourante depuis la fonction intérieure.

4.3 Syntaxe

Les figures 4.4, 4.5 et 4.6 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-mêmes 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 9.

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 4.6.

4.4 Mémoire et valeurs

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

La mémoire est constituée de variables, qui contiennent des valeurs. Ces variables sont organisées, d’une part en un ensemble de variables globales, et d’autre part en une pile de contextes d’appel (qu’on appellera donc aussi cadres de pile, ou *stack frames* en anglais). Cette structure empilée permet de représenter les différents contextes à chaque appel de fonction : par exemple, si une fonction s’appelle récursivement, plusieurs instances de ses variables locales sont présentes dans le programme.

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.

Pour identifier de manière non ambiguë une variable, on note simplement x la variable globale nommée x , et (n, x) la variable locale nommée x dans le n^{e} cadre de pile³.

3. Les paramètres de fonction sont traités comme des variables locales et se retrouvent dans le cadre correspondant.

Constantes	$c ::= n$	Entier
	d	Flottant
	NULL	Pointeur nul
	$()$	Valeur unité
Expressions	$e ::= c$	Constante
	lv	Accès mémoire
	$\Box e$	Opération unaire
	$e \boxplus e$	Opération binaire
	$\&lv$	Pointeur
	$lv \leftarrow e$	Affectation
	$\{l_1 : e_1; \dots; l_n : e_n\}$	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
	$lv.l_S$	Accès à un champ
	$lv[e]$	Accès à un élément
Fonctions	$f ::= \text{fun}(x_1, \dots, x_n)\{i\}$	Arguments, corps

FIGURE 4.4: Syntaxe des expressions

Instructions	$i ::= \text{PASS}$	Instruction vide
	$i; i$	Séquence
	e	Expression
	$\text{DECL } x = e \text{ IN } \{i\}$	Déclaration de variable
	$\text{IF}(e)\{i\}\text{ELSE}\{i\}$	Alternative
	$\text{WHILE}(e)\{i\}$	Boucle
	$\text{RETURN}(e)$	Retour de fonction
Phrases	$p ::= x = e$	Variable globale
	e	Évaluation d'expression
Programme	$P ::= (p_1, \dots, p_n)$	Phrases

FIGURE 4.5: Syntaxe des 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 4.6: Syntaxe des opérateurs

Les affectations peuvent avoir la forme $x \leftarrow e$ où x est une variable et e est une expression, mais pas seulement. En effet, à gauche de \leftarrow on trouve en général non pas une variable mais une left-value. Pour représenter quelle partie de la mémoire doit être accédée par cette left-value, on introduit la notion de chemin φ . Un chemin est en quelque sorte une left-value symbolique évaluée : les cas sont similaires, sauf que tous les indices sont évalués. Par exemple, $\varphi = (5, x).p$ représente le champ “ p ” de la variable x dans le 5^e cadre de pile. C’est à ce moment qu’on évalue les déréréfencements qui peuvent apparaître dans une left-value.

Les valeurs, quant à elles, peuvent avoir les formes suivantes (résumé sur la figure 4.7) :

- \hat{c} : une constante. La notation circonflexe permet de distinguer les constructions syntaxique des constructions sémantiques. Par exemple, à la syntaxe 3 correspond la valeur $\hat{3}$.

Les valeurs entières sont les entiers signés sur 32 bits, c’est-à-dire entre -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. Les flottants sont les flottants IEEE 754 de 32 bits [oEE08].

- φ : une référence mémoire. Ce chemin correspond à un pointeur sur une left-value. Par exemple, l’expression $\&x$ s’évalue en $\varphi = (5, x)$ si x désigne lexicalement une variable dans le 5^e cadre de pile.
- $\{l_1 : \widehat{v_1}; \dots; l_n : \widehat{v_n}\}$: une structure. Comme précédemment, on note $\{\cdot\}$ pour dénoter les valeurs.
- $[\widehat{v_1}; \dots; \widehat{v_n}]$: un tableau. Pareillement, $[\cdot]$ permet de désigner les valeurs. Par exemple, si x vaut 2 et y vaut 3, l’expression $[x, y]$ s’évaluera en valeur $[\widehat{2}, \widehat{3}]$
- \hat{f} : une fonction. Les valeurs fonctions comportent l’intégralité de la définition de la fonction (liste de paramètres, de variables locales et corps). Remarquons que contrairement à certains langages, l’environnement n’est pas capturé (il n’y a pas de clôture lexicale).
- Ω : une erreur. Par exemple le résultat de $5/0$ est Ω_{div} .

Les erreurs peuvent être de plusieurs types :

Valeurs	$v ::= \hat{c}$	Constante
	$\hat{\&} \varphi$	Référence mémoire
	$\{l_1 : \widehat{v_1}; \dots; l_n : \widehat{v_n}\}$	Structure
	$[\widehat{v_1}; \dots; \widehat{v_n}]$	Tableau
	\hat{f}	Fonction
	Ω	Erreur
Chemins	$\varphi ::= a$	Adresse
	$\varphi.l$	Accès à un champ
	$\varphi[n]$	Accès à un élément
Adresses	$a ::= (n, x)$	Variable locale
	(x)	Variable globale
Erreur	$\Omega ::= \Omega_{array}$	Débordement de tableau
	Ω_{ptr}	Erreur de pointeur
	Ω_{div}	Division par zéro
	Ω_{field}	Erreur de champ
	Ω_{var}	Variable inconnue

FIGURE 4.7: Valeurs

- Ω_{div} correspond à une division par zéro. Cela arrive si la seconde opérande de $/$, $\%$ ou $/.$ est nulle.
- Ω_{array} signale un débordement de tableau.
- Ω_{ptr} signale une manipulation ou un déréférencement de pointeur invalide. Par exemple, déréférencer le pointeur nul ou appliquer l'arithmétique en dehors d'un tableau provoque cette erreur.
- Ω_{field} arrive lorsqu'on accède à un champ inexistant d'une structure. Ce cas sera entièrement évité grâce au typage dans le chapitre 5.
- Ω_{var} se produit lorsqu'on accède à une variable qui n'est pas visible dans un état mémoire. Ce cas est également évité dans les programmes bien typés.

La figure 4.8 résume comment ces valeurs sont organisées. Une pile est une liste de cadre de piles, et un cadre de pile est une liste de couples (nom, valeur). Un état mémoire m est un couple (s, g) où s est une pile et g un cadre de pile (qui représente les variables globales). On note $|m| = |s|$ la hauteur de la pile (en nombre de cadres).

Enfin, l'interprétation est définie comme une relation $\cdot \rightarrow \cdot$ entre états Ξ ; ces états sont d'une des formes suivantes :

- un couple $\langle e, m \rangle$ où e est une expression et m un état mémoire. m est l'état mémoire sous lequel l'évaluation sera réalisée. Par exemple $\langle 3, ([], [x \mapsto 3]) \rangle \rightarrow \langle \hat{3}, ([], [x \mapsto 3]) \rangle$. L'évaluation des expressions est détaillée dans la section 4.10.

Pile	$s ::= []$ $\{x_1 \mapsto v_1; \dots; x_n \mapsto v_n\} :: s$	Pile vide Ajout d'un cadre
État mémoire	$m ::= (s, \{x_1 \mapsto v_1; \dots; x_n \mapsto v_n\})$	Pile, globales
État d'inter-préteur	$\Xi ::= \langle e, m \rangle$ $\langle i, m \rangle$ Ω	Expression, mémoire Instruction, mémoire Erreur

FIGURE 4.8: Composantes d'un état mémoire

- un couple $\langle i, m \rangle$ où i est une instruction et m un état mémoire. La réduction des instructions est traitée dans la section 4.11.
- un couple $\langle lv, m \rangle$ où lv est une left-value et m un état mémoire. L'évaluation des left-values est décrite en section 4.9.
- une erreur Ω . La propagation des erreurs est détaillée dans la section 4.12.

L'évaluation des expressions, left-values et instruction se fait à petits pas. C'est-à-dire qu'on simplifie d'étape en étape leur forme, jusqu'à arriver à un cas de base :

- pour les expressions, une valeur v .
- pour les instructions, l'instruction PASS ou RETURN(v) où v est une valeur.
- pour les left-values, un chemin φ .

On considère en fait la clôture transitive de cette relation. Cela revient à ajouter une règle :

$$\frac{\Xi_1 \rightarrow \Xi_2 \quad \Xi_2 \rightarrow \Xi_3}{\Xi_1 \rightarrow \Xi_3} \text{ (TRANS)}$$

4.5 Opérations sur les valeurs

Un certain nombre d'opérations est possibles sur les valeurs (figure 4.6) :

- les opérations arithmétiques $+$, $-$, \times , $/$ et $\%$ sur les entiers. L'opérateur $\%$ correspond au modulo (reste de la division euclidienne). En cas de division par zéro, l'erreur Ω_{div} est levée.
- les versions “pointées” $+$., $-$., \times ., et $/$., sur les flottants
- les opérations d'arithmétique de pointeur $+_p$ et $-_p$ qui à un chemin mémoire et un entier associent un chemin mémoire.
- les opérations d'égalité $=$ et \neq . L'égalité entre entiers ou entre flottants est immédiate. Deux valeurs composées (tableaux ou structures) sont égales si elles ont la même “forme” (même taille pour les tableaux, et même champs pour les structures) et que toutes leurs sous-valeurs sont égales deux à deux. Les références mémoire φ sont égales lorsque les chemins qu'ils décrivent sont syntaxiquement égaux⁴.

4. Il est donc possible que deux pointeurs pointent sur la même adresse mais soient considérés différents. La raison pour ce choix est que la comparaison doit pouvoir se faire sans accéder à la mémoire.

- les opérations de comparaison $\leq, \geq, <, >$ sont définies avec leur sémantique habituelle sur les entiers et les flottants. Sur les références mémoires, elles sont définies dans le cas où les deux opérarandes sont de la forme $\varphi[\cdot]$ par : $\varphi[n] \boxplus \varphi[m] \stackrel{\text{def}}{=} n \boxplus m$. Dans les autres cas, l'erreur Ω_{ptr} est renvoyée.
- les opérateurs bit à bit sont définis sur les entiers. $\&$, $|$ et \wedge représentent respectivement la conjonction, la disjonction et la disjonction exclusive (XOR).
- des versions logiques de la conjonction ($\&\&$) et de la disjonction ($||$) sont également présentes. Leur sémantique est donnée par le tableau suivant :

n	m	$n \&\& m$	$n m$
0	0	0	0
0	$\neq 0$	0	1
$\neq 0$	0	0	1
$\neq 0$	$\neq 0$	1	1

- des opérateurs de décalage à gauche (\ll) et à droite (\gg) sont présents. Eux aussi ne s'appliquent qu'aux entiers.
- les opérateurs arithmétiques unaires $+$, $-$, $+$. et $-$. sont équivalents à l'opération binaire correspondante avec 0 ou 0.0 comme première opérande.
- \sim inverse tous les bits de son opérande. $!$ est une version logique, c'est-à-dire que $!0 = 1$ et si $n \neq 0$, $!n = 0$.

4.6 Opérations sur les états mémoire

Définition 4.4 (Recherche de variable). *La recherche de variable permet d'associer à une variable x une adresse a .*

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

$$\text{Lookup}((s, g), x) = \begin{cases} (|s|, x) & \text{si } |s| > 0 \text{ et } (x \mapsto v) \in s_1 \\ x & \text{si } (x \mapsto v) \in g \\ \Omega_{var} & \text{sinon} \end{cases}$$

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 4.5 (Manipulations de pile). *On définit l'empilement d'un cadre de pile $c = ((x_1 \mapsto v_1), \dots, (x_n \mapsto v_n))$ sur un état mémoire $m = (s, g)$ (figure 4.9(a)) :*

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

On définit aussi l'extension du dernier cadre de pile (figure 4.9(b)) :

$$\text{Extend}((c :: s, g), x \mapsto v) = (((x \mapsto v) :: c) :: s, g)$$

De même on définit le dépilement (figure 4.9(c)) :

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

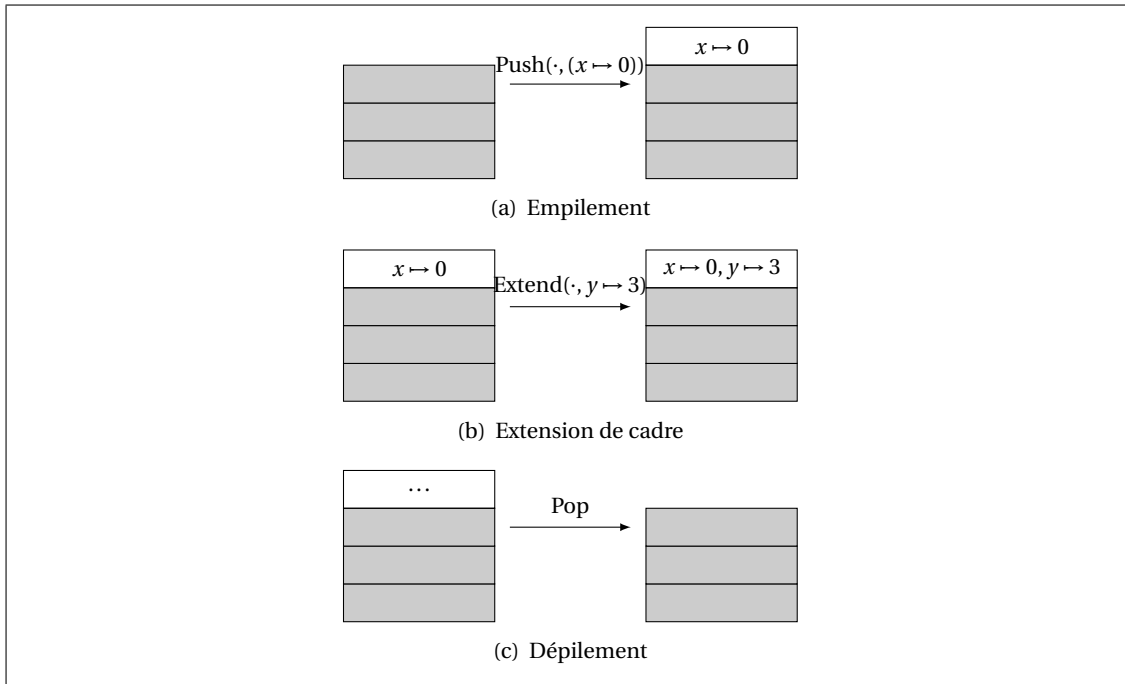


FIGURE 4.9: Opérations de pile

Les opérations Extend et Pop ne sont définies que pour une pile non vides. Néanmoins cela ne pose pas de problème, puisque lors de l'exécution, la pile n'est vide que lors de l'évaluation d'expressions dans les phrases de programme. À cet endroit, seules des expressions peuvent apparaître, et leur évaluation ne manipule jamais la pile avec ces opérations.

On définit aussi une opération de nettoyage de pile, qui sera utile pour les retours de fonction.

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 :

$$\begin{aligned}
 \text{Cleanup}(s, g) &= (s', g') \\
 \text{où } g' &= \text{CleanupList}(|s|, g) \\
 s' &= [\text{CleanupList}(|s|, s_1), \dots, \text{CleanupList}(|s|, s_n)] \\
 \text{CleanupList}(p, u) &= \{(x, v) \in u / v \text{ n'est pas une adresse}\} \\
 &\quad \cup \{(x, \varphi) \in u / \text{Live}(p, \varphi)\} \\
 \text{Live}(p, (n, x)) &= n < p \\
 \text{Live}(p, (x)) &= \text{Vrai} \\
 \text{Live}(p, \varphi.l) &= \text{Live}(p, \varphi) \\
 \text{Live}(p, \varphi[n]) &= \text{Live}(p, \varphi)
 \end{aligned}$$

Sans cette règle, examinons le programme suivant :

f ()	g (p)	h ()
(x=0)	(x=0.0)	(p=f ())
{	{	{
return (&x);	*p = 1;	g(p);
}	}	}

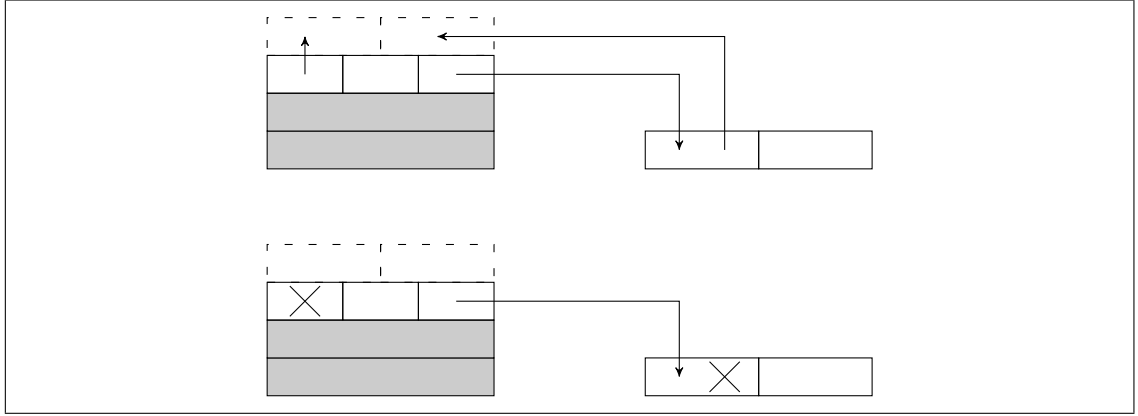


FIGURE 4.10: Nettoyage d'un cadre de pile

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

4.7 Accesseurs

On définit ici quelques lentilles. Le but est de définir rigoureusement les accès à la mémoire. À partir d'un état mémoire m et d'une left-value φ , on veut pouvoir obtenir :

- la valeur accessible au chemin $\varphi : m[\varphi]_\Phi$
- l'état mémoire obtenu en remplaçant celle-ci par une nouvelle valeur $v' : m[\varphi \leftarrow v']_\Phi$

Pour définir cette lentille indexée Φ , on commence par définir des lentilles élémentaires, et on les compose pour pouvoir définir des lentilles entre valeurs.

Accès à une liste par indice : I

On définit une lentille indexée $I : \mathbb{N} \rightarrow \text{LENS}_{\text{List}(\alpha), \alpha}$ permettant d'accéder aux éléments d'une liste par leur indice.

$$\begin{aligned}
 l[n]_I &= l_n \text{ si } n \in [1; |l|] \\
 l[n \leftarrow x]_I &= l' \\
 &\text{où } l'_n = x \\
 &\quad \forall i \neq n, l'_i = l_i
 \end{aligned}$$

I n'est définie que pour $n \in [1; |l|]$.

Accès à une liste d'associations : L

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.

$$l[x]_L = \begin{cases} v & \text{si } \exists! n \in [1; |l|]. \exists v. l_n = (x \mapsto v) \\ \Omega_{var} & \text{sinon} \end{cases}$$

$$l[x \leftarrow v]_L = \begin{cases} l[n \leftarrow (x \mapsto v)]_I & \text{si } \exists! n \in [1; |l|]. \exists v_0. l_n = (x \mapsto v_0) \\ \Omega_{var} & \text{sinon} \end{cases}$$

Accès par adresse : A

Les états mémoire sont constitués des listes d'association (nom, valeur).

L'accessneur 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.

On pose $m = (s, g)$.

Les accès aux globales se font de la manière suivante. Si la variable n'existe pas, notons que L retourne Ω_{var} .

$$m[x]_A = g[x]_L$$

$$m[x \leftarrow v]_A = (s, g[x \leftarrow v]_L)$$

Les accès aux locales reviennent à accéder à la bonne variable du bon cadre de pile. Cela revient naturellement à composer les lentilles L et I . On définit donc une lentille $\mathcal{L}_{n,x} = I(|s| - n + 1) \gg L(x)$ qui accède à la variable x du n^e cadre de pile.

$$m[(n, x)]_A = \begin{cases} \text{get}_{\mathcal{L}_{n,x}}(s) & \text{si } n \in [1; |s|] \\ \Omega_{var} & \text{sinon} \end{cases}$$

$$m[(n, x) \leftarrow v]_A = \begin{cases} (\text{put}_{\mathcal{L}_{n,x}}(v, s), g) & \text{si } n \in [1; |s|] \\ \Omega_{var} & \text{sinon} \end{cases}$$

Les numéros de cadre qui permettent d'identifier les locales (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 $|s| - n + 1$ apparaît dans la définition précédente.

Accès par champ : F

Les valeurs qui sont des structures possèdent des sous-valeurs, associées à des noms de champ.

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

L'erreur Ω_{field} est levée si on accède à un champ non existant.

$$\begin{aligned}
\{l_1 : v_1; \dots; l_n : v_n\}[l]_F &= v_i \text{ si } \exists i \in [1; n], l = l_i \\
\{l_1 : v_1; \dots; l_n : v_n\}[l \leftarrow v]_F &= \{l_1 : v_1 \\
&\quad ; \dots \\
&\quad ; l_{p-1} : v_{p-1} \\
&\quad ; l_p : v \\
&\quad ; l_{p+1} : v_{p+1} \\
&\quad ; \dots \\
&\quad ; l_n : v_n\} \text{ si } \exists p \in [1; n], l = l_p \\
\{l_1 : v_1; \dots; l_n : v_n\}[l]_F &= \Omega_{field} \text{ sinon} \\
\{l_1 : v_1; \dots; l_n : v_n\}[l \leftarrow v]_F &= \Omega_{field} \text{ sinon}
\end{aligned}$$

Accès par indice de tableau : T

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

$$\begin{aligned}
[v_1; \dots; v_n][i]_T &= v_{i+1} \text{ si } i \in [0; n-1] \\
[v_1; \dots; v_n][i]_T &= \Omega_{array} \text{ sinon} \\
[v_1; \dots; v_n][i \leftarrow v]_T &= [v'_1; \dots; v'_n] \text{ si } i \in [0; n-1] \\
&\quad \text{où } \begin{cases} v'_i = v \\ \forall j \neq i, v'_j = v_j \end{cases} \\
[v_1; \dots; v_n][i \leftarrow v]_T &= \Omega_{array} \text{ sinon}
\end{aligned}$$

Accès par chemin : Φ

L'accès par chemin Φ permet de lire et de modifier la mémoire en profondeur. On peut accéder directement à une variable :

$$\Phi(a) = A(a)$$

Les accès à des sous-valeurs se font en composant les accesseurs (définition 4.3) :

$$\begin{aligned}
\Phi(\varphi.l) &= \Phi(\varphi) \ggg F(l) \\
\Phi(\varphi[i]) &= \Phi(\varphi) \ggg T(i)
\end{aligned}$$

Enfin, l'accès à la mémoire par le pointeur nul provoque une erreur :

$$\begin{aligned}
m[\text{NULL}]_\Phi &= \Omega_{ptr} \\
m[\text{NULL} \leftarrow v]_\Phi &= \Omega_{ptr}
\end{aligned}$$

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

En revanche, la modification est plus complexe. Le premier cas ($\varphi = a$) modifie directement une valeur complète (en modifiant une association), mais les deux suivants ($\varphi = \varphi'.l$ 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 $[l \leftarrow v]_F$)
- affecter cette valeur au même chemin (le $m[\varphi \leftarrow \dots]_\Phi$ externe)

4.8 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 contexte
- 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 4.11).

Le deuxième est résolu par la règle d'inférence suivante.

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

Enfin, le troisième revient à définir l'opérateur de substitution $\cdot(\cdot)$ présent dans la règle précédente. Afin de pouvoir appliquer des substitution au niveau des left-values et des instructions, on définit aussi respectivement $\cdot(\cdot)_L$ et $\cdot(\cdot)_I$.

Notons que puisque $i ::= e$ et $e ::= lv$, on peut aussi l'appliquer aux expressions et aux left-values : l'opération $\cdot(\cdot)$ est purement syntaxique.

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 : $C(X)$ est l'objet syntaxique (instruction, expression ou left-value) obtenu en remplaçant l'unique \bullet dans C par X . Par exemple, $\text{DECL } x = 2 + \bullet \text{ IN}\{\{\}5\} = \text{DECL } x = 2 + 5 \text{ IN}\{\}$

À titre d'illustration, décomposons l'évaluation de $e_1 \boxplus e_2$ en $v = v_1 \boxplus v_2$ depuis un état mémoire m :

1. on commence par évaluer, d'une manière ou d'une autre, l'expression e_1 en une valeur v_1 . Le nouvel état mémoire est noté m' . Soit donc $\langle e_1, m \rangle \rightarrow \langle v_1, m' \rangle$.
2. En appliquant la règle CTX avec $C = \bullet \boxplus e_2$ (qui est une des formes possibles pour un contexte d'évaluation), on déduit de 1. que $\langle e_1 \boxplus e_2, m \rangle \rightarrow \langle v_1 \boxplus e_2, m' \rangle$
3. D'autre part, on évalue e_2 depuis m' . En supposant encore que l'évaluation converge, notons v_2 la valeur calculée et m'' l'état mémoire résultant : $\langle e_2, m' \rangle \rightarrow \langle v_2, m'' \rangle$.

Contextes	$C ::= \bullet$
	$C \boxplus e$
	$v \boxplus C$
	$\boxplus C$
	$C \leftarrow e$
	$\varphi \leftarrow C$
	$\{l_1 : v_1; \dots; l_i : C; \dots; l_n : e_n\}$
	$[v_1; \dots; C; \dots; e_n]$
	$C(e_1, \dots, e_n)$
	$f(v_1, \dots, C, \dots, e_n)$
	$C.l_S$
	$C[e]$
	$\varphi[C]$
	$C; i$
	$\text{IF}(C)\{i_1\}\text{ELSE}\{i_2\}$
	$\text{RETURN}(C)$
	$\text{DECL } x = C \text{ IN}\{i\}$

FIGURE 4.11: Contextes d'évaluation

x	$.l_S$	$[2 * n]$	$\cdot g$
\downarrow	\downarrow	\downarrow	\downarrow
$(2, x)$	$.l$	$[4]$	$\cdot g$

FIGURE 4.12: Évaluation des left-values.

4. Appliquons la règle CTX à 3. avec $C = v_1 \boxplus \bullet$. On obtient $\langle v_1 \boxplus e_2, m \rangle \rightarrow \langle v_1 \boxplus v_2, m' \rangle$.
5. En combinant les résultats de 2. et 4. on en déduit que $\langle e_1 \boxplus e_2, m \rangle \rightarrow \langle v_1 \boxplus v_2, m'' \rangle$.
6. D'après la règle EXP-BINOP, $\langle v_1 \boxplus v_2, m'' \rangle \rightarrow \langle v_1 \hat{\boxplus} v_2, m'' \rangle$
7. D'après 5. et 6., on a par combinaison $\langle e_1 \boxplus e_2, m \rangle \rightarrow \langle v, m'' \rangle$ en posant $v = v_1 \hat{\boxplus} v_2$.

4.9 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 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. Dans le cas des champs de structure, la définition de la structure S n'est pas prise en compte.

```

Decl s0 = { .f = 0 } in
Decl s1 = { .f = 1 } in
Decl x  = & s0 in
Decl p  = & ((*x).f) in
/* (a) */
x <- & s1
/* (a) */

```

FIGURE 4.13: Évaluation stricte ou paresseuse des left-values

$$\frac{}{\langle lv.l_S, m \rangle \rightarrow \langle lv.\widehat{l}, m \rangle} \text{ (PHI-STRUCT)}$$

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 .

La règle portant sur le déréférencement est particulière. On peut penser que la bonne définition de φ consiste à se calquer sur la définition de lv , en remplaçant les noms de variable par leur adresse résolue et en évaluant les indices de tableau, et à ajouter une règle qui transforme $*\varphi$ en $\widehat{*}\varphi$.

Cela ne fonctionne pas, car alors les déréférencements sont évalués trop tard : au moment de l'affectation dans la left-value plutôt qu'à sa définition. La figure 4.13 illustre ce problème. On s'intéresse à l'évaluation de l'expression $*p$ aux points (a) et (b). Avec une sémantique paresseuse (en ajoutant un $\widehat{*}\varphi$), la valeur de p est $\widehat{*}((0, x)).f$, donc $*p$ est évalué à 0 en (a) et 1 en (b). Au contraire, avec une sémantique stricte (correcte), p vaut $\widehat{*}(((0, s0)).f)$ et donc $*p$ est évalué à 0 en (a) et en (b). La règle d'évaluation est donc la suivante.

$$\frac{\widehat{*}\varphi' = m[\varphi]_\Phi}{\langle *\varphi, m \rangle \rightarrow \langle \varphi', m \rangle} \text{ (PHI-DEREF)}$$

Un exemple d'évaluation est donné dans la figure 4.12.

4.10 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 est déjà évaluée :

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

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

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

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 $\widehat{\boxplus}$ qui agit sur les valeurs. Par exemple, l'opérateur $\widehat{+}$ est l'addition classique entre entiers. Comme précisé dans la section 4.5, la division par zéro via $/$, $\%$ ou $/.$ provoque l'erreur Ω_{div} .

$$\frac{}{\langle \boxminus v, m \rangle \rightarrow \langle \widehat{\boxminus} v, m \rangle} \text{ (EXP-UNOP)} \quad \frac{}{\langle v_1 \boxplus v_2, m \rangle \rightarrow \langle v_1 \widehat{\boxplus} v_2, m \rangle} \text{ (EXP-BINOP)}$$

Il est nécessaire de dire un mot sur les opérations $\widehat{+}_p$ et $\widehat{-}_p$ 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 :

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

Cela implique qu'on ne peut pas faire faire d'arithmétique de pointeurs au sein d'une même structure. Autrement c'est une erreur de manipulation de pointeurs ⁵ :

$$\begin{aligned} \varphi +_p m &= \Omega_{ptr} \text{ si } \nexists(\varphi', n), \varphi = \varphi'[n] \\ \varphi -_p m &= \Omega_{ptr} \text{ si } \nexists(\varphi', n), \varphi = \varphi'[n] \\ \text{NULL} +_p m &= \Omega_{ptr} \\ \text{NULL} -_p m &= \Omega_{ptr} \end{aligned}$$

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

$$\frac{}{\langle \& \varphi, m \rangle \rightarrow \langle \widehat{\&} \varphi, m \rangle} \text{ (EXP-ADDR OF)}$$

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.

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

5. Cela est cohérent avec la norme C99 : "If the pointer operand points to an element of an array object, and the array is large enough, [...]; otherwise, the behavior is undefined." [ISO99, 6.5.6 §8]

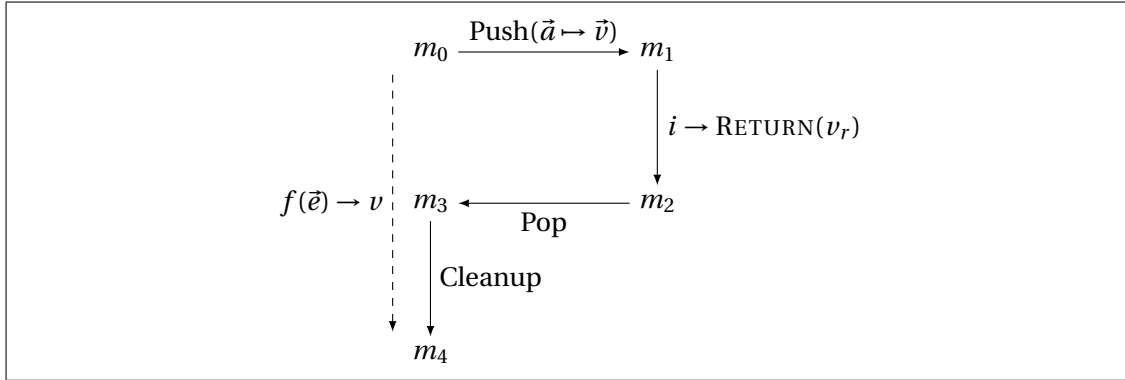


FIGURE 4.14: Appel d'une fonction. La taille de la pile croît de gauche à droite, et les réductions se font de haut en bas.

Expressions composées

Les littéraux de structures sont évalués en leurs constructions syntaxiques respectives. Puisque les contextes d'évaluation sont de la forme $[v_1; \dots; C; \dots; e_n]$, l'évaluation se fait toujours de gauche à droite.

$$\frac{}{\langle \{l_1 : v_1; \dots; l_n : v_n\}, m \rangle \rightarrow \langle \{l_1 : \overline{v_1}; \dots; l_n : \overline{v_n}\}, m \rangle} \text{ (EXP-STRUCT)}$$

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

L'appel de fonction est traité de manière similaire à la déclaration de fonction : on ne peut pas facilement relier un pas d'évaluation de i à un pas d'évaluation de $\text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n)$, donc un contexte $C ::= \text{fun}(a_1, \dots, a_n)\{\bullet\}(v_1, \dots, v_n)$ n'est pas à considérer. Il faut donc transmettre le pas à la main. C'est le but des deux règles EXP-CALL-CTX et EXP-CALL-ERR. D'autre part, les cas terminaux sont EXP-CALL-RETURN-UNIT et EXP-CALL-RETURN :

$$\frac{\begin{array}{l} m_1 = \text{Push}(m_0, (a_1 \mapsto v_1, \dots, a_n \mapsto v_n)) \\ \langle i, m_1 \rangle \rightarrow \langle i', m_2 \rangle \quad \forall i \in [1; n], v'_i = m_2[(0, a_i)]_A \quad m_3 = \text{Pop}(m_2) \end{array}}{\langle \text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n), m_0 \rangle \rightarrow \langle \text{fun}(a_1, \dots, a_n)\{i'\}(v'_1, \dots, v'_n), m_3 \rangle} \text{ (EXP-CALL-CTX)}$$

$$\frac{m' = \text{Push}(m, (a_1 \mapsto v_1, \dots, a_n \mapsto v_n)) \quad \langle i, m' \rangle \rightarrow \Omega}{\langle \text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n), m \rangle \rightarrow \Omega} \text{ (EXP-CALL-ERR)}$$

$$\frac{m' = \text{Cleanup}(m)}{\langle \text{fun}(a_1, \dots, a_n)\{\text{PASS}\}(v_1, \dots, v_n), m \rangle \rightarrow \langle (), m' \rangle} \text{ (EXP-CALL-RETURN-UNIT)}$$

$$\frac{m' = \text{Cleanup}(m)}{\langle \text{fun}(a_1, \dots, a_n)\{\text{RETURN}(v)\}(v_1, \dots, v_n), m \rangle \rightarrow \langle v, m' \rangle} \text{ (EXP-CALL-RETURN)}$$

Cette évaluation est décrite dans la figure 4.14.

4.11 Instructions

Définition 4.6. *Évaluation d'une instruction*

Comme pour 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.

$$\frac{}{\langle (\text{PASS}; i), m \rangle \rightarrow \langle i, m \rangle} \text{ (SEQ)} \qquad \frac{}{\langle v, m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{ (EXP)}$$

L'évaluation de $\text{DECL } x = v \text{ IN } \{i\}$ sous m peut se faire de plusieurs manières. La règle principale est DECL-CTX qui relie un pas d'évaluation sous une déclaration à un pas d'évaluation externe : pour ce faire, on étend l'état mémoire en ajoutant x , on effectue le pas, puis on enlève x . L'instruction résultante est la déclaration de x avec la nouvelle valeur v' de x après le pas d'exécution⁶.

Si i est terminale (PASS ou RETURN(v)), alors on peut s'évaluer en i en nettoyant l'espace mémoire des références à x qui peuvent subsister.

Enfin, si une erreur se produit elle est propagée.

$$\begin{aligned} &\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{PASS}\}, m \rangle \rightarrow \langle \text{PASS}, m' \rangle} \text{ (DECL-PASS)} \\ &\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{RETURN}(v')\}, m \rangle \rightarrow \langle \text{RETURN}(v'), m' \rangle} \text{ (DECL-RETURN)} \\ &\frac{\begin{array}{c} m' = \text{Extend}(m, x \mapsto v) \\ \langle i, m' \rangle \rightarrow \langle i', m'' \rangle \quad v' = m''[(0, x)]_A \quad m''' = m'' - x \end{array}}{\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \langle \text{DECL } x = v' \text{ IN } \{i'\}, m''' \rangle} \text{ (DECL-CTX)} \\ &\frac{\langle i, m \rangle \rightarrow \Omega}{\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \Omega} \text{ (DECL-ERR)} \end{aligned}$$

6. On peut remarquer qu'il est impossible de définir un contexte d'évaluation $C ::= \text{DECL } x = v \text{ IN } \{C\}$. En effet, puisque celui-ci nécessiterait d'ajouter une variable, il ne préserve pas la mémoire.

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.

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

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.

$$\begin{array}{lcl} \text{while}(e) \{ & & \text{if}(e) \{ \\ \quad i & \cong & \quad i; \\ \} & & \quad \text{if}(e) \{ \\ & & \quad \quad i; \\ & & \quad \quad \text{if}(e) \{ \\ & & \quad \quad \quad i; \\ & & \quad \quad \quad \dots \\ & & \quad \quad \} \\ & & \quad \} \\ & & \} \end{array}$$

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

$$\begin{array}{lcl} \text{while}(e) \{ & & \text{if}(e) \{ \\ \quad i & \cong & \quad i; \\ \} & & \quad \text{while}(e) \{ \\ & & \quad \quad i \\ & & \quad \} \\ & & \} \end{array}$$

Enfin, si un RETURN(-) apparaît dans une séquence, on peut supprimer la suite :

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

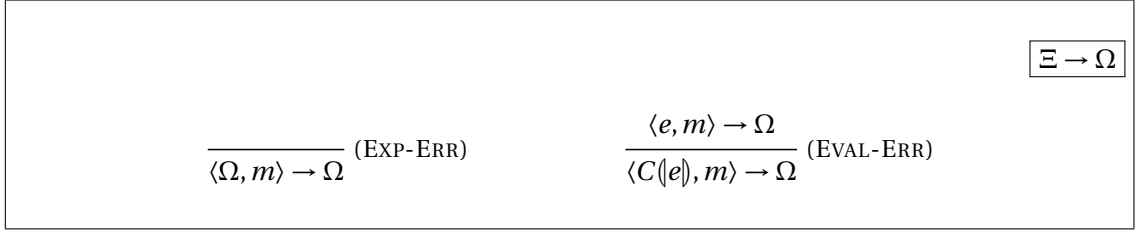


FIGURE 4.15: Cas d'erreur lors d'une évaluation

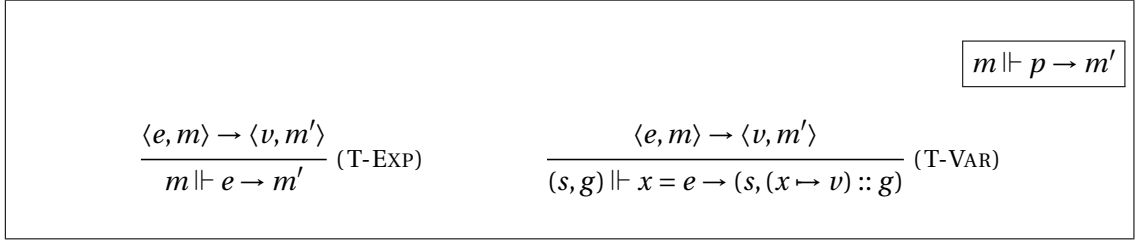


FIGURE 4.16: Évaluation des phrases d'un programme

4.12 Erreurs

Les erreurs se propagent des données vers l'interprète ; c'est-à-dire que si une expression ou instruction est réduite en une valeur d'erreur Ω , alors une transition est faite vers cet état d'erreur.

Cela est aussi vrai d'une sous-expression ou sous-instruction : si l'évaluation de e_1 provoque une erreur, l'évaluation de $e_1 + e_2$ également. La notion de sous-expression ou sous instruction est définie en fonction des contextes C (figure 4.15).

4.13 Phrases

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

Donc l'évaluation d'une phrase p fait passer d'un état mémoire m à un autre m' , ce que l'on note $m \Vdash p \rightarrow m'$.

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()`. La déclaration d'une variable globale (avec un initialiseur), quant à elle consiste à évaluer cet initialiseur et à étendre l'état mémoire avec ce couple (variable, valeur) (figure 4.16).

Notons que ces évaluations se font à grands pas.

4.14 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 :

$$\frac{}{m \Vdash [] \rightarrow^* m} \text{ (T*-NIL)} \qquad \frac{m \Vdash p \rightarrow m' \quad m' \Vdash ps \rightarrow^* m''}{m \Vdash p :: ps \rightarrow^* m''} \text{ (T*-CONS)}$$

L'exécution d'un programme est définie par :

$$\frac{([], []) \Vdash P \rightarrow^* m}{\Vdash P \rightarrow^* m} \text{ (PROG)}$$

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 5.1(a) des programmes mal formés comme celui de la figure 5.1(b).

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 5.1(a), la variable x est initialisée avec la valeur 0, c'est donc un entier. Cela signifie que dans 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 `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 .

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.

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

FIGURE 5.1: Programmes bien et mal formés

Type	$t ::= \text{INT}$	Entier
	FLOAT	Flottant
	UNIT	Unité
	t^*	Pointeur
	$t[]$	Tableau
	S	Structure
	$(t_1, \dots, t_n) \rightarrow t$	Fonction
Structure	$S ::= \{l_1 : t_1; \dots; l_n : t_n\}$	Structure simple
Environnement de typage	$\Gamma ::= []$	Environnement vide
	$(a, t) :: \Gamma'$	Extension

FIGURE 5.2: Types et environnements de typage

Ceci est valable pour tout t (INT, FLOAT ou 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 Environnements et notations

Les types associés aux expressions sont décrits dans la figure 5.2. Tous sont des types concrets : il n'y a pas de polymorphisme.

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), c'est à dire une suite finie $(\gamma_i)_{i \in [1; n]}$ avec $\gamma_i = (x_i, t_i)$. On utilise les notations suivantes :

$$\begin{aligned}
 x : t \in \Gamma &\stackrel{\text{def}}{=} \exists i \in [1; n], \gamma_i = (x, t) \\
 \text{dom}(\Gamma) &\stackrel{\text{def}}{=} \{x_i / i \in [1; n]\} \\
 \Gamma, x : t &\stackrel{\text{def}}{=} (\gamma'_i)_{i \in [1; n+1]} \text{ tel que } \begin{cases} \forall i \in [1; n], \gamma'_i = \gamma_i \\ \gamma'_{n+1} = (x, t) \end{cases} \\
 \Gamma - x &\stackrel{\text{def}}{=} (\gamma'_i)_{i \in [1; n-1]} \text{ tel que } \begin{cases} \forall i \in [1; p-1], \gamma'_i = \gamma_i \\ \forall i \in [p; n-1], \gamma'_i = \gamma_{i+1} \end{cases} \quad \text{où } \gamma_p = (x, t_p)
 \end{aligned}$$

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

Typage d'une expression : on note de la manière suivante le fait qu'une expression e (telle que définie dans la figure 4.4) ait pour type t dans le contexte Γ .

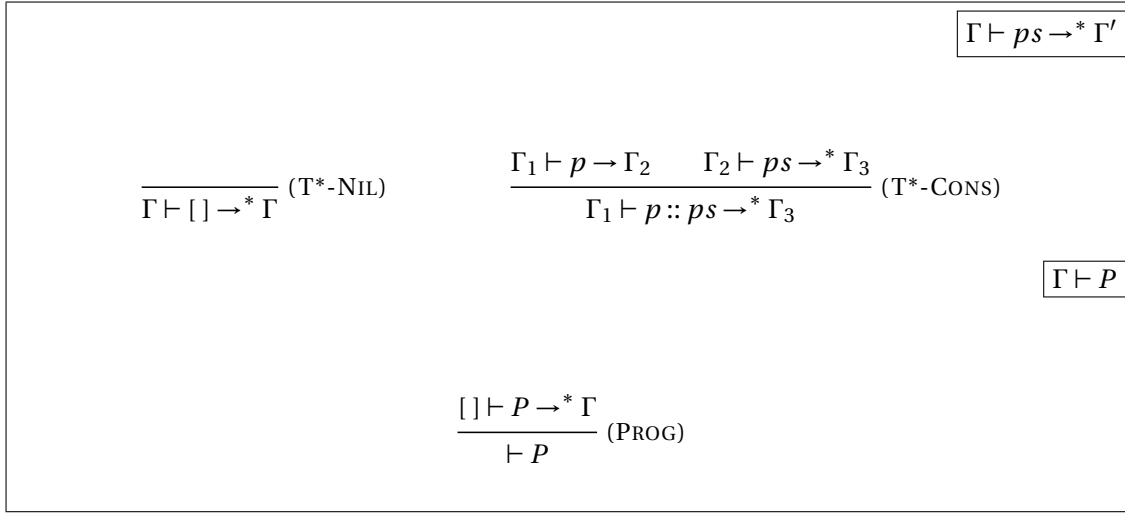


FIGURE 5.3: Typage d'une suite de phrases et d'un programme

$$\Gamma \vdash e : t$$

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 :

$$\Gamma \vdash i$$

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, ce qui définit le typage d'un programme, ce que l'on note $\vdash P$ (figure 5.3).

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.

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

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

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

Enfin, le littéral unité a le type UNIT.

$$\frac{}{\Gamma \vdash () : \text{UNIT}} \text{ (CST-UNIT)}$$

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 $t[]$).

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

Le typage de l'accès à un champ est facilité par le fait que dans le programme, le type complet de la structure est accessible sur le champ.

Dans la définition de cette règle on utilise la notation :

$$(l, t) \in \{l_1 : t_1; \dots; l_n : t_n\} \stackrel{\text{def}}{=} \exists i \in [1; n], l = l_i \wedge t = t_i$$

$$\frac{(l, t) \in S \quad \Gamma \vdash lv : S}{\Gamma \vdash lv.l_S : t} \text{ (LV-FIELD)}$$

Opérateurs

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

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

$\text{EQ}(t)$		
$\frac{t \in \{\text{INT}, \text{FLOAT}\}}{\text{EQ}(t)} \text{ (EQ-NUM)}$	$\frac{}{\text{EQ}(t*)} \text{ (EQ-PTR)}$	$\frac{\text{EQ}(t)}{\text{EQ}(t[])} \text{ (EQ-ARRAY)}$
$\frac{\forall i \in [1; n]. \text{EQ}(t_i)}{\text{EQ}(\{l_1 : t_1; \dots l_n : t_n\})} \text{ (EQ-STRUCT)}$		

FIGURE 5.4: Jugements d'égalité sur les types

De même sur `Float`.

$$\frac{\boxplus \in \{+, -, \times, /, \leq, \geq, <, >\} \quad \Gamma \vdash e_1 : \text{Float} \quad \Gamma \vdash e_2 : \text{Float}}{\Gamma \vdash e_1 \boxplus e_2 : \text{Float}} \text{ (OP-FLOAT)}$$

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, ainsi que les types composés si les types de leurs composantes (figure 5.4). Les opérateurs `=` et `≠` renvoient alors un `INT` :

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

Les opérateurs unaires “+” et “−” appliquent aux entiers, et leurs équivalents “+.” et “−.” aux flottants.

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

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

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

L'arithmétique de pointeurs préserve le type des pointeurs.

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

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.

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

Un littéral de structure est bien typé si ses champs sont bien typés.

$$\frac{\forall i \in [1; n], \Gamma \vdash e_i : t_i}{\Gamma \vdash \{l_1 : e_1; \dots; l_n : e_n\} : \{l_1 : t_1; \dots; l_n : t_n\}} \text{ (STRUCT)}$$

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.

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

Une déclaration de variable est bien typée si son bloc interne est bien typé quand on ajoute à l'environnement la variable avec le type de son initialiseur.

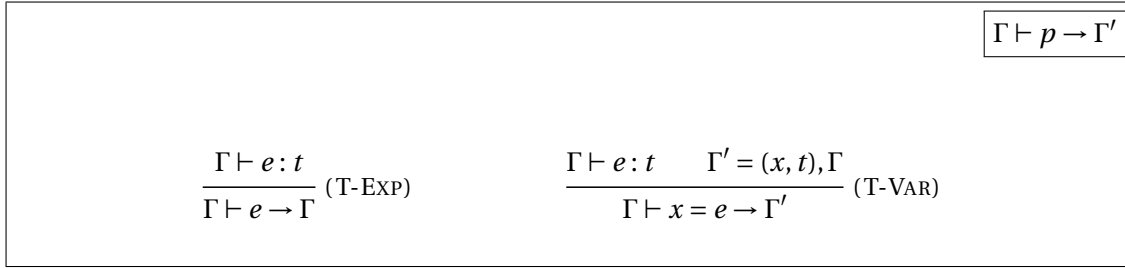


FIGURE 5.5: Typage des phrases

$$\frac{\Gamma \vdash e : t \quad \Gamma, x : t \vdash i}{\Gamma \vdash \text{DECL } x = e \text{ IN } \{i\}} \text{ (DECL)}$$

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 le type de la valeur de retour de la fonction appelée, $\underline{R} : t$

Si le corps de la fonction est bien typé sous Γ' , alors la fonction est typable en $(t_1, \dots, t_n) \rightarrow t$ sous Γ .

$$\frac{\Gamma' = (\Gamma - \underline{R}), a_1 : t_1, \dots, a_n : t_n, \underline{R} : t \quad \Gamma' \vdash i}{\Gamma \vdash \text{fun}(a_1, \dots, a_n)\{i\} : (t_1, \dots, t_n) \rightarrow t} \text{ (FUN)}$$

5.6 Phrases

Le typage des phrases est détaillé dans la figure 5.5. Le typage 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. En revanche, la déclaration d'une variable globale commence de la même manière, mais on enrichit l'environnement de cette nouvelle valeur.

Type sémantique	$\tau ::= \text{INT}$	Entier
	FLOAT	Flottant
	UNIT	Unité
	$\tau *$	Pointeur
	$\tau []$	Tableau
	$\{l_1 : \tau_1; \dots; l_n : \tau_n\}$	Structure
	FUN_n	Fonction

FIGURE 5.6: Types sémantiques

5.7 Sûreté du typage

Comme nous l'évoquions au début de ce chapitre, le but du typage est de rejeter certains programmes afin de ne garder que ceux qui ne provoquent pas un certain type d'erreurs à l'exécution.

Dans cette section, nous donnons des propriétés que respectent tous les programmes bien typés. Il est traditionnel de rappeler l'adage de Robin Milner :

Well-typed programs don't go wrong.

To go wrong reste bien sûr à définir ! Cette sûreté du typage repose sur les deux théorèmes :

- progrès : si un terme est bien typé, il y a toujours une règle d'évaluation qui s'applique.
- préservation (ou *subject reduction*) : l'évaluation transforme un terme bien typé en un terme du même type.

5.8 Typage des valeurs

Puisque nous allons manipuler les propriétés statiques et dynamiques des programmes, nous allons avoir à traiter des environnements de typage Γ et des états mémoires m . La première chose à faire est donc d'établir une correspondance entre ces deux mondes.

Étant donné un état mémoire m , on associe un type de valeur (ou type sémantique) τ aux valeurs v . Cela est fait sous la forme d'un jugement $m \models v : \tau$.

Ces types sémantiques ne sont pas exactement les mêmes que les types statiques. Pour les calculer, on n'a pas accès au code du programme, seulement à ses données. Il est par exemple possible de reconnaître le type des constantes, mais pas celui des fonctions. Celles-ci sont en fait le seul cas qu'il est impossible de déterminer statiquement. On le remplace donc par un cas plus simple où seul l'arité est conservée.

Le cas des références (règle S-REF) utilise le typage des left-values, codéfini par :

$$m \models_{\Phi} \varphi : \tau \stackrel{\text{def}}{=} m \models m[\varphi]_{\Phi} : \tau$$

Les règles de définition du typage sémantique sont données dans la figure 5.6.

Les règles sont détaillées dans la figure 5.7 : les types des constantes sont simples à retrouver car il y a assez d'information en mémoire. Pour les références, ce qui peut être déréférencé en une valeur de type τ est un $\tau *$. Le typage des valeurs composées se fait en profondeur. Enfin, la seule information restant à l'exécution sur les fonctions est son arité.

$m \models v : \tau$		
$\frac{}{m \models n : \text{INT}}$ (S-INT)	$\frac{}{m \models d : \text{FLOAT}}$ (S-FLOAT)	$\frac{}{m \models () : \text{UNIT}}$ (S-UNIT)
$\frac{}{m \models \text{NULL} : \tau *}$ (S-NULL)	$\frac{m \models_{\Phi} \varphi : \tau}{m \models \hat{\&} \varphi : \tau *}$ (S-REF)	$\frac{\forall i \in [1; n]. m \models v_i : \tau}{m \models [\overline{v_1; \dots; v_n}] : \tau []}$ (S-ARRAY)
$\frac{\forall i \in [1; n]. m \models v_i : \tau_i}{m \models \{\overline{l_1 : v_1; \dots; l_n : v_n} : \{\overline{l_1 : \tau_1; \dots; l_n : \tau_n}\}}$ (S-STRUCT)		
$\frac{}{m \models \text{fun}(x_1, \dots, x_n)\{i\} : \text{FUN}_n}$ (S-FUN)		

FIGURE 5.7: Règles de typage sémantique

$\text{Comp}(\tau, t)$		
$\frac{t \in \{\text{INT}, \text{FLOAT}, \text{UNIT}\}}{\text{Comp}(t, t)}$ (COMP-GROUND)	$\frac{\text{Comp}(\tau, t)}{\text{Comp}(\tau *, t *)}$ (COMP-PTR)	
$\frac{\text{Comp}(\tau, t)}{\text{Comp}(\tau [], t [])}$ (COMP-ARRAY)	$\frac{\forall i \in [1; n]. \text{Comp}(\tau_i, t_i)}{\text{Comp}(\{\overline{l_1 : \tau_1; \dots; l_n : \tau_n}, \{\overline{l_1 : t_1; \dots; l_n : t_n}\})}$ (COMP-STRUCT)	
$\frac{}{\text{Comp}(\text{FUN}_n, (t_1, \dots, t_n) \rightarrow t)}$ (COMP-FUN)		

FIGURE 5.8: Compatibilité entre types sémantiques et statiques

La prochaine étape est de définir une relation de compatibilité entre les types sémantiques τ et statiques t . Nous noterons ceci sous la forme d'un jugement $\text{Comp}(\tau, t)$. Les règles sont décrites dans la figure 5.8, la règle importante étant COMP-FUN.

Grâce à ce jugement, on peut donner la définition suivante.

Définition 5.1 (État mémoire bien typé). *On dit qu'un état mémoire m est bien typé sous un environnement Γ , ce que l'on note $\Gamma \models m$, si les domaines de Γ et m coïncident sur la partie visible, et que les types sémantiques des variables visibles coïncident avec leurs types statiques.*

$$\Gamma \models m \stackrel{\text{def}}{=} \left\{ \begin{array}{l} \text{dom}(\Gamma) = \{\text{Varname}(a) / a \in \text{Visible}(m)\} \\ \forall a \in \text{Visible}(m). \exists \tau, t. \left\{ \begin{array}{l} \Gamma \vdash \text{Varname}(a) : t \\ m \models m[a]_A : \tau \\ \text{Comp}(\tau, t) \end{array} \right. \end{array} \right.$$

$\text{Visible}(m)$ désigne l'ensemble des adresses du cadre de pile le plus récents ainsi que des variables globales :

$$\text{Visible}((s, g)) \stackrel{\text{def}}{=} \{(x_i) / (x_i \mapsto v_i) \in g\} \cup \{(|s|, x_i) / (x_i \mapsto v_i) \in s_0\}$$

$\text{Varname}(a)$ désigne le nom simple d'une variable :

$$\begin{aligned}\text{Varname}((x)) &\stackrel{\text{def}}{=} x \\ \text{Varname}((n, x)) &\stackrel{\text{def}}{=} x\end{aligned}$$

Définition 5.2 (Hauteur d'une valeur). *La hauteur d'une valeur est l'indice du plus haut cadre qu'elle référence, ou -1 sinon.*

$$\begin{aligned}\mathcal{H}(\hat{c}) &= -1 \\ \mathcal{H}(\hat{f}) &= -1 \\ \mathcal{H}(\hat{\&} \varphi) &= \mathcal{H}_\Phi(\varphi) \\ \mathcal{H}(\{l_1 : v_1; \dots; l_n : v_n\}) &= \max_{i \in [1; n]} \mathcal{H}(v_i) \\ \mathcal{H}([v_1, \dots, v_n]) &= \max_{i \in [1; n]} \mathcal{H}(v_i)\end{aligned}$$

où :

$$\begin{aligned}\mathcal{H}_\Phi((x)) &= -1 \\ \mathcal{H}_\Phi((n, x)) &= n \\ \mathcal{H}_\Phi(\varphi.l) &= \mathcal{H}_\Phi(\varphi) \\ \mathcal{H}_\Phi(\varphi[\hat{n}]) &= \mathcal{H}_\Phi(\varphi)\end{aligned}$$

Définition 5.3 (État mémoire bien typé). *On définit par induction une relation entre les environnements de typage Γ et les états mémoire m .*

$$\begin{array}{c} \frac{}{[] \models ([], [])} \text{ (M-Empty)} \qquad \frac{\Gamma \models (s, g) \quad m \models v : \tau \quad \text{Comp}(\tau, t)}{\Gamma, x : t \models (s, ((x \mapsto v) :: g))} \text{ (M-Global)} \\[10pt] \frac{\Gamma \models m \quad ?}{\Gamma \models \text{Push}(m, (x_1 \mapsto v_1, \dots, x_n \mapsto v_n))} \text{ (M-Push)} \qquad \frac{\Gamma \models m \quad ?}{\Gamma \models \text{Pop}(m)} \text{ (M-Pop)} \\[10pt] \frac{\Gamma \models m \quad \Gamma \vdash \varphi : t \quad m \models v : \tau \quad \text{Comp}(\tau, t) \quad m' = m[\varphi \leftarrow v]_\Phi}{\Gamma \models m'} \text{ (M-Write)}\end{array}$$

Lemme 5.1 (Hauteur des chemins typés). *Si $m \models v : \tau$, alors $\mathcal{H}(v) \leq |m|$.*

Démonstration. On procède par induction sur la forme de v .

\hat{c} : Alors $\mathcal{H}(v) = -1$. Comme $|m| \geq 0$, ce cas est établi.

\hat{f} : Idem.

$\hat{\&}(x)$: Idem.

$\widehat{\&}(n, x)$: D'après la forme de v , la dernière règle appliquée dans la dérivation de $m \models v : \tau$ est S-REF, donc $m[(n, x)]_\Phi$ est une valeur. D'après la définition de φ , $n \leq |m|$.

$\widehat{\&}\varphi.l$: On procède par induction sur $v' = \widehat{\&}\varphi$. Comme $\mathcal{H}(\widehat{\&}\varphi) \leq |m|$ et $\mathcal{H}(\widehat{\&}\varphi.l) = \mathcal{H}(\widehat{\&}\varphi)$, on en déduit que $\mathcal{H}(\widehat{\&}\varphi.l) \leq |m|$.

$\widehat{\&}\varphi[n]$: Idem

$\{l_1 : \widehat{v_1}; \dots; l_n : \widehat{v_n}\}$: Par induction, $\forall i \in [1; n], \mathcal{H}(v_i) \leq |m|$. Donc il en est de même pour leur maximum, et $\mathcal{H}(v) \leq |m|$.

$\widehat{[v_1, \dots, v_n]}$: Idem.

□

Lemme 5.2 (Accès à une mémoire bien typée). *Si $\Gamma \models m$, et $\Gamma \vdash \varphi : t$, alors :*

- $m[\varphi]_\Phi$ est une valeur v et $m \models v : \tau$ où $\text{Comp}(\tau, t)$.
- si v est une valeur telle que $m \models v : \tau$ où $\text{Comp}(\tau, t)$, alors en posant $m' = m[\varphi \leftarrow v]_\Phi$, on a $\Gamma \models m'$.

Démonstration. On procède par induction sur une dérivation de $\Gamma \models m$.

M-EMPTY : Il n'existe pas de φ tel que $\Gamma \vdash \varphi : t$. En effet, on montre par induction sur la forme de φ que cela implique qu'il existe a tel que $\text{Varnome}(a) \in \text{dom}(\Gamma)$.

M-GLOBAL : Soit φ tel que $\Gamma, x : t' \vdash \varphi : t$. Soit $m' = (s, ((x \mapsto v) :: g))$.

On veut :

$m'[\varphi]_\Phi$ est une valeur v tq $m' \models v : \tau$ où $\text{Comp}(\tau, t)$.

et

si v est une valeur telle que $m' \models v : \tau$ où $\text{Comp}(\tau, t)$, alors en posant $m'' = m'[\varphi \leftarrow v]_\Phi$, on a $\Gamma, x : t' \models m''$.

M-PUSH :

M-POP :

M-WRITE :

□

5.9 Progrès et préservation

On commence par énoncer quelques lemmes utiles dans la démonstration de ces théorèmes.

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

- Constantes

- $si \Gamma \vdash n : t, \text{ alors } t = \text{Int}$
- $si \Gamma \vdash d : t, \text{ alors } t = \text{Float}$
- $si \Gamma \vdash \text{Null} : t, \text{ alors } \exists t', t = t' *$
- $si \Gamma \vdash () : t, \text{ alors } t = \text{Unit}$
- *Références mémoire :*
 - $si \Gamma \vdash x : t, \text{ alors } x : t \in \Gamma$
 - $si \Gamma \vdash *lv : t, \text{ alors } \Gamma \vdash lv : t *$
 - $si \Gamma \vdash lv[e] : t, \text{ alors } \Gamma \vdash lv : t[] \text{ et } \Gamma \vdash e : \text{Int}$
 - $si \Gamma \vdash lv.l_S : t, \text{ alors } \Gamma \vdash lv : S$
- *Opérations :*
 - $si \Gamma \vdash \boxplus e : t, \text{ alors on est dans un des cas suivants :}$
 - $\boxplus \in \{+, -, \sim, !\}, t = \text{Int}, \Gamma \vdash e : \text{Int}$
 - $\boxplus \in \{+., -.\}, t = \text{Float}, \Gamma \vdash e : \text{Float}$
 - $si \Gamma \vdash e_1 \boxplus e_2 : t, \text{ un des cas suivants se présente :}$
 - $\boxplus \in \{+, -, \times, /, \&, |, ^, \&\&, ||, \ll, \gg, \leq, \geq, <, >\}, \Gamma \vdash e_1 : \text{Int}, \Gamma \vdash e_2 : \text{Int}, t = \text{Int}$
 - $\boxplus \in \{+., -., \times., /., \leq., \geq., <., >.\}, \Gamma \vdash e_1 : \text{Float}, \Gamma \vdash e_2 : \text{Float}, t = \text{Float}$
 - $\boxplus \in \{=, \neq\}, \Gamma \vdash e_1 : t', \Gamma \vdash e_2 : t', \text{Eq}(t'), t = \text{Int}$
 - $\boxplus \in \{\leq, \geq, <, >\}, t = \text{Int}, \Gamma \vdash e_1 : t', \Gamma \vdash e_2 : t', t' \in \{\text{Int}, \text{Float}\}$
 - $\boxplus \in \{+_p, -_p\}, \exists t', t = t' *, \Gamma \vdash e_1 : t' *, \Gamma \vdash e_2 : \text{Int}$
- *Appel de fonction :* $si \Gamma \vdash e(e_1, \dots, e_n) : t, \text{ il existe } (t_1, \dots, t_n) \text{ tels que}$

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

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

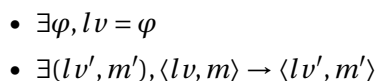
Démonstration. Pour chaque jugement, on considère les règles qui peuvent amener à cette conclusion.

- *Références mémoire :*
 - $\Gamma \vdash x : t$
La seule règle de cette forme est LV-VAR. Puisque sa prémisse est vraie, on en conclut que $x : t \in \Gamma$.
 - $\Gamma \vdash \varphi[] : t$
Idem avec LV-INDEX.
 - $\Gamma \vdash \varphi.l : t$
Idem avec LV-FIELD.
- *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 $\Gamma \vdash \text{fun}(a_1, \dots, a_n)\{i\} : t$ est FUN.

□

Lemme 5.4 (Formes canoniques). *Il est possible de déterminer la forme syntaxique d'une valeur étant donné son type, comme décrit dans le tableau suivant. Par exemple, d'après la première ligne, si $\Gamma \vdash v : \text{Int}$, alors $\exists n, v = d$.*

Ces lemmes étant établis, on énonce maintenant le théorème de progrès. Contrairement aux langages où tout est expression, il faut traiter séparément les trois constructions principales de SAFESPEAK : les expressions, les left-values et les instructions. Celles-ci sont mutuellement dépendantes de par le graphe suivant.



- $\exists \Omega \in \{\Omega_{div}, \Omega_{array}, \Omega_{ptr}\}, \langle l v, m \rangle \rightarrow \Omega$

C'est-à-dire, soit :

- l'entité (instruction, expression ou left-value) est complètement évaluée.
- un pas d'évaluation est possible.
- une erreur de division, tableau ou pointeur se produit.

La preuve du théorème 5.1 se trouve en annexe D.2.

Lemme 5.5 (Représentabilité). *On définit un opérateur de représentation d'un type statique à l'exécution :*

$$\begin{aligned}
 \text{Repr}(\text{Int}) &= \text{Int} \\
 \text{Repr}(\text{Float}) &= \text{Float} \\
 \text{Repr}(\text{Unit}) &= \text{Unit} \\
 \text{Repr}(t' *) &= \text{Repr}(t') * \\
 \text{Repr}(t' []) &= \text{Repr}(t') [] \\
 \text{Repr}(\{l_1 : t_1; \dots; l_n : t_n\}) &= \{l_1 : \text{Repr}(t_1); \dots; l_n : \text{Repr}(t_n)\} \\
 \text{Repr}((t_1, \dots, t_n) \rightarrow t) &= \text{Fun}_n
 \end{aligned}$$

Supposons que $\Gamma \vdash v : t$ et $\Gamma \models m$. On pose $\tau = \text{Repr}(t)$. Alors $m \models v : \tau$ et $\text{Comp}(\tau, t)$.

Démonstration. On procède par induction sur la forme de t .

- **INT** : D'après le lemme des formes canoniques, $v = n$. On conclut avec S-INT et COMP-GROUND.
- **FLOAT** : Idem avec $v = d$ et S-FLOAT.
- **UNIT** : Idem avec $v = ()$ et S-UNIT.
- $t = t' *$: Soient $\tau' = \text{Repr}(t')$ et $\tau = \tau' *$. D'après le lemme des formes canoniques, deux cas sont possibles :
 - $v = \hat{\&} \varphi$:
 Par inversion, $\Gamma \vdash \varphi : t'$
 Puisque $\Gamma \models m$ et $\Phi \vdash \varphi : t', m[\varphi]_\Phi$ est une valeur.
 Par induction, $m \models m[\varphi]_\Phi : \tau'$ et $\text{Comp}(\tau', t')$. D'après S-REF, $m \models \hat{\&} \varphi : \tau$.
 De plus par COMP-PTR $\text{Comp}(\tau, t)$.
 - $v = \text{NULL}$:
 Par induction, $\text{Comp}(\tau', t')$. Alors par COMP-PTR, $\text{Comp}(\tau *, t *)$.
 En outre, grâce à S-NUL, on obtient $m \models \text{NULL} : \tau *$.
- $t = t' []$: Par le lemme des formes canoniques, $v = [\widehat{v_1}, \dots, \widehat{v_n}]$. Par inversion on obtient que $\forall i, \Gamma \vdash v_i : t'$.
 Soient $\tau' = \text{Repr}(t')$ et $\tau = \tau' []$.
 Alors par induction $\forall i, m \models v_i : \tau'$ et $\text{Comp}(\tau, t)$. De la première propriété il vient (via S-ARRAY) $m \models v : \tau$, et de la seconde (via COMP-ARRAY) $\text{Comp}(\tau, t)$.

- $\{l_1 : t_1; \dots; l_n : t_n\}$: Par le lemme des formes canoniques, $v = \{l_1 : \widehat{v_1}; \dots; l_n : \widehat{v_n}\}$. Et par inversion, $\forall i, \Gamma \vdash v_i : t_i$.
Soient $\tau_i = \text{Repr}(t_i)$ et $\tau = \{l_1 : \tau_1; \dots; l_n : \tau_n\}$.
Alors par induction, $\forall i, m \models v_i : \tau_i$ et $\text{Comp}(\tau_i, t_i)$.
On déduit de S-STRUCT que $m \models v : \tau$, et de COMP-STRUCT que $\text{Comp}(\tau, t)$.
- $t = (t_1, \dots, t_n) \rightarrow t'$: Par formes canoniques, on a $v = \text{fun}(\widehat{a_1}, \dots, \widehat{a_n})\{i\}$.
Soit $\tau = \text{FUN}_n$: par S-FUN on obtient que $m \models v : \tau$. On conclut d'autre par que $\text{Comp}(\tau, t)$ grâce à COMP-FUN.

□

Théorème 5.2 (Préservation). *Soient Γ un environnement de typage, et m un état mémoire tels que $\Gamma \models m$.*

Alors :

- Si $\Gamma \vdash e : t$ et $\langle e, m \rangle \rightarrow \langle e', m' \rangle$, alors $\Gamma \models m'$ et $\Gamma \vdash e' : t$.
- Si $\Gamma \vdash e : t$ et $\langle e, m \rangle \rightarrow \langle v, m' \rangle$, alors $\Gamma \models m'$ et $m \models v : \tau$ où $\text{Comp}(\tau, t)$.
- Si $\Gamma \vdash i$ et $\langle i, m \rangle \rightarrow \langle i', m' \rangle$, alors $\Gamma \models m'$ et $\Gamma \vdash i'$.
- Si $\Gamma \vdash lv : t$ et $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$, alors $\Gamma \models m'$ et $\Gamma \vdash lv' : t$.
- Si $\Gamma \vdash lv : t$ et $\langle lv, m \rangle \rightarrow \langle \varphi, m' \rangle$, alors $\Gamma \models m'$ et $m \models_{\Phi} \varphi : \tau$ où $\text{Comp}(\tau, t)$.

Autrement dit, si une construction est typable, alors un pas d'évaluation ne modifie pas son type et préserve le typage de la mémoire.

La preuve de ce théorème se trouve en annexe D.3.

Cela prouve qu'aucun terme ne reste “bloqué” parce qu'aucune règle ne s'applique, et que la sémantique respecte le typage. En quelque sorte, les types sont un contrat entre les expressions et les fonctions : si leur évaluation converge, alors une valeur du type inféré sera produite.

Enfin, on donne une version de ces propriétés pour les phrases de programme.

Théorème 5.3 (Progrès pour les phrases). *Soient Γ un environnement de typage, m un état mémoire et p une phrase de programme. Supposons que $\Gamma \vdash p \rightarrow \Gamma'$ et $\Gamma \models m$.*

On suppose en outre que l'évaluation de p termine.

Alors $\exists m'. m \Vdash p \rightarrow m'$.

Démonstration. On distingue selon la dernière règle de typage appliquée dans la dérivation de $\Gamma \vdash p \rightarrow \Gamma'$.

T-EXP : On applique le théorème de progrès à e . Puisque l'évaluation termine, il existe v et m' tels que $\langle e, m \rangle \rightarrow \langle v, m' \rangle$. Alors on peut appliquer T-EXP.

T-VAR : Idem avec T-VAR.

□

Théorème 5.4 (Préservation pour les phrases). *On suppose que les trois propriétés suivantes sont vérifiées :*

$$\left\{ \begin{array}{l} \Gamma \models m \\ \Gamma \vdash p \rightarrow \Gamma' \\ m \Vdash p \rightarrow m' \end{array} \right.$$

Alors $\Gamma' \models m'$.

Démonstration.



ANALYSE DE PROVENANCE DE POINTEURS

Jusqu'ici SAFESPEAK est un langage impératif généraliste, ne prenant pas en compte les spécificités de l'adressage utilisé dans les systèmes d'exploitation.

Dans ce chapitre, on commence par l'étendre en ajoutant des constructions modélisant les variables présentes dans l'espace utilisateur (cf. chapitre 2). Pour accéder à celles-ci, on ajoute un opérateur de déréréfencement sûr qui vérifie à l'exécution que l'invariant suivant est respecté :

Les pointeurs dont la valeur est contrôlée par l'utilisateur, pointent vers l'espace utilisateur.

La terminologie mérite d'être détaillée :

Un pointeur contrôlé par l'utilisateur, ou *pointeur utilisateur*, est référence mémoire dont la valeur est modifiable par le code utilisateur (opposé au code noyau, que nous analysons ici). Ceci correspond à des données provenant de l'extérieur du système vérifié. Pour faire un parallèle avec la programmation réseau, cela correspond à des variables dont la valeur proviendrait d'un client distant. C'est une propriété statique, qui peut être déterminée à la compilation à partir de considérations syntaxiques. Par exemple, l'adresse d'une variable locale est toujours considérée comme étant contrôlée par le noyau.

Un pointeur pointant vers l'espace utilisateur fait référence à une variable allouée en espace utilisateur. Cela veut dire qu'y accéder ne risque pas de mettre en péril l'isolation du noyau en faisant fuiter des informations confidentielles ou en déjouant son intégrité. Cette propriété est dynamique : un pointeur utilisateur peut *a priori* pointer vers l'espace utilisateur, ou non.

Pour reprendre le parallèle avec la programmation web, une chaîne provenant du réseau peut être correctement échappée ou non. Pour s'assurer de la sécurité du traitement d'une telle chaîne, il faut s'assurer que celles-ci sont correctement échappées.

Pour prouver que cette approche est bien fondée, on procède en plusieurs étapes.

Tout d'abord, on définit un nouveau type d'erreur Ω_{taint} , déclenché lorsqu'un pointeur contrôlé par l'utilisateur et pointant vers le noyau est déréréfencé (le cas que l'on cherche à éviter)¹.

Ensuite, on montre qu'avec cet ajout, si on étend naïvement le système de types en donnant le même type aux pointeurs contrôlés par l'utilisateur et le noyau, le théorème de progrès (5.1) n'est plus valable.

1. Il est important de noter que ce cas d'erreur est "virtuel", la distinction entre pointeur contrôlé par l'utilisateur et par le noyau n'étant normalement pas visible à l'exécution.

L'étape suivante est d'étendre, à son tour, le système de types de SAFESPEAK en distinguant les types des pointeurs contrôlés par l'utilisateur des pointeurs contrôlés par le noyau. Puisqu'on veut interdire le déréférencement des premiers par l'opérateur $*$, on modifie également le typage de celui-ci.

Enfin, une fois ces modifications faites, on prouve que les propriétés de progrès et de préservation sont rétablies.

6.1 Extensions noyau pour SAFESPEAK

On ajoute à SAFESPEAK la notion de valeur provenant de l'espace utilisateur. Pour marquer la séparation entre les deux espaces d'adressage, on ajoute une construction $\varphi ::= \hat{\diamond} \varphi'$. Le chemin interne φ' désigne une variable classique (un pointeur noyau) et l'opérateur $\hat{\diamond} \cdot$ permet de l'interpréter comme un pointeur vers l'espace utilisateur, contrairement à $\&\cdot$ qui crée des pointeurs vers l'espace noyau. En quelque sorte, on ne classe pas les valeurs selon la variable pointée mais selon la construction du pointeur.

En plus du déréférencement par $*$ (qui devra donc renvoyer Ω_{taint} pour les valeurs de la forme $\hat{\diamond} \varphi'$), il faut aussi ajouter des constructions de lecture et d'écriture à travers les pointeurs utilisateurs. Ceci sera fait sous forme de deux fonctions, `copy_from_user` et `copy_to_user`. Celles-ci prennent deux pointeurs en paramètre et renvoient un booléen indiquant si la copie pu être faite (si le paramètre contrôlé par l'utilisateur pointe en espace noyau, les fonctions ne font pas la copie et signalent l'erreur).

Illustrons ceci par un exemple. Imaginons un appel système fictif qui renvoie la version du noyau, en remplissant par pointeur une structure contenant les champs entiers `major`, `minor` et `patch` (un équivalent dans Linux est l'appel système `uname()`). Celui-ci peut être alors écrit comme dans la figure 6.1. Une fois la structure noyau v remplie, il faut la copier vers l'espace utilisateur. La fonction `copy_to_user` va réaliser cette copie (de la même manière qu'avec un `memcpy()`), mais après avoir testé dynamiquement que p pointe en espace utilisateur (dans le cas contraire, la copie n'est pas faite).

```
sys_getver = fun(p){
    DECL  $v = \{ \text{major} : 3; \text{minor} : 14; \text{patch} : 15 \}$  IN
    copy_to_user(p, & $v$ )
}
```

FIGURE 6.1: Implantation d'un appel système qui remplit une structure par pointeur

On peut remarquer que contrairement aux fonctions de Linux, les fonctions `copy_from_user` et `copy_to_user` n'ont pas de paramètre "taille à copier". Cela est dû au fait que le modèle mémoire de SAFESPEAK est plus haut niveau. L'information de taille est déjà présente dans chaque valeur.

On commence donc par ajouter aux instructions des constructions `copy_to_user(\cdot , \cdot)` et `copy_from_user(\cdot , \cdot)` de copie sûre. Afin de leur donner une sémantique, il faut étendre l'ensemble des valeurs pointeur φ aux constructions de la forme $\hat{\diamond} \varphi'$. Pour créer des termes s'évaluant en de telles valeurs, il faut une construction syntaxique $\hat{\diamond} e$ telle que si e s'évalue en φ , $\hat{\diamond} e$ s'évalue en $\hat{\diamond} \varphi$. Cela correspond à 2 ajouts : une un nouveau contexte d'évaluation $\hat{\diamond} \bullet$ et une règle d'évaluation. Enfin, on ajoute une nouvelle erreur Ω_{taint} à déclencher

Expressions	$e ::= \dots$ $\diamond lv$	Adresse utilisateur
Contextes	$C ::= \dots$ $\diamond C$	
Chemins	$\varphi ::= \dots$ $\hat{\diamond} \varphi$	Valeur souillée
Erreurs	$\Omega ::= \dots$ Ω_{taint}	Erreur de souillure

FIGURE 6.2: Ajouts liés aux pointeurs utilisateurs (par rapport à l'évaluateur du chapitre 4)

lorsqu'on déréfère directement un pointeur utilisateur. Ces étapes sont résumées dans la figure 6.2.

6.2 Extensions sémantiques

En ce qui concerne l'évaluation des expressions $\diamond \cdot$, on ajoute la règle suivante :

$$\frac{}{\langle \diamond \varphi, m \rangle \rightarrow \langle \hat{\diamond} \varphi, m \rangle} \text{ (EXPR-TAINTED)}$$

Ensuite, il est nécessaire d'adapter les règles d'accès à la mémoire pour déclencher une erreur Ω_{taint} en cas de déréférencement d'un pointeur utilisateur. Les accès mémoire en lecture proviennent de la règle EXP-LV et ceux en lecture, de la règle EXP-SET, rappelées ici :

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

Les accès à la mémoire sont en effet faits par le biais de la lentille Φ . Il suffit donc d'adapter sa définition (page 50) de celle-ci en rajoutant le cas :

$$\Phi(\hat{\diamond} \varphi) = \Omega_{taint}$$

Enfin, il est nécessaire de donner une sémantique aux fonctions `copy_from_user` et `copy_to_user`. L'idée est que celles-ci testent dynamiquement la *valeur* du paramètre contrôlé par l'utilisateur afin de vérifier que celui-ci pointe vers l'espace utilisateur (c'est-à-dire, qu'il est de la forme $\hat{\diamond} \varphi$).

Deux cas peuvent se produire. Soit la partie à vérifier a la forme $\hat{\diamond} \varphi'$, soit non (et dans ce cas $\nexists \varphi', \varphi = \hat{\diamond} \varphi'$). Dans le premier cas (règles USER-* -OK), alors la copie est faite et l'opération de copie retourne la valeur entière 0. Dans le second (règles USER-* -ERR), aucune copie n'est faite et la valeur -1 est retournée. Ce comportement est calé celui des fonctions `copy_{from,to}_user` du noyau Linux : en cas de succès elles renvoient 0, et en cas d'erreur -EFAULT (= -14).

$$\begin{array}{c}
\frac{\langle \varphi_d \leftarrow \varphi_s, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle \text{copy_from_user}(\varphi_d, \hat{\diamond} \varphi_s), m \rangle \rightarrow \langle 0, m' \rangle} \text{ (USER-GET-OK)} \\
\\
\frac{\nexists \varphi_s. \varphi = \hat{\diamond} \varphi_s}{\langle \text{copy_from_user}(\varphi_d, \varphi), m \rangle \rightarrow \langle -1, m \rangle} \text{ (USER-GET-ERR)} \\
\\
\frac{\langle \varphi_d \leftarrow \varphi_s, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle \text{copy_to_user}(\hat{\diamond} \varphi_d, \varphi_s), m \rangle \rightarrow \langle 0, m' \rangle} \text{ (USER-PUT-OK)} \\
\\
\frac{\nexists \varphi_d. \varphi = \hat{\diamond} \varphi_d}{\langle \text{copy_to_user}(\varphi, \varphi_s), m \rangle \rightarrow \langle -1, m \rangle} \text{ (USER-PUT-ERR)}
\end{array}$$

Ces règles sont à appliquer en priorité de la règle d'appel de fonction classique, puisque ces fonctions n'ont pas de corps en tant que tel. En effet celles ci ne sont pas implantables directement en SAFESPEAK, puisqu'il n'y a pas par exemple d'opérateur permettant d'extraire φ depuis une valeur $\hat{\diamond} \varphi$. L'opération en "boîte noire" de ces deux fonctions permet d'assurer que l'accès à l'espace utilisateur est toujours couplé à un test dynamique.

6.3 Insuffisance des types simples

Étant donné SAFESPEAK augmenté de cette extension sémantique, on peut étendre trivialement le système de types avec la règle suivante :

$$\frac{\Gamma \vdash lv : t}{\Gamma \vdash \hat{\diamond} lv : t * } \text{ (TAINT-IGNORE)}$$

Cette règle est compatible avec l'extension, sauf qu'elle introduit des termes qui sont bien typables mais dont l'évaluation provoque une erreur autre que Ω_{div} , Ω_{array} ou Ω_{ptr} , violant ainsi le théorème 5.1. Posons :

$$\begin{cases} e = * \hat{\diamond} x \\ \Gamma = x : \text{INT} \\ m = ([x \mapsto 0], []) \end{cases}$$

Les hypothèses du théorème de progrès sont bien vérifiées, mais cependant la conclusion n'est pas vraie :

- On a bien $\Gamma \models m$. En effet $\text{Visible}(m) = \{(0, x)\}$ et :

$$\begin{cases} \Gamma \vdash x : \text{INT} \\ m \models 0 : \text{INT} \\ \text{Comp}(\text{INT}, \text{INT}) \end{cases}$$

- e est bien typée sous Γ :

$$\begin{array}{c}
\frac{x : \text{INT} \in \Gamma}{\Gamma \vdash x : \text{INT}} \text{ (LV-VAR)} \\
\frac{\Gamma \vdash x : \text{INT}}{\Gamma \vdash \&x : \text{INT}^*} \text{ (LV-DEREF)} \\
\frac{\Gamma \vdash \&x : \text{INT}^*}{\Gamma \vdash \diamond x : \text{INT}^*} \text{ (TAINT-IGNORE)} \\
\frac{\Gamma \vdash \diamond x : \text{INT}^*}{\Gamma \vdash * \diamond x : \text{INT}} \text{ (LV-DEREF)}
\end{array}$$

- L'évaluation de e sous m provoque une erreur différente de Ω_{div} , Ω_{array} , ou Ω_{ptr} :

$$\frac{\frac{\overline{m[* \diamond x] = \Omega_{taint}}}{\langle * \diamond x, m \rangle \rightarrow \langle \Omega_{taint}, m \rangle} \text{ (EXP-LV)} \quad \frac{}{\langle \Omega_{taint}, m \rangle \rightarrow \Omega_{taint}} \text{ (EVAL-ERR)}}{\langle e, m \rangle \rightarrow \Omega_{taint}}$$

Cela montre que le typage n'apporte plus de garantie de sûreté sur l'exécution : le système de types naïvement étendu par une règle comme Taint-Ignore n'est pas en adéquation avec les extensions présentées dans la section 6.1.

6.4 Extensions du système de types

On présente ici un système de types plus expressif permettant de capturer les extensions de sémantique. *In fine*, cela permettra de prouver le théorème 6.1 qui est l'équivalent du théorème 5.1 mais pour le nouveau jugement de typage.

Définir un nouveau système de types revient à étendre le jugement de typage $\cdot \vdash \cdot : \cdot$, en modifiant certaines règles et en en ajoutant d'autres. Naturellement, la plupart des différences porteront sur le traitement des pointeurs.

Pointeurs utilisateurs

Le changement clef est l'ajout de *pointeurs utilisateurs*. En plus des types pointeurs t^* , on ajoute des types pointeurs $t@$. La différence entre les deux représente *qui* contrôle leur valeur (section 2.4).

Les différences sont les suivantes (figure 6.3) :

- Les types " t^* " s'appliquent aux pointeurs contrôlés par le noyau. Par exemple, prendre l'adresse d'un objet de la pile noyau donne un pointeur noyau.

Type	$t ::= \dots$	
	$t@$	Pointeur utilisateur
Type sémantique	$\tau ::= \dots$	
	$\tau@$	Pointeur utilisateur

FIGURE 6.3: Ajouts liés aux qualificateurs de types (par rapport aux figures 5.2 et 5.6)

	Noyau	Utilisateur
Syntaxe	$\& x$	$\diamond x$
Valeur	$\hat{\&} x$	$\hat{\diamond} x$
Type	$t *$	$t @$
Accès	$* x$	copy_*_user

FIGURE 6.4: Notations utilisées pour les pointeurs

- Les types “ $t @$ ”, quant à eux, s’appliquent aux pointeurs qui proviennent de l’espace utilisateur. Ces pointeurs proviennent toujours d’interfaces particulières, comme les appels système ou les paramètres passés aux implantations de la fonction `ioctl`.

L’ensemble des notations est résumé dans le tableau de la figure 6.4.

Puisqu’on s’intéresse à la provenance des pointeurs, détaillons les règles qui créent, manipulent et utilisent des pointeurs.

Sources de pointeurs

La source principale de pointeurs est l’opérateur $\&$ qui prend l’adresse d’une variable. Celle-ci est bien entendue contrôlée par le noyau (dans le sens où son déréférencement est toujours sûr). Cette construction crée donc des pointeurs noyau :

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

Manipulations de pointeurs

L’avantage du typage est que celui-ci suit le flot de données : si à un endroit une valeur de type t est affectée à une variable, que le contenu de cette variable est placé puis retiré d’une structure de données, il conserve ce type t . En particulier un pointeur utilisateur reste un pointeur utilisateur.

Une seule règle consomme un pointeur et en retourne un. Elle concerne l’arithmétique des pointeurs. On l’étend aux pointeurs utilisateur :

$$\frac{\boxplus \in \{+, -\} \quad \Gamma \vdash e_1 : t @ \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxplus e_2 : t @} \text{ (PTR-ARITH-USER)}$$

Utilisations de pointeurs

La principale restriction, au cœur de ce chapitre et de cette thèse, est que seuls les pointeurs noyau peuvent être déréférencés de manière sûre. La règle capitale est donc la suivante (déjà introduite dans le chapitre 5) :

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

Ainsi, on interdit le déréférencement des expressions de type $t @$ à la compilation.

L’opérateur \diamond transforme un pointeur selon la règle suivante :

$$\frac{\Gamma \vdash lv : t}{\Gamma \vdash \Diamond lv : t@} \text{ (TAINT)}$$

Les “fonctions” `copy_from_user` et `copy_to_user` sont typées de la manière suivante. Il est à remarquer que ce ne sont pas vraiment des fonctions et qu’elles n’ont pas un type en $(t_1, t_2) \rightarrow t$, car il faudrait un type polymorphe pour pouvoir les appliquer n’importe quel type de pointeurs. Leur typage est donc plus proche de celui d’un opérateur.

$$\frac{\Gamma \vdash e_d : t * \quad \Gamma \vdash e_s : t@}{\Gamma \vdash \text{copy_from_user}(e_d, e_s) : \text{INT}} \text{ (GETU)} \quad \frac{\Gamma \vdash e_d : t@ \quad \Gamma \vdash e_s : t *}{\Gamma \vdash \text{copy_to_user}(e_d, e_s) : \text{INT}} \text{ (PUTU)}$$

Remarquons qu’il n’y a pas de sous-typage : les pointeurs noyau ne peuvent être utilisés qu’en tant que pointeurs noyau, et les pointeurs utilisateurs qu’en tant que pointeurs utilisateurs.

6.5 Sûreté du typage

Typage sémantique

La définition du typage sémantique doit aussi être étendue au cas $\varphi = \Diamond \varphi'$.

$$\frac{m \models \varphi' : \tau}{m \models \Diamond \varphi' : \tau@} \text{ (S-TAINTED)} \quad \frac{\text{Comp}(\tau, t)}{\text{Comp}(\tau @, t @)} \text{ (COMP-PTR)}$$

Progrès et préservation

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.

Théorème 6.1 (Progrès pour les types qualifiés). *Supposons que $\Gamma \vdash e : t$. Soit m un état mémoire tel que $\Gamma \models m$. Alors l’un des cas suivants est vrai :*

- $\exists v \neq \Omega, e = v$
- $\exists (e', m'), \langle e, m \rangle \rightarrow \langle e', m' \rangle$
- $\exists \Omega \in \{\Omega_{div}, \Omega_{array}, \Omega_{ptr}\}, \langle e, m \rangle \rightarrow \Omega$

La preuve de ce théorème est en annexe D.4.

Et nous donnons un équivalent du théorème 5.2.

Théorème 6.2 (Préservation pour les types qualifiés). *Soient Γ un environnement de typage, e une expression et m un état mémoire.*

On suppose que les hypothèses suivantes sont vérifiées :

- $\Gamma \vdash e : t$
- $\Gamma \models m$
- $\langle e, m \rangle \rightarrow \langle e', m' \rangle$

Alors $\Gamma \vdash e' : t$ et $\Gamma \models m'$.

Autrement dit, si une expression est typable, alors un pas d’évaluation ne modifie pas son type et préserve le typage de la mémoire.

La preuve de ce théorème est en annexe D.5.

La sûreté du typage étant à nouveau établie, on a montré que les qualificateurs de types sur les pointeurs suffisent pour avoir une adéquation entre les extensions de sémantique de la section 6.1 et les extensions du système de type de la section 6.4.

CONCLUSION DE LA PARTIE II

Comme nous l'avons fait remarquer dans le chapitre 3, utiliser une technique de typage pour étudier des propriétés sur les données a déjà été explorée dans l'outil CQual [FFA99], en particulier sur les problèmes de pointeurs utilisateurs [JW04].

En effet, si on remplace “ $t *$ ” par “ $\text{KERNEL } t *$ ” et “ $t @$ ” par “ $\text{USER } t *$ ”, on obtient un début de système de types qualifiés.

En revanche, il y a une différence fondamentale : CQual modifie fondamentalement l'ensemble du système de types, pas SAFESPEAK. Le jugement de typage de CQual a pour forme générale $\Gamma \vdash e : q \ t$ (où Γ est un environnement de typage, e une expression, q un qualificateur et t un type), alors que le nôtre a la forme plus classique $\Gamma \vdash e : t$.

En intégrant q à la relation de typage, on ajoute un qualificateur à chaque type, même les expressions pour lesquels il n'est pas directement pertinent de déterminer qui les contrôle (comme par exemple, un entier). Dans CQual, ceci permet de traiter de manière correcte le transtypage. Par exemple, si e a pour type qualifié USER INT , alors $(\text{FLOAT } *) \ e$ aura pour type qualifié $\text{USER FLOAT } *$, et déréférencer cette expression produira une erreur de typage. SAFESPEAK, dans son état actuel, ne permet pas de traiter ce cas.

Au contraire, nous ne modifions le système de types que là où cela est nécessaire, c'est-à-dire sur les types pointeurs. Cela permet de ne pas avoir à modifier en profondeur un système de types existant.

Une autre différence fondamentale est que le système de types de CQual fait intervenir une relation de sous-typage. Dans le cas particulier du problème déréférencement des pointeurs utilisateurs peut être traité dans ce cadre en posant $\text{KERNEL} \leq \text{USER}$ pour restreindre certaines opérations aux pointeurs KERNEL .

Notre approche, au contraire, n'utilise pas de sous-typage, mais consiste à définir un type abstrait $t @$ partageant certaines propriétés avec $t *$ (comme la taille et la représentation) mais incompatible avec certaines opérations. C'est à rapprocher du système de masquage par interface des langages ML, ou encore de la définition de types abstraits par le mot-clef `newtype` en Haskell.

Les perspectives de travaux futurs sont également très différentes. Dans le cas des pointeurs, même si le noyau Linux (et la plupart des systèmes d'exploitation) ne comportent que deux espaces d'adressage, il est commun dans les systèmes embarqués de manipuler des pointeurs provenant d'espaces mémoire indépendants : par exemple, de la mémoire flash, de la RAM, ou une EEPROM de configuration. Ces différentes mémoires possèdent des adresses, et un pointeur est interprété comme faisant référence à une ou l'autre selon le code dont il est tiré. Dans ces cas, la notion de sous-typage n'est pas adaptée.

Troisième partie

Expérimentation

On décrit ici la démarche expérimentale liée à l'implémentation des analyses décrites dans la partie II.

Le chapitre 7 décrit l'implémentation en elle-même : comment le code source C est compilé vers SAFESPEAK, et comment les types du programme sont vérifiés.

Ensuite, dans le chapitre 8, le cas d'un bogue de pilote graphique dans le noyau Linux est étudié. On montre que les analyses précédentes permettent de distinguer statiquement entre le cas incorrect et le cas corrigé.

Enfin, le chapitre 9 conclut : les limitations de cette approche sont présentées, ainsi qu'un résumé des contributions de cet ouvrage.

IMPLANTATION

Dans ce chapitre, nous décrivons la mise en œuvre des analyses statiques précédentes. Nous commençons par décrire la représentation intermédiaire utilisée : Newspeak. La chaîne de compilation est explicitée, partant de C pour aller jusqu'à Tyspeak, sur lequel nous décrivons d'un algorithme d'inférence de types à la Hindley-Milner, reposant sur l'unification et le partage de références.

En réalité, Newspeak (et donc Tyspeak) utilise une représentation de plus bas niveau que SAFESPEAK donc les règles sont légèrement différentes.

7.1 Newspeak

Newspeak [HL08] est un langage conçu pour être à la fois :

- Précis : sa sémantique est définie formellement dans [HL08]
- Expressif : la plupart des primitives présentes dans les langages de bas niveau sont compilables en Newspeak.
- Simple : peu de primitives sont présentes.
- Minimal : aucun élément syntaxique ne peut être exprimé comme combinaison d'autres.
- Explicite : les constructions sont toutes indépendantes du contexte.
- Orienté analyse : les primitives sont décorées d'informations reflétant leur validité (tests de bornes, etc)
- Indépendant de l'architecture : toutes les caractéristiques comme la taille des types ou l'alignement des structures sont rendues explicites.

7.2 Chaîne de compilation

La compilation vers C est faite en trois étapes (figure 7.1) : prétraitement du code source, compilation de C prétraité vers NEWSPEAK, puis compilation de NEWSPEAK vers ce langage.

La première étape consiste à prétraiter les fichiers C source avec le logiciel cpp, comme pour une compilation normale. Cette étape interprète les directives comme `#include`, `#ifdef`, À cet étape, les commentaires sont aussi supprimés.

Une fois cette passe effectuée, le résultats est un ensemble de fichiers C prétraités, c'est-à-dire des unités de compilation.

Puisque la directive `#include` est textuelle, ces fichiers sont très grands et donnent lieu à beaucoup de duplication dans les passes suivantes.

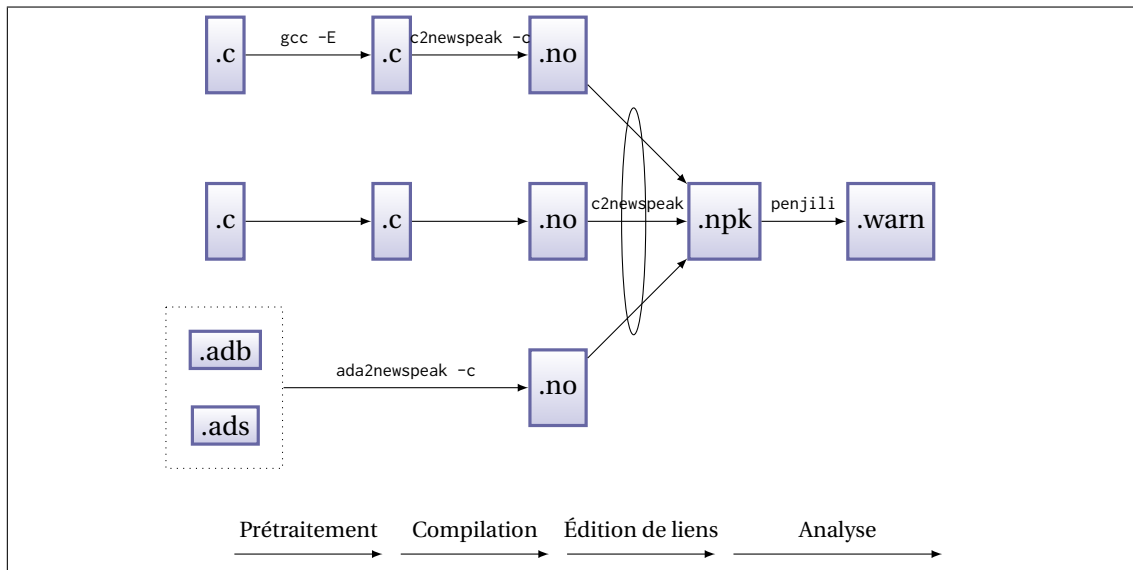


FIGURE 7.1: Compilation depuis Newspeak

Sur cette représentation (du C prétraité), il est possible d'ajouter des annotations de la forme `/*!npk [...] */` qui pourront être accessibles dans l'arbre de syntaxe abstrait des passes suivantes.

À ce niveau, les fichiers sont passés à l'outil `c2newspeak` qui les traduit vers Newspeak. Dans cette étape, les types et les noms sont résolus, et le programme est annoté de manière à rendre les prochaines étapes indépendante du contexte. Par exemple, chaque déclaration de variable est adjointe d'une description complète du type.

Lors de cette étape, le flût de contrôle est également simplifié. C en effet propose de nombreuses constructions ambiguës ou redondantes.

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` [EH94], 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].

Newspeak est conçu pour l'analyse statique par interprétation abstraite. Il a donc une vue de bas niveau sur les programmes. Par exemple, aucune distinction n'est faite entre l'accès à un champ et l'accès à un élément d'un tableau (tous deux sont traduits par un décalage numérique depuis le début de la zone mémoire). Pour supprimer cette ambiguïté, il faut s'interfacer dans les structures internes de `c2newspeak`, où les informations nécessaires sont encore présentes.

Ensuite, les différents fichiers sont liés ensemble. Cet étape consiste principalement à s'assurer que les hypothèses faites par les différentes unités de compilation sont cohérentes entre elles. Les objets marqués `static`, invisibles à l'extérieur de leur unité de compilation, sont renommés afin qu'ils aient un nom unique.

Enfin, l'implantation d'un algorithme d'inférence pour les systèmes de types décrits dans les chapitres 5 et 6 est assez simple. Puisqu'ils sont suffisamment proches du lambda calcul simplement typé, on peut utiliser une variante de l'algorithme W de Damas et Milner [DM82]. On utilise l'optimisation classique qui consiste à se reposer sur le partage de références pour

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

FIGURE 7.2: Exemple de programme C à compiler

réaliser l'unification, plutôt que de faire des substitutions explicites. Puisque ces systèmes de types sont monomorphes, on présente une erreur si des variable de type libres sont présentes.

À la fin de cette étape, on obtient soit un programme complètement annoté, soit une erreur de type.

Prenons l'exemple de la figure 7.2 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 $Ptr(t_1)$. 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 $Ptr(Int)$.

```
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
```

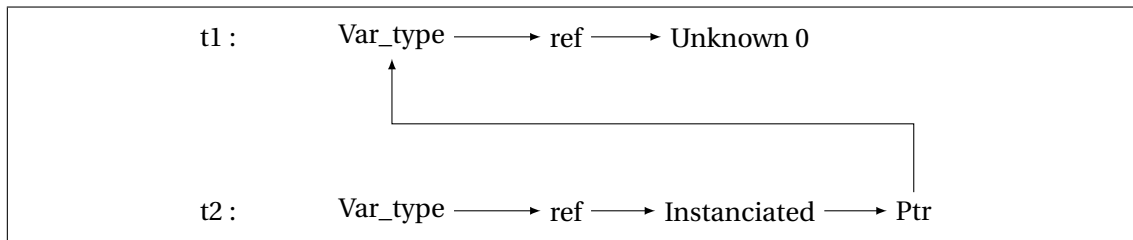


FIGURE 7.3: Partage lors de l'unification

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 faut 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

```

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 7.3).

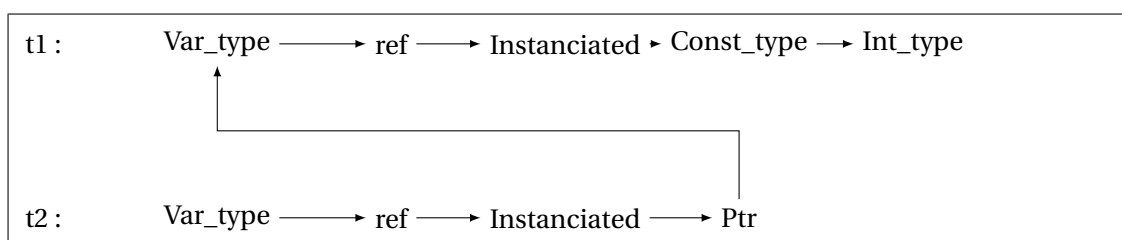


FIGURE 7.4: Unification par mutation de références

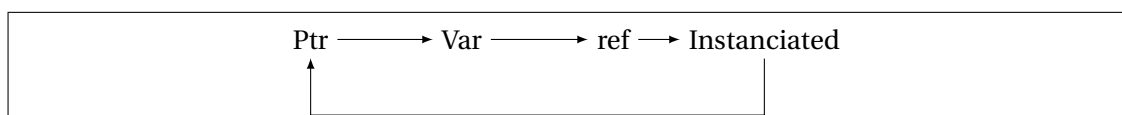


FIGURE 7.5: Cycle dans le graphe de types

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 7.4).

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`.

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 (Instanciated 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 7.5). 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.

```

let process_npk npk =
  let tpk = Npk2tpk.convert_unit npk in
  let order = Topological.topological_sort (Topological.make_graph npk) in

  let function_is_defined f =
    Hashtbl.mem tpk.Tyspeak.fundecs f
  in

  let (internal_funcs, external_funcs) =
    List.partition function_is_defined order
  in

  let exttbl = Printer.parse_external_type_annotations tpk in

  let env =
    env_add_external_fundecs exttbl external_funcs Env.empty
  in
  let s = Infer.infer internal_funcs env tpk in
  begin
    if !Options.do_checks then
      Check.check env s
  end;
  Printer.dump s

```

FIGURE 7.6: Fonction principale de ptrtype

- 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.
- si les deux types sont des types structures, on étudie les types de leurs champs. On commence par unifier les types des champs de même nom, et on ajoute à l'un les champs exclusifs à l'autre.
- dans les autres cas, l'algorithme échoue.

7.3 Architecture de ptrtype

L'outil ptrtype lit un programme Newspeak (ou un fichier C), et réalise l'inférence de qualificateurs. En sortie, il affiche soit le programme complètement annoté, soit une erreur. Le cœur de l'outil est dans la

Si l'argument à ptrtype est un fichier C, il est tout d'abord compilé vers Newspeak grâce à l'utilitaire c2newspeak. Ensuite, les autres passes travaillent sur une représentation intermédiaire proche de Newspeak, mais où des étiquettes de type supplémentaires sont ajoutées. Ce type de représentation intermédiaire (polymorphe en le type des étiquettes) est 'a Tyspeak.t.

Le reste de l'outil est résumé dans la fonction process_npk (figure 7.6) :

- Grâce à la fonction convert_unit : Newspeak.t -> unit Tyspeak.t, on ajoute des étiquettes “vides” (toutes égales à () : unit).

- L'ensemble des fonctions du programme est trié topologiquement selon la relation \leq définie par $f \leq g \stackrel{\text{def}}{=} \text{“} g \text{ apparaît dans la définition de } f \text{”}$. Cela est fait en construisant une représentation de \leq sous forme de graphe, puis en faisant un parcours en largeur de celui-ci.
- Les annotations extérieures sont alors lues (variable `ext tbl`), ce qui permet de créer un environnement initial.
- Les types de chaque fonction sont inférés, par le biais de la fonction suivante :

```
val infer : Newspeak.fid list
    -> Types.simple Env.t
    -> 'a Tyspeak.t
    -> Types.simple Tyspeak.t
```

- Le programme obtenu, de type `Types.simple Tyspeak.t`, est affiché sur le terminal.

7.4 Inférence de types

L'inférence de types consiste à remplacer les étiquettes de type `unit t` par des étiquettes de type `simple` (autrement dit de vraies représentations de types).

Cette étape se fait de manière impérative : on peut créer de nouveaux types avec `new_unknown` et unifier deux types avec `unify`. Leurs types sont :

```
val new_unknown : unit -> Types.simple
val unify : Types.simple -> Types.simple -> unit
```

La fonction `infer` s'appuie sur un ensemble de fonctions récursivement définies portant sur chaque type de fragment : `infer_fdec` pour les déclarations de fonction, `infer_exp` pour les expressions, `infer_stmtkind` pour les instructions, etc.

Les règles de typage sont implantées par `new_unknown` et `unify`. Par exemple, pour typer une déclaration (figure 7.7), on crée un nouveau type `t0`. On étend l'environnement courant avec cette nouvelle association et sous ce nouvel environnement, on type le bloc de portée de la déclaration.

De même, pour typer un appel de fonction, on infère le type de ses arguments et left-values de retour. On obtient également le type de la fonction (à partir du type de la fonction présent dans l'environnement, ou du type du pointeur de fonction qui est déréférencé), et on unifie ces deux informations.

On peut donner quelques exemples, comme :

- addition de deux flottants (dans `infer_binop`) :

```
let infer_binop op (_, a) (_, b) =
  match op with
  (* [...] *)
  | N.PlusF _ ->
    unify a Float;
    unify b Float;
    Float
```

- adresse d'une left-value (dans `lval_type`) :

```

let rec infer_stmtkind env sk =
  match sk with
  (* [...] *)
  | T.Decl (n, nty, _ty, blk) ->
    let var = T.Local n in
    let t0 = new_unknown () in
    let new_env = Env.add (VLocal n) (Some nty) t0 env in
    let blk' = infer_blk new_env blk in
    let ty = lval_type new_env var in
    T.Decl (n, nty, ty, blk')
  (* [...] *)
  | T.Call (args, fexp, rets) ->
    let infer_arg (e, nt) =
      let et = infer_exp env e in
      (et, nt)
    in

    let infer_ret (lv, nt) =
      (infer_lv env lv, nt)
    in

    let args' = List.map infer_arg args in
    let rets' = List.map infer_ret rets in

    let t_args = List.map (fun (_, t), _ -> t) args' in
    let t_rets = List.map (fun (lv, _) -> lval_type env lv) rets' in

    let (fexp', tf) = infer_funexp env fexp in
    let call_type = Fun (t_args, t_rets) in
    unify tf call_type;
    T.Call (args', fexp', rets')

```

FIGURE 7.7: Inférence des déclarations de variable et appels de fonction

```

| T.AddrOf lv ->
  let lv' = infer_lv env lv in
  let ty = lval_type env lv in
  (T.AddrOf lv', Ptr (Kernel, ty))

```

- déréférencement d'une left-value (dans lval_type) :

```

| T.Deref(e, _sz) ->
  let (_, te) = infer_exp env e in
  let t = new_unknown () in
  unify (Ptr (Kernel, t)) te;
  t

```

7.5 Unification

```

type unknown = { id : int }

type 'a var_type =
  | Unknown of unknown
  | Instanciated of 'a

and qual =
  | Kernel
  | User
  | QVar of qual var_type ref

and simple =
  | Int
  | Float
  | Fun of simple list * simple list
  | Ptr of qual * simple
  | Array of simple
  | Struct of (int * simple) list ref
  | Var of simple var_type ref

```

FIGURE 7.8: Représentation des types

Les types de données utilisés sont donnés dans la figure 7.8. Les types SAFESPEAK sont représentés soit par un constructeur de type “résolu” immutable comme Int, soit par une référence dans le cas d’une variable inconnue (placée alors derrière le constructeur Var).

Ces références contiennent une valeur du type simple var_type¹, c’est-à-dire :

- soit un numéro d’inconnue (constructeur Unknown)².
- soit un type résolu (constructeur Instanciated) si cette inconnue a été unifiée avec un type concret.

Ce système peut créer des représentations de types arbitrairement longues, comme par exemple :

```
Var (ref (Instanciated (Var (ref (Instanciated Int))))))
```

Cela est dû au fait que fun x -> Var (ref (Instanciated x)) est typée simple -> simple et peut donc être appliquée à loisir. Pour éviter cet effet d’allongement, on définit une fonction shorten qui supprime ces chaînes (figure 7.9).

La fonction d’unification, quant à elle, commence par raccourcir ses deux arguments puis faire une analyse de cas par filtrage.

Les cas principaux pour unifier deux types réduits sta et stb sont :

- sta et stb sont deux inconnues. Alors on modifie l’un pour pointer sur l’autre. Les références étant uniques et partagées, cela revient à substituer l’un par l’autre dans toutes les représentations de types.
- Un type est concret, l’autre une inconnue. Dans ce cas on modifie le second comme étant un Instanciated du premier. Il faut vérifier que le type inconnu remplacé n’apparaît pas dans le type concret, sinon on crée un cycle. Ce cas a lieu par exemple quand

1. Ce type est polymorphe pour pouvoir être repris dans l’unification des qualificateurs (figure 7.12).

2. On place le numéro dans un enregistrement pour abstraire l’entier sous-jacent, en empêchant par exemple de faire de l’arithmétique sur celui-ci.

```

let rec shorten = function
| Var ({contents = Instanciated (Var _ as t)} as vt) ->
    let t2 = shorten t in
    vt := Instanciated t;
    t2
| Var {contents = Instanciated t} -> t
| t -> t

```

FIGURE 7.9: Fonction de raccourcissement des représentations de types

on cherche à unifier t avec $t *$. Il faut alors signaler une erreur, ce que fait la fonction `occurs_check_failed`.

- Les deux types sont des types concrets. Alors ils sont de la forme respective $C(t_1, \dots, t_n)$ et $D(u_1, \dots, u_m)$ où C et D sont des constructeurs de type avec respectivement n et m arguments³. Si $C = D$ et $n = m$, alors on unifie récursivement chaque t_i avec u_i . Sinon on lève une erreur (fonction `type_clash`).

Cette analyse de cas est implantée dans la fonction `unify_now` (figure 7.10).

`cas structure` (figure 7.11).

Pour les pointeurs, il est nécessaire de définir des fonctions similaires sur les qualificateurs (figure 7.12).

La fonction appelée directement par le reste du code, appelée `unify`, peut retarder l'unification (figure 7.13). Dans ce cas, la paire de types à unifier est mise dans une liste d'attente qui sera unifiée après le parcours du programme. Le but est d'instrumenter l'inférence de types afin de pouvoir en faire une exécution “pas à pas”.

7.6 Exemple

Lançons l'analyse sur un petit exemple :

```
int f(int *x) { return (*x + 1); }
```

L'exécution de notre analyseur affiche un programme complètement annoté :

```

% ptrtype example.c
f : (Ptr (Int)) -> (Int)
Int (example.c:1#4)^f(Ptr (Int) x) {
  (.c:3#4)^!return =(int32)
  (coerce[-2147483648,2147483647]
    ( ( [(x_Ptr (Int) : Ptr (Int))]32_Int
      : Int
    )
    + (1 : Int)
  ) : Int
) : Int
};
}

```

3. Le cas $n = 0$ ou $m = 0$ correspond aux types concrets comme `INT` ou `FLOAT`.

```

let rec unify_now ta tb =
  let sta = shorten ta in
  let stb = shorten tb in
  match (sta, stb) with
  | ((Var ({contents = Unknown na} as ra)),
      (Var ({contents = Unknown nb}))) ->
    begin
      if na <> nb then
        ra := Instanciated stb
      end
    | ((Var ({contents = Unknown {id = n}} as r)), t)
      ->
        begin
          if occurs n t then
            occurs_check_failed sta stb
          else
            r := Instanciated t
          end
        | (_, (Var ({contents = Unknown _}))) -> unify_now stb sta

        | Int, Int
        | Float, Float -> ()

        | Ptr (qa, ta), Ptr (qb, tb) ->
          unify_equals qa qb;
          unify_now ta tb

        | Array a, Array b -> unify_now a b

        | Fun (args_a, rets_a), Fun (args_b, rets_b) ->
          List.iter2 unify_now args_a args_b;
          List.iter2 unify_now rets_a rets_b

        | Struct rfa, Struct rfb ->
          unify_structs rfa rfb

        | _ -> type_clash sta stb

```

FIGURE 7.10: Fonction d'unification

```

let unify_structs rfa rfb =
  let fa = !rfa in
  let fb = !rfb in

  let new_a = ref [] in
  let new_b = ref [] in

  let unify_fields = function
    | _, InBoth (ta, tb) -> unify_now ta tb
    | k, OnlyL f -> new_b := (k,f) :: !new_b
    | k, OnlyR f -> new_a := (k,f) :: !new_a
  in

  List.iter unify_fields (compare_lists fa fb);
  let by_offset (x, _) (y, _) =
    compare x y
  in
  rfa := List.sort by_offset (!new_a @ !rfa);
  rfb := List.sort by_offset (!new_b @ !rfb)

```

FIGURE 7.11: Structures

L'opérateur `coerce[a,b]` est un artefact de Newspeak, destiné à détecter les débordements d'entiers lors d'une analyse de valeurs par interprétation abstraite. Dans le cas de notre analyse, les valeurs ne sont pas pertinentes et cet opérateur peut être vu que comme un "plus" unaire typé ($\text{INT} \rightarrow \text{INT}$).

A contrario, lorsqu'il n'est pas possible d'inférer des types compatibles, l'analyseur s'arrête avec une erreur.

```

void f(int *p) {
  /*!nkp userptr p */
  *p = 3;
}

```

Le commentaire `/*!nkp userptr p */` est interprété par l'analyseur et le fait unifier le type de `p` avec `t @`, c'est-à-dire qu'il force son qualificateur à être `USER`.

```

04-addrof.c:4#4 - Cannot unify qualifiers:
Kernel
User

```

On signale l'emplacement où la dernière unification a échoué.

```

let rec shorten_q = function
| QVar ({contents = Instanciated (QVar _ as t)} as vt) ->
  let t2 = shorten_q t in
  vt := Instanciated t;
  t2
| QVar {contents = Instanciated t} -> t
| t -> t

let rec unify_qual a b =
  let sa = shorten_q a in
  let sb = shorten_q b in
  match (sa, sb) with
  | User, User -> ()
  | Kernel, Kernel -> ()
  | ((QVar ({contents = Unknown na} as ra)),
    (QVar {contents = Unknown nb})) ->
    begin
      if na <> nb then
        ra := Instanciated sb
      end
    end
  | ((QVar ({contents = Unknown {id = n}} as r)), q)
  ->
    begin
      if occurs_q n q then
        occurs_q_check_failed sa sb
      else
        r := Instanciated q
      end
    end
  | (_, (QVar ({contents = Unknown _}))) -> unify_qual sb sa
  | _ -> Utils.error "Cannot unify qualifiers:\n %s\n %s\n"
    (string_of_qual sa)
    (string_of_qual sb)

```

FIGURE 7.12: Qualificateurs

```

let unify a b =
  if !Options.lazy_unification then
    Queue.add (Unify (a, b)) unify_queue
  else
    unify_now a b

```

FIGURE 7.13: Unification directe ou retardée

ÉTUDE DE CAS : UN PILOTE DE CARTE GRAPHIQUE

Dans ce chapitre, un exemple de mise en œuvre du système de type décrit dans le chapitre 6 et implanté dans le chapitre 7.

8.1 Linux

Le noyau Linux, abordé dans le chapitre 2, est un noyau de système d'exploitation développé depuis le début des années 90 et "figure de proue" du mouvement open-source. Au départ écrit par Linus Torvalds sur son ordinateur personnel, il a au fil des années été porté sur de nombreuses architectures et s'est enrichi de nombreux pilotes de périphériques. Dans la version 3.9.8 (2013), son code source comporte 11.2 millions de lignes de code (en grande majorité du C) dont 57% de pilotes.

Même si le noyau est monolithique (la majeure partie des traitements s'effectue au sein d'un même fichier objet), les sous-systèmes sont indépendants. C'est ce qui permet d'écrire des pilotes de périphériques et des modules.

8.2 GNU C

Linux est écrit dans le langage C, mais pas sa version normalisée. Il utilise le dialecte GNU C qui est celui que supporte GCC.

Pour traduire ce dialecte, il a été nécessaire d'adapter `c2newspeak`. La principale particularité est la notation `__attribute__((...))` qui peut décorer les déclarations de fonctions, de variables ou de types. De nouvelles fonctionnalités sont aussi présentes.

Par exemple, il est possible de manipuler des étiquettes de première classe : si `lbl` est présent avant une instruction, on peut capturer l'adresse de celle-ci avec `void *p = &lbl` et y sauter indirectement avec `goto *p`.

Une autre fonctionnalité est le concept d'instruction-expression : `{(bloc)}` est une expression, dont la valeur est celle de la dernière expression évaluée lors de `i`.

Les attributs, quant à eux, rentrent dans trois catégories :

- les annotations de compilation ; par exemple, `used` désactive l'avertissement "cette variable n'est pas utilisée".
- les optimisations ; par exemple, les objets marqués `hot` sont groupés de telle manière qu'ils se retrouvent en cache ensemble.
- les annotations de bas niveau ; par exemple, `aligned(n)` spécifie qu'un objet doit être aligné sur au moins `n` bits.

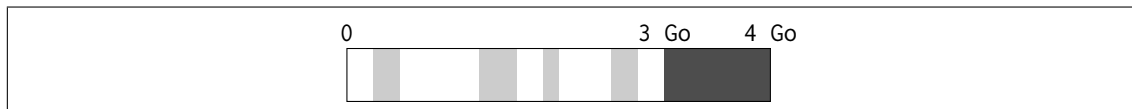


FIGURE 8.1: L'espace d'adressage d'un processus. En gris clair, les zones accessibles à tous les niveaux de privilèges : code du programme, bibliothèques, tas, pile. En gris foncé, la mémoire du noyau, réservée au mode privilégié.

Dans notre cas, toutes ces annotations peuvent être ignorées.

8.3 Configuration

Pour que le code noyau soit compilable, il est nécessaire de définir certaines macros. En particulier, le système de configuration de Linux utilise des macros nommées `CONFIG_*` pour inclure ou non certaines fonctionnalités. Il a donc fallu faire un choix ; nous avons choisi la configuration par défaut. Pour analyser des morceaux plus importants du noyau, il faudrait définir un fichier de configuration plus important.

8.4 Appels systèmes sous 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 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 8.1), et donc chacun un ensemble de tables 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.

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 [Int10].

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)
    RING0_INT_FRAME                # can't unwind into user space anyway
    pushl %eax                     # save orig_eax
    CFI_ADJUST_CFA_OFFSET 4
    SAVE_ALL
    GET_THREAD_INFO(%ebp)

                                # system call tracing in operation / emulation
    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` :

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*.

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

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.

8.5 Bug

On décrit le cas d'un pilote video qui contenait un bug de pointeur utilisateur. Il est répertorié sur <http://freedesktop.org> en tant que bug #29340.

Pour changer de mode graphique, les pilotes de GPU peuvent supporter le *Kernel Mode Setting* (KMS).

Pour configurer un périphérique, l'utilisateur communique avec le pilote noyau avec le mécanisme d'*ioctl*². Ils sont similaires à des appels système, mais spécifique à un périphérique particulier. Les fonctions implantant un *ioctl* sont donc vulnérables à la même classe d'attaques que les appels système, et donc doivent être écrits avec une attention particulière.

Le code suivant est présent dans le pilote KMS pour les GPU AMD Radeon :

```
/* 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;
```

2. Ce nom vient de la fonction `ioctl()` pour *Input/Output Control*.

```

--- a/drivers/gpu/drm/radeon/radeon_kms.c
+++ b/drivers/gpu/drm/radeon/radeon_kms.c
@@ -112,7 +112,9 @@

    info = data;
    value_ptr = (uint32_t *)((unsigned long)info->value);
-   value = *value_ptr;
+   if (DRM_COPY_FROM_USER(&value, value_ptr, sizeof(value)))
+       return -EFAULT;
+
    switch (info->request) {
    case RADEON_INFO_DEVICE_ID:
        value = dev->pci_device;

```

FIGURE 8.2: Patch résolvant le problème de pointeur utilisateur.

```

    [...]
}

```

On peut voir que l'argument `data` est converti en un struct `drm_radeon_info *`. Un pointeur `value_ptr` est extrait de son champ `value`, et finalement ce pointeur est déréféréncé.

Cependant, l'argument `data` est un pointeur vers une structure (allouée en espace noyau) du type suivant, dont les champs proviennent d'un appel utilisateur de `ioctl()`.

```

/* from include/drm/radeon_drm.h */
struct drm_radeon_info {
    uint32_t    request;
    uint32_t    pad;
    uint64_t    value;
};

```

Pour mettre ce problème en évidence, nous avons annoté la fonction `radeon_info_ioctl` de telle manière que son second paramètre soit un pointeur noyau vers une structure contenant un champ `USER`, `value`. Ceci est possible puisqu'avant la traduction, on efface les types présents dans le programme C. Ainsi, un pointeur ne peut pas être distingué d'un entier transtypé en un pointeur. Avec cette configuration, nous obtenons une erreur de type à la ligne 16.

L'intégralité de ce code peut être trouvée en annexe A.

La bonne manière de faire a été publiée avec le numéro de *commit* `d8ab3557` (figure 8.2) (`DRM_COPY_FROM_USER` étant une simple macro pour `copy_from_user`). Dans ce cas, on n'obtient pas d'erreur de typage.

8.6 Détails

Pour utiliser notre système de types, plusieurs étapes sont nécessaires en plus de traduire le noyau linux en SAFESPEAK.

Afin de réaliser l'analyse, il faut annoter les sources pour créer un environnement initial (via la variable `exttbl` décrite en section 7.3).

Ensuite, il faut réécrire les fonctions de manipulation de pointeurs fournies par le noyau : `get_user`, `put_user`, `copy_from_user`, `copy_to_user`, etc. Leur implantation revient à réaliser un test puis à faire la copie. À leur place on utilise les primitives `copy_from_user(·, ·)` et `copy_to_user(·, ·)`.

Enfin, on peut lancer l'inférence de type.



Quatrième partie

Conclusion

CONCLUSION

On fait ici un bilan des travaux présentés, en commençant par un bilan des contributions réalisées. On fait ensuite un tour des aspects posant problème, ou traités de manière incomplète, en évoquant les travaux possibles pour enrichir l'expressivité de ce système.

9.1 Contributions

Un langage impératif bien typable Le système de types de C est trop rudimentaire pour permettre de d'obtenir des garanties sur l'exécution des programmes bien typés. En interdisant certaines constructions dangereuses et en en annotant certaines, nous avons isolé un langage impératif bien typable, pour lequel on peut définir un système de types sûr.

Une sémantique basée sur les lentilles Une des particularités de ce langage impératif est qu'il utilise un état mémoire structuré, modélisant les cadres de piles présents dans le langage. Pour décrire la sémantique des accès mémoire, nous utilisons le concept de lentilles issues de la programmation fonctionnelle et des systèmes de base de données. Cela permet de définir de manière déclarative la modification en profondeur de valeurs dans la mémoire, sans avoir à distinguer le cas de la lecture et celui de l'écriture.

Un système de types abstraits En partant de ce système de types, on a décrit une extension permettant de créer des pointeurs pour lesquels l'opération de déréréfencement est restreinte à certaines fonctions. Dans le contexte d'un noyau de système d'exploitation, cette restriction permet de vérifier statiquement qu'à aucun moment le noyau ne déréréfence un pointeur dont la valeur est contrôlée par l'espace utilisateur, évitant ainsi un problème de sécurité. Cette approche peut s'étendre à d'autres classes de problèmes comme par exemple éviter l'utilisation de certaines opérations sur les types entiers lorsqu'ils sont utilisés comme identificateurs ou masque de bits.

9.2 Limitations et travaux futurs

Arithmétique En C, on dispose de plusieurs types entiers, pouvant avoir plusieurs tailles et être signés ou non signés, ainsi que des types flottants qui diffèrent par leur taille. Au contraire, en SAFESPEAK on ne conserve qu'un seul type d'entier. La raison pour cela est que nous ne nous intéressons pas du tout aux problématiques de sémantique entière : les débordements, dénormalisations, etc, sont supposés ne pas arriver.

Comme en C il n'est pas possible de définir de nouveaux types de base (des entiers de 24 bits par exemple), il est facile d'étendre le système de types de SAFESPEAK pour ajouter tous ces nouveaux types. La traduction depuis Newspeak insère déjà des opérateurs de transtypage pour lesquels il est facile de donner une sémantique (pouvant lever une erreur en cas de débordement, comme en Ada) et un typage. Les littéraux numériques peuvent poser problème, puisqu'ils deviennent alors polymorphes. Une solution peut être de leur donner le plus grand type entier et d'insérer un opérateur de transtypage à chaque littéral. Haskell utilise une solution similaire : les littéraux entier ont le type de précision arbitraire `Integer` et est converti dans le bon type en appelant la fonction `fromInteger` du type synthétisé à partir de l'environnement.

Transtypage Puisque l'approche retenue est basée sur le typage statique, il est impossible de capturer de nombreuses constructions qui sont possibles, ou même idiomatiques en C : les unions, les conversions de types (explicites ou implicites) et le *type punning* (défini ci-dessous). Les deux premières sont équivalentes. Bien qu'on puisse remplacer chaque conversion explicite d'un type t_1 vers un type t_2 par l'appel à une fonction `castt1,t2`, on ajoute alors un "trou" dans le système de types. Cette fonction devrait en effet être typée $(t_1) \rightarrow t_2$, autrement dit le type "maudit" $\alpha \rightarrow \beta$ de `Obj.magic` en OCaml ou `unsafeCoerce` en Haskell.

Le *type punning* consiste à modifier directement la suite de bits de certaines données pour la manipuler d'une manière efficace. Par exemple, il est commun de définir un ensemble de macros pour accéder à la mantisse et à l'exposant de flottants IEEE754. Ceci peut être fait avec des unions ou des masques de bits.

Dans de tels cas, le typage statique est bien sûr impossible. Pour traiter ces cas, il faudrait encapsuler la manipulation dans une fonction et y ajouter une information de type explicite, comme `float_exponent : (FLOAT) → INT`.

En fait, les casts entre types entiers ne posent pas de problèmes. Ceux entre types pointeurs ne posent pas de problème non plus puisqu'au moment de la traduction en SAFESPEAK on efface les étiquettes de types sur les pointeurs.

Mais les casts entre entiers et pointeurs posent problème. Si on autorise les casts, on peut néanmoins trouver la solution suivante.

$$\frac{\Gamma \vdash e : t *}{\Gamma \vdash (\text{INT}) e : \text{INT}} \text{ (PTRINT-BAD)} \qquad \frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash (\text{PTR}) e : t * } \text{ (INTPTR-BAD)}$$

Cela pose problème car il est alors possible de créer une fonction qui convertit n'importe quel type pointeur en n'importe quel autre type pointeur :

$$\vdash \text{fun}(p)\{\text{RETURN}((\text{PTR}) (\text{INT}) p)\} \\ : (t_a *) \rightarrow t_b *$$

Si on crée une variable du type t_a , prend son adresse, la convertit à l'aide de cette fonction, puis déréférence le résultat, on obtient une valeur du type t_b (remarquons que ce genre d'opération est tout à fait possible en C).

Outre ce problème de typage, il faudrait pouvoir donner une sémantique à ces opérations. Convertir un pointeur en entier revient à spécifier l'environnement d'exécution, c'est-à-dire qu'il faut une fonction de placement en mémoire beaucoup plus précise que notre modèle mémoire actuel. Celle-ci dépend de beaucoup de paramètres : dans quel sens croît la pile, quelle est la taille des types, etc.

La conversion dans le sens inverse, d'entier vers pointeur, est encore plus complexe. D'une part cela suppose qu'on puisse retrouver la taille des valeurs à partir de leur adresse. Dans de nombreux langages, on résout ce problème en ajoutant la taille de chaque valeur avec elle.

Mais cela fait s'éloigner du modèle mémoire de C, où le déréférencement porte sur une adresse mais également sur une taille (portée implicitement par le type du pointeur). D'ailleurs le langage NEWSPEAK conserve cette distinction, que nous avons éliminée dans SAFESPEAK. Il y a une incompatibilité entre ces deux approches : pour pouvoir typer la mémoire, il faut associer un type à chaque zone mémoire, alors qu'avec une sémantique de bas niveau les valeurs sont créées à la volée par le déréférencement.

Environnement d'exécution La sémantique opérationnelle utilise un environnement d'exécution pour certains cas. Contrairement à C, les débordements de tampon et les déréférencements de pointeurs sont vérifiés dynamiquement. Mais ce n'est pas une caractéristique cruciale de cette approche : en effet, si on suppose que les programmes que l'on analyse sont corrects de ce point de vue, on peut désactiver ces vérifications et le reste des propriétés est toujours valable.

Un autre endroit, plus problématique, où des tests dynamiques sont faits est lorsqu'on recherche en mémoire des pointeurs référençant un cadre d'appel qui n'est plus valide (à l'aide de l'opérateur Cleanup(\cdot)). Supprimer ce test rend l'analyse incorrecte, car il est alors possible de faire référence à une variable avec un type différent.

De même, si on peut avoir une garantie statique que les adresses des variables locales ne seront plus accessibles au retour d'une fonction, alors on peut supprimer le nettoyage en posant Cleanup(m) = m . Cette garantie peut être obtenue avec une analyse statique préalable [DDMPn10].

Allocation dynamique La plupart des programmes, et le noyau Linux en particulier, utilisent la notion d'allocation dynamique de mémoire. C'est une manière de créer dynamiquement une zone de mémoire qui restera accessible après l'exécution de la fonction courante. Cette mémoire pourra être libérée à l'aide d'une fonction dédiée. Dans l'espace utilisateur, les programmes peuvent utiliser les fonctions malloc(), calloc(), et realloc() pour allouer des zones de mémoire et free() pour les libérer. Dans le noyau Linux, ces fonctions existent sous la forme de kmalloc(), kfree(), etc. Une explication détaillée de ces mécanismes peut être trouvée dans [Gor04].

Ces fonctions manipulent les données en tant que zones mémoires opaques, en renvoyant un pointeur vers une zone mémoire d'un nombre d'octets donnés. Cela présuppose un modèle mémoire plus bas niveau. Pour se rapprocher de la sémantique de SAFESPEAK, une manière de faire est de définir un opérateur de plus haut niveau prenant une expression et retournant l'adresse d'une cellule mémoire contenant cette valeur (la taille de chaque valeur fait partie de celle-ci). Le typage est alors direct :

$$\frac{\Gamma \vdash e : t}{\Gamma \vdash \text{ALLOC}(e) : t *} \text{ (MALLOC)} \qquad \frac{\Gamma \vdash e : t *}{\Gamma \vdash \text{FREE}(e)} \text{ (FREE)}$$

En ce qui concerne l'exécution, on peut ajouter une troisième composante aux états mémoire : $m = (s, g, h)$ où h associe à une étiquette numérique une valeur (comme les globales mais indexées par des entiers plutôt que par leur nom). La libération de mémoire est problématique parce qu'il faut faire confiance au programmeur pour ne pas accéder aux zones mémoires libérées. Encore une fois, on peut utiliser une analyse préalable comme [DDMPn10]. Il est sûr (mais pas possible en pratique) d'ignorer les commandes de libérations de mémoire.

Structures récursives Une autre limitation est que seules les structures simples sont supportées. Les structures récursives, contenant un pointeur vers un objet du même type structure, ne sont pas typables. La raison est syntaxique : l'accès à un champ est noté $lv.l_s$ afin d'avoir à disposition le type complet de la structure.

Par exemple, une liste chaînée d'entiers peut être décrite de la manière suivante en C :

```
struct int_list {
    int value;
    struct int_list *next;
};
```

La compilation est possible car la définition complète de `struct int_list` n'est pas nécessaire pour connaître la taille d'un pointeur vers un objet de ce type.

Pour compiler `l.next` en SAFESPEAK, il faudrait écrire le type en extension :

$$l.next_{value:INT;next:\{value:INT;next:...\} *}$$

On peut fermer cette récursion infinie en ajoutant directement un opérateur de point fixe qui rend accessible le type d'un pointeur vers le type structure en train d'être défini.

$$S ::= \dots \mid \text{fix}(p \rightarrow S)$$

On écrit alors :

$$l.next_{\text{fix}(p \rightarrow \{value:INT;next:p\})}$$

Les structures mutuellement récursives peuvent être également exprimées de cette manière.

```
struct a {
    int value;
    struct b *next;
};

struct b {
    float value;
    struct a *next;
};
```

Ces structures sont respectivement traduites en :

$$\begin{aligned} A &= \text{fix}(p \rightarrow \{value:INT;next:\{value:FLOAT;next:p\} * \}) \\ B &= \text{fix}(p \rightarrow \{value:FLOAT;next:\{value:INT;next:p\} * \}) \end{aligned}$$

Ce schéma suffit à compiler toutes les types de structures possibles en C. En effet, les seuls cas posant problème arrivent lorsqu'on fait référence à la structure en cours de définition ; et il n'est alors possible que d'y accéder par pointeur.

Assembleur Comme la plupart des outils d'analyse de code C, il est impossible de traiter l'assembleur en ligne qui peut se trouver entre deux instructions. Dans le cas de Linux, le code est fait pour être portable, et les parties dépendantes d'une certaine architecture (et donc le code assembleur) sont isolées explicitement. On peut alors si nécessaire ajouter une annotation de type sur une fonction dont l'implantation est faite en assembleur, mais au sein de cette implantation on ne peut bien sûr rien dire.

Analyse du noyau Linux Ici nous avons présenté l'analyse expérimentale d'une fonction problématique d'une interface programmation particulière. Mais le principe de l'analyse est applicable à toutes les fonctions de traitement d'octets des pilotes KMS, et même à toutes les fonctions faisant partie des différentes interfaces recevant un pointeur de l'espace utilisateur.

Autres types abstraits Notre approche est basée sur le fait de rendre abstrait un type de C et d'y interdire certaines opérations : ici, on marque un certain type de pointeur comme "utilisateur" et on interdit l'opérateur `*` dessus.

Le langage C n'ayant pas de types abstraits, on ne peut pas séparer la représentation d'un type (par exemple : entier signé de 32 bits) des opérations qui y sont attachées. La seule abstraction possible est lorsqu'on manipule une structure par pointeur. Il n'est alors pas nécessaire de connaître sa définition totale. L'idiome est alors de placer uniquement une déclaration en avance (`struct s;`) dans l'en-tête (`.h`) et de renseigner la définition complète dans l'implantation (`.c`). Cette technique de *pointeurs opaques* n'est pas applicable aux autres types.

Dans de nombreuses interfaces, on emploie des types entiers qui servent d'étiquettes. Par exemple, les descripteurs de fichiers renvoyés par la fonction `open()` et passés aux fonctions `read()` et `write()` ont le type `int`. Le langage autorise donc par exemple de multiplier entre eux deux descripteurs de fichiers, ce qui ne correspond pas à une opération concrète.

On peut aussi distinguer plusieurs types abstraits entre eux. Par exemple, si un encodeur vidéo manipule des numéros de *frame* et des identifiants de codec tous les deux entiers, on peut interdire d'utiliser avec un identifiant de codec une fonction prenant en paramètre un numéro de *frame*.

Annexes



MODULE RADEON KMS

```

/* 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 include/drm/radeon_drm.h */
struct drm_radeon_info {
    uint32_t          request;

```

```

        uint32_t          pad;
        uint64_t          value;
};

/* 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

$\frac{\langle i, m \rangle \rightarrow \langle i', m' \rangle}{\langle C\langle i \rangle, m \rangle \rightarrow \langle C\langle i' \rangle, m' \rangle} \text{ (CTX)}$	
Contextes	$C ::= \bullet$ $ C \boxplus e$ $ \nu \boxplus C$ $ \boxplus C$ $ C \leftarrow e$ $ \varphi \leftarrow C$ $ \{l_1 : \nu_1; \dots; l_i : C; \dots; l_n : e_n\}$ $ [\nu_1; \dots; C; \dots; e_n]$ $ C(e_1, \dots, e_n)$ $ f(\nu_1, \dots, C, \dots, e_n)$ $ C.l_s$ $ C[e]$ $ \varphi[C]$ $ C; i$ $ \text{IF}(C)\{i_1\}\text{ELSE}\{i_2\}$ $ \text{RETURN}(C)$ $ \text{DECL } x = C \text{ IN}\{i\}$

FIGURE B.1: Règles d'évaluation des contextes

		$\Xi \rightarrow \Omega$
$\frac{}{\langle \Omega, m \rangle \rightarrow \Omega} \text{ (EXP-ERR)}$	$\frac{\langle e, m \rangle \rightarrow \Omega}{\langle C\langle e \rangle, m \rangle \rightarrow \Omega} \text{ (EVAL-ERR)}$	

FIGURE B.2: Règles d'évaluation des erreurs

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

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

$$\frac{\widehat{\&} \varphi' = m[\varphi]_{\Phi}}{\langle * \varphi, m \rangle \rightarrow \langle \varphi', m \rangle} \text{ (PHI-DEREF)}$$

$$\frac{}{\langle lv.l_S, m \rangle \rightarrow \langle lv.\widehat{l}, m \rangle} \text{ (PHI-STRUCT)}$$

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

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

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

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

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

$$\frac{}{\langle \boxminus v, m \rangle \rightarrow \langle \widehat{\boxminus} v, m \rangle} \text{ (EXP-UNOP)}$$

$$\frac{}{\langle v_1 \boxplus v_2, m \rangle \rightarrow \langle v_1 \widehat{\boxplus} v_2, m \rangle} \text{ (EXP-BINOP)}$$

$$\frac{}{\langle \& \varphi, m \rangle \rightarrow \langle \widehat{\&} \varphi, m \rangle} \text{ (EXP-ADDR OF)}$$

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

$$\frac{}{\langle \{l_1 : v_1; \dots; l_n : v_n\}, m \rangle \rightarrow \langle \{l_1 : \widehat{v_1}; \dots; l_n : \widehat{v_n}\}, m \rangle} \text{ (EXP-STRUCT)}$$

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

$$\frac{m_1 = \text{Push}(m_0, (a_1 \mapsto v_1, \dots, a_n \mapsto v_n)) \quad \langle i, m_1 \rangle \rightarrow \langle i', m_2 \rangle \quad \forall i \in [1; n], v'_i = m_2[(0, a_i)]_A \quad m_3 = \text{Pop}(m_2)}{\langle \text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n), m_0 \rangle \rightarrow \langle \text{fun}(a_1, \dots, a_n)\{i'\}(v'_1, \dots, v'_n), m_3 \rangle} \text{ (EXP-CALL-CTX)}$$

$$\frac{m' = \text{Push}(m, (a_1 \mapsto v_1, \dots, a_n \mapsto v_n)) \quad \langle i, m' \rangle \rightarrow \Omega}{\langle \text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n), m \rangle \rightarrow \Omega} \text{ (EXP-CALL-ERR)}$$

$$\frac{m' = \text{Cleanup}(m)}{\langle \text{fun}(a_1, \dots, a_n)\{\text{PASS}\}(v_1, \dots, v_n), m \rangle \rightarrow \langle (), m' \rangle} \text{ (EXP-CALL-RETURN-UNIT)}$$

$$\frac{m' = \text{Cleanup}(m)}{\langle \text{fun}(a_1, \dots, a_n)\{\text{RETURN}(v)\}(v_1, \dots, v_n), m \rangle \rightarrow \langle v, m' \rangle} \text{ (EXP-CALL-RETURN)}$$

FIGURE B.3: Règles d'évaluation des left-values et expressions

$\langle i, m \rangle \rightarrow \langle i', m \rangle$	
$\frac{}{\langle \text{PASS}; i, m \rangle \rightarrow \langle i, m \rangle} \text{ (SEQ)}$	$\frac{}{\langle v, m \rangle \rightarrow \langle \text{PASS}, m \rangle} \text{ (EXP)}$
$\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{PASS}\}, m \rangle \rightarrow \langle \text{PASS}, m' \rangle} \text{ (DECL-PASS)}$	
$\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{RETURN}(v')\}, m \rangle \rightarrow \langle \text{RETURN}(v'), m' \rangle} \text{ (DECL-RETURN)}$	
$\frac{\langle i, m' \rangle \rightarrow \langle i', m'' \rangle \quad v' = m''[(0, x)]_A \quad m''' = m'' - x}{\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \langle \text{DECL } x = v' \text{ IN } \{i'\}, m''' \rangle} \text{ (DECL-CTX)}$	
$\frac{\langle i, m \rangle \rightarrow \Omega}{\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \Omega} \text{ (DECL-ERR)}$	$\frac{}{\langle \text{IF}(0)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_f, m \rangle} \text{ (IF-FALSE)}$
$\frac{v \neq 0}{\langle \text{IF}(v)\{i_t\}\text{ELSE}\{i_f\}, m \rangle \rightarrow \langle i_t, m \rangle} \text{ (IF-TRUE)}$	
$\frac{}{\langle \text{WHILE}(e)\{i\}, m \rangle \rightarrow \langle \text{IF}(e)\{i; \text{WHILE}(e)\{i\}\}, m \rangle} \text{ (WHILE)}$	
$\frac{}{\langle \text{RETURN}(v); i, m \rangle \rightarrow \langle \text{RETURN}(v), m \rangle} \text{ (RETURN)}$	
$m \Vdash p \rightarrow m'$	
$\frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{m \Vdash e \rightarrow m'} \text{ (T-EXP)}$	$\frac{\langle e, m \rangle \rightarrow \langle v, m' \rangle}{(s, g) \Vdash x = e \rightarrow (s, (x \mapsto v) :: g)} \text{ (T-VAR)}$
$m \Vdash p \rightarrow^* m'$	
$\frac{}{m \Vdash [] \rightarrow^* m} \text{ (T*-NIL)}$	$\frac{m \Vdash p \rightarrow m' \quad m' \Vdash ps \rightarrow^* m''}{m \Vdash p :: ps \rightarrow^* m''} \text{ (T*-CONS)}$
$\Vdash P \rightarrow^* m$	
$\frac{([], []) \Vdash P \rightarrow^* m}{\Vdash P \rightarrow^* m} \text{ (PROG)}$	

FIGURE B.4: Règles d'évaluation des instructions et phrases

$$\begin{array}{c}
\frac{}{\langle \diamond \varphi, m \rangle \rightarrow \langle \widehat{\diamond} \varphi, m \rangle} \text{ (EXPR-TAINTED)} \\
\\
\frac{\langle \varphi_d \leftarrow \varphi_s, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle \text{copy_from_user}(\varphi_d, \widehat{\diamond} \varphi_s), m \rangle \rightarrow \langle 0, m' \rangle} \text{ (USER-GET-OK)} \\
\\
\frac{\nexists \varphi_s. \varphi = \widehat{\diamond} \varphi_s}{\langle \text{copy_from_user}(\varphi_d, \varphi), m \rangle \rightarrow \langle -1, m \rangle} \text{ (USER-GET-ERR)} \\
\\
\frac{\langle \varphi_d \leftarrow \varphi_s, m \rangle \rightarrow \langle \text{PASS}, m' \rangle}{\langle \text{copy_to_user}(\widehat{\diamond} \varphi_d, \varphi_s), m \rangle \rightarrow \langle 0, m' \rangle} \text{ (USER-PUT-OK)} \\
\\
\frac{\nexists \varphi_d. \varphi = \widehat{\diamond} \varphi_d}{\langle \text{copy_to_user}(\varphi, \varphi_s), m \rangle \rightarrow \langle -1, m \rangle} \text{ (USER-PUT-ERR)}
\end{array}$$

FIGURE B.5: Règles d'évaluation des extensions noyau

RÈGLES DE TYPAGE

			$\boxed{\Gamma \vdash c : t}$
$\frac{}{\Gamma \vdash n : \text{INT}} \text{ (CST-INT)}$	$\frac{}{\Gamma \vdash d : \text{FLOAT}} \text{ (CST-FLOAT)}$	$\frac{}{\Gamma \vdash \text{NULL} : t*} \text{ (CST-NULL)}$	
$\frac{}{\Gamma \vdash () : \text{UNIT}} \text{ (CST-UNIT)}$			
			$\boxed{\Gamma \vdash lv : t}$
$\frac{x : t \in \Gamma}{\Gamma \vdash x : t} \text{ (LV-VAR)}$	$\frac{\Gamma \vdash lv : t*}{\Gamma \vdash *lv : t} \text{ (LV-DEREF)}$	$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash lv : t[]}{\Gamma \vdash lv[e] : t} \text{ (LV-INDEX)}$	
$\frac{(l, t) \in S \quad \Gamma \vdash lv : S}{\Gamma \vdash lv.l_S : t} \text{ (LV-FIELD)}$			

FIGURE C.1: Règles de typage des constantes et variables

		$\boxed{\Gamma \vdash \boxminus e : t}$
$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash +e : \text{INT}}$	(UNOP-PLUS-INT)	$\frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash +.e : \text{FLOAT}}$
		(UNOP-PLUS-FLOAT)
$\frac{\Gamma \vdash e : \text{INT}}{\Gamma \vdash -e : \text{INT}}$	(UNOP-MINUS-INT)	$\frac{\Gamma \vdash e : \text{FLOAT}}{\Gamma \vdash -.e : \text{FLOAT}}$
		(UNOP-MINUS-FLOAT)
$\frac{\boxminus \in \{\sim, !\} \quad \Gamma \vdash e : \text{INT}}{\Gamma \vdash \boxminus e : \text{INT}}$		(UNOP-NOT)
		$\boxed{\Gamma \vdash e_1 \boxplus e_2 : t}$
$\frac{\boxplus \in \{+, -, \times, /, \&, , \wedge, \&\&, , \ll, \gg, \leq, \geq, <, >\}}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}}$		(OP-INT)
$\frac{\boxplus \in \{+., -., \times., /., \leq., \geq., <., >.\}}{\Gamma \vdash e_1 \boxplus e_2 : \text{FLOAT}}$		(OP-FLOAT)
$\frac{\boxplus \in \{=, \neq\} \quad \Gamma \vdash e_1 : t \quad \Gamma \vdash e_2 : t \quad \text{EQ}(t)}{\Gamma \vdash e_1 \boxplus e_2 : \text{INT}}$		(OP-EQ)
$\frac{\boxplus \in \{+_p, -_p\} \quad \Gamma \vdash e_1 : t* \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxplus e_2 : t*}$		(PTR-ARITH)
		$\boxed{\text{EQ}(t)}$
$\frac{t \in \{\text{INT}, \text{FLOAT}\}}{\text{EQ}(t)}$	(EQ-NUM)	$\frac{}{\text{EQ}(t*)}$
		(EQ-PTR)
		$\frac{\text{EQ}(t)}{\text{EQ}(t[])}$
		(EQ-ARRAY)
$\frac{\forall i \in [1; n]. \text{EQ}(t_i)}{\text{EQ}(\{l_1 : t_1; \dots l_n : t_n\})}$		(EQ-STRUCT)

FIGURE C.2: Règles de typage des opérateurs

$\boxed{\Gamma \vdash e : t}$	
$\frac{\Gamma \vdash lv : t}{\Gamma \vdash \&lv : t*} \text{ (ADDR)}$	$\frac{\forall i \in [1; n], \Gamma \vdash e_i : t_i}{\Gamma \vdash \{l_1 : e_1; \dots; l_n : e_n\} : \{l_1 : t_1; \dots; l_n : t_n\}} \text{ (STRUCT)}$
$\frac{\Gamma \vdash e : (t_1, \dots, t_n) \rightarrow t \quad \forall i \in [1; n], \Gamma \vdash e_i : t_i}{\Gamma \vdash e(e_1, \dots, e_n) : t} \text{ (CALL)}$	$\frac{\Gamma \vdash lv : t \quad \Gamma \vdash e : t}{\Gamma \vdash lv \leftarrow e : t} \text{ (SET)}$
$\frac{\forall i \in [1; n], \Gamma \vdash e_i : t}{\Gamma \vdash [e_1; \dots; e_n] : t[]} \text{ (ARRAY)}$	$\frac{\Gamma' = (\Gamma - \underline{R}), a_1 : t_1, \dots, a_n : t_n, \underline{R} : t \quad \Gamma' \vdash i}{\Gamma \vdash \text{fun}(a_1, \dots, a_n)\{i\} : (t_1, \dots, t_n) \rightarrow t} \text{ (FUN)}$
$\boxed{\Gamma \vdash i}$	
$\frac{}{\Gamma \vdash \text{PASS}} \text{ (PASS)}$	$\frac{\Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash i_1; i_2} \text{ (SEQ)}$
$\frac{\Gamma \vdash e : t \quad \Gamma, x : t \vdash i}{\Gamma \vdash \text{DECL } x = e \text{ IN } \{i\}} \text{ (DECL)}$	$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i_1 \quad \Gamma \vdash i_2}{\Gamma \vdash \text{IF}(e)\{i_1\}\text{ELSE}\{i_2\}} \text{ (IF)}$
$\frac{\Gamma \vdash e : \text{INT} \quad \Gamma \vdash i}{\Gamma \vdash \text{WHILE}(e)\{i\}} \text{ (WHILE)}$	$\frac{\Gamma \vdash \underline{R} \leftarrow e}{\Gamma \vdash \text{RETURN}(e)} \text{ (RETURN)}$
$\boxed{\Gamma \vdash p \rightarrow \Gamma'}$	
$\frac{\Gamma \vdash e : t}{\Gamma \vdash e \rightarrow \Gamma} \text{ (T-EXP)}$	$\frac{\Gamma \vdash e : t \quad \Gamma' = (x, t), \Gamma}{\Gamma \vdash x = e \rightarrow \Gamma'} \text{ (T-VAR)}$

FIGURE C.3: Règles de typage

$\boxed{m \models v : \tau}$	
$\frac{}{m \models n : \text{INT}} \text{ (S-INT)}$	$\frac{}{m \models d : \text{FLOAT}} \text{ (S-FLOAT)}$
$\frac{}{m \models () : \text{UNIT}} \text{ (S-UNIT)}$	
$\frac{}{m \models \text{NULL} : \tau*} \text{ (S-NULL)}$	$\frac{m \models_{\Phi} \varphi : \tau}{m \models \widehat{\&} \varphi : \tau*} \text{ (S-REF)}$
	$\frac{\forall i \in [1; n]. m \models v_i : \tau}{m \models [\widehat{v_1; \dots; v_n}] : \tau[]} \text{ (S-ARRAY)}$
$\frac{\forall i \in [1; n]. m \models v_i : \tau_i}{m \models \{l_1 : \widehat{v_1; \dots; v_n}; l_n : v_n\} : \{l_1 : \tau_1; \dots; l_n : \tau_n\}} \text{ (S-STRUCT)}$	
$\frac{}{m \models \text{fun}(x_1, \dots, x_n)\{i\} : \text{FUN}_n} \text{ (S-FUN)}$	

FIGURE C.4: Règles de typage sémantique

$$\begin{array}{c}
\frac{\Gamma \vdash lv : t}{\Gamma \vdash \diamond lv : t@} \text{ (TAINT)} \qquad \frac{\Gamma \vdash e_d : t * \quad \Gamma \vdash e_s : t@}{\Gamma \vdash \text{copy_from_user}(e_d, e_s) : \text{INT}} \text{ (GETU)} \\
\\
\frac{\Gamma \vdash e_d : t@ \quad \Gamma \vdash e_s : t *}{\Gamma \vdash \text{copy_to_user}(e_d, e_s) : \text{INT}} \text{ (PUTU)} \qquad \frac{\Gamma \vdash lv : t}{\Gamma \vdash \&lv : t * } \text{ (ADDR0F-KERNEL)} \\
\\
\frac{\boxplus \in \{+p, -p\} \quad \Gamma \vdash e_1 : t@ \quad \Gamma \vdash e_2 : \text{INT}}{\Gamma \vdash e_1 \boxplus e_2 : t@} \text{ (PTR-ARITH-USER)}
\end{array}$$

FIGURE C.5: Règles de typage des extensions noyau

D.1 Composition de lentilles

Démonstration. Pour prouver que $\mathcal{L}_1 \ggg \mathcal{L}_2 \in \text{LENS}_{A,C}$, il suffit de prouver les trois propriétés caractéristiques.

GetPut

$$\begin{aligned}
 \text{put}_{\mathcal{L}}(\text{get}_{\mathcal{L}}(r), r) &= \text{put}_{\mathcal{L}}(\text{get}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_1}(r)), r) && \{ \text{définition de } \text{get}_{\mathcal{L}} \} \\
 &= \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_1}(r)), \text{get}_{\mathcal{L}_1}(r)), r) && \{ \text{définition de } \text{put}_{\mathcal{L}} \} \\
 &= \text{put}_{\mathcal{L}_1}(\text{get}_{\mathcal{L}_1}(r), r) && \{ \text{GETPUT sur } \mathcal{L}_2 \} \\
 &= r && \{ \text{GETPUT sur } \mathcal{L}_1 \}
 \end{aligned}$$

PutGet

$$\begin{aligned}
 \text{get}_{\mathcal{L}}(\text{put}_{\mathcal{L}}(a, r)) &= \text{get}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}}(a, r))) && \{ \text{définition de } \text{get}_{\mathcal{L}} \} \\
 &= \text{get}_{\mathcal{L}_2}(\text{get}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r))) && \{ \text{définition de } \text{put}_{\mathcal{L}} \} \\
 &= \text{get}_{\mathcal{L}_2}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r))) && \{ \text{PUTGET sur } \mathcal{L}_1 \} \\
 &= a && \{ \text{PUTGET sur } \mathcal{L}_2 \}
 \end{aligned}$$

PutPut

$$\begin{aligned}
& \text{put}_{\mathcal{L}}(a', \text{put}_{\mathcal{L}}(a, r)) \\
&= \text{put}_{\mathcal{L}}(a', \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r)) \\
&\quad \{ \text{définition de } \text{put}_{\mathcal{L}} \} \\
&= \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a', \text{get}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r))), \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r)) \\
&\quad \{ \text{définition de } \text{put}_{\mathcal{L}} \} \\
&= \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a', \text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r))), \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r)) \\
&\quad \{ \text{GETPUT sur } \mathcal{L}_1 \} \\
&= \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a', \text{get}_{\mathcal{L}_1}(r)), \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a, \text{get}_{\mathcal{L}_1}(r)), r)) \\
&\quad \{ \text{PUTPUT sur } \mathcal{L}_2 \} \\
&= \text{put}_{\mathcal{L}_1}(\text{put}_{\mathcal{L}_2}(a', \text{get}_{\mathcal{L}_1}(r)), r) \\
&\quad \{ \text{PUTPUT sur } \mathcal{L}_1 \} \\
&= \text{put}_{\mathcal{L}}(a', r) \\
&\quad \{ \text{définition de } \ggg \}
\end{aligned}$$

□

D.2 Progrès

(théorème 5.1)

Démonstration. On procède par induction sur la dérivation du jugement de typage. Puisque les jugements $\Gamma \vdash i$, $\Gamma \vdash e : t$ et $\Gamma \vdash lv : t$ sont interdépendants, on traite tous les cas par récursion mutuelle.

Le squelette de cette preuve est une analyse de cas selon la dernière règle utilisée. La plupart des cas ont la même forme : , et on utilise l'hypothèse de récurrence sur les sous-éléments syntaxiques (en appliquant éventuellement le lemme 5.3 d'inversion pour établir qu'ils sont bien typés). Dans le cas “valeur”, on appelle une règle qui permet de transformer une opération syntaxique en opération sémantique (par exemple, on transforme le + unaire en un $\hat{+}$ sémantique). Dans le cas “évaluation”, on applique la règle CTX avec un contexte particulier qui permet de passer d'un jugement $\langle a, m \rangle \rightarrow \langle a', m' \rangle$ à un jugement $\langle b, m \rangle \rightarrow \langle b', m' \rangle$ (où a apparaît dans b). Enfin, dans le cas “erreur”, on utilise EVAL-ERR avec ce même contexte C .

Ceci est valable pour la majorité des cas. Il faut faire attention en particulier aux opérations sémantiques qui peuvent produire des erreurs (comme la division, ou l'opérateur Lookup(\cdot, \cdot)).

PASS : Ce cas est trivial.

RETURN : Partant de $i = \text{RETURN}(e)$, on applique le lemme d'inversion. Il vient : $\exists t_e, \Gamma \vdash e : t_e$. On applique l'hypothèse de récurrence à e .

- $e = v$. Alors $i = \text{RETURN}(v)$, ce qui nous permet de conclure.
- $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. Alors en appliquant CTX avec $C = \text{RETURN}(\bullet)$, on conclut que $\langle \text{RETURN}(e), m \rangle \rightarrow \langle \text{RETURN}(e'), m' \rangle$.

- $\langle e, m \rangle \rightarrow \Omega$. On applique EVAL-ERR avec ce même C .

SEQ : Avec $i = i_1; i_2$, on applique l'hypothèse de récurrence à i_1 .

- $i_1 = \text{PASS}$. On peut donc appliquer la règle SEQ et donc $\langle i, m \rangle \rightarrow \langle i_2, m \rangle$.
- $i_1 = \text{RETURN}(v)$. Alors on peut appliquer la règle RETURN : $\langle i, m \rangle \rightarrow \langle \text{RETURN}(v), m \rangle$.
- $\langle i_1, m \rangle \rightarrow \langle i'_1, m' \rangle$. Soit $C = \bullet; i_2$. Par CTX il vient $\langle i, m \rangle \rightarrow \langle i'_1; i_2, m' \rangle$.
- $\langle i_1, m \rangle \rightarrow \Omega$. Avec ce même C dans EVAL-ERR on trouve $\langle i, m \rangle \rightarrow \Omega$.

EXP : Ici $i = e$. On peut appliquer l'hypothèse de récurrence à e qui est "plus petit" que i ($i ::= e$ introduit un constructeur implicite).

- $e = v$. Alors on peut appliquer EXP : $\langle e, m \rangle \rightarrow \langle \text{PASS}, m \rangle$.
- $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. Alors $\langle i, m \rangle \rightarrow \langle e', m' \rangle$ (cela revient à appliquer CTX au constructeur implicite mentionné ci-dessus).
- $\langle e, m \rangle \rightarrow \Omega$. C'est-à-dire $\langle i, m \rangle \rightarrow \Omega$.

DECL : Ici $i = \text{DECL } x = e \text{ IN } \{i'\}$. On commence par appliquer l'hypothèse de récurrence à e .

- $e = v$. On applique alors l'hypothèse de récurrence à i sous $\Gamma' = \Gamma, x : t$ et avec $m' = \text{Extend}(m, x \mapsto v)$.
 - $i = \text{PASS}$. Dans ce cas la règle DECL-PASS s'applique.
 - $i = \text{RETURN}(v)$. Idem avec DECL-RETURN.
 - $\langle i, m' \rangle \rightarrow \langle i', m'' \rangle$. Par DECL-CTX il vient $\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \langle \text{DECL } x = v' \text{ IN } \{i'\}, m''' \rangle$ où $v' = m''[(0, x)]_A$ et $m''' = m'' - x$.
 - $\langle i', m' \rangle \rightarrow \Omega$. On applique DECL-ERR.
- $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. On pose $C = \text{DECL } x = \bullet \text{ IN } \{ \}$ et on conclut avec la règle CTX.
- $\langle e, m \rangle \rightarrow \Omega$. Idem avec EVAL-ERR.

IF : Ici $i = \text{IF}(e)\{i_1\}\text{ELSE}\{i_2\}$. On applique l'hypothèse de récurrence à e .

- $e = v$.
Si $v \neq 0$, on applique IF-TRUE. Dans le cas contraire, on applique IF-FALSE.
- $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. On pose $C = \text{IF}(\bullet)\{i_1\}\text{ELSE}\{i_2\}$ et on conclut avec CTX.
- $\langle e, m \rangle \rightarrow \Omega$. Avec ce même C et EVAL-ERR.

WHILE : Ce cas est direct : la règle WHILE peut toujours s'appliquer.

XXXXXXXXXXXXXXXXXXXXXXXXXXXX

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

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

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

CST-UNIT : e est alors égale à $()$, qui est une valeur.

FUN : Ce cas est direct : la règle EXP-FUN s'applique.

OP-INT : Cela implique que $e = e_1 \boxplus e_2$. Par le lemme 5.3, 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, sauf dans le cas d'une division par zéro ($\boxplus \in \{/, \%, /. \}$ et $v_2 = 0$) où alors $v_1 \hat{\boxplus} v_2 = \Omega_{div}$. Dans ce cas, on a alors par EXP-ERR $\langle e, m \rangle \rightarrow \Omega_{div}$.
 - $\exists (e'_2, m''), \langle e_2, m' \rangle \rightarrow \langle e'_2, m'' \rangle$.
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 $\langle e, m \rangle \rightarrow \langle v_1 \boxplus e'_2, m'' \rangle$.
 - $\langle e_2, m' \rangle \rightarrow \Omega$. De EVAL-ERR avec $C = v_1 \boxplus \bullet$ vient alors $\langle e, m \rangle \rightarrow \Omega$.
- $\exists (e'_1, m'), \langle e_1, m \rangle \rightarrow \langle e'_1, m' \rangle$. 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 $\langle e, m \rangle \rightarrow \langle e'_1 \boxplus e_2, m' \rangle$.
- $\langle e_1, m \rangle \rightarrow \Omega$. D'après EVAL-ERR avec $C = \bullet \boxplus e_2$, on a $\langle e, m \rangle \rightarrow \Omega$.

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

OP-EQ : 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 $\langle e_1, m \rangle \rightarrow \Omega$. De EVAL-ERR avec $C = + \bullet$ il vient $\langle e, m \rangle \rightarrow \Omega$.

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

ADDR : On applique l'hypothèse de récurrence à lv .

Les cas d'évaluation et d'erreur sont traités en appliquant respectivement CTX et EVAL-ERR avec $C = \&\bullet$. Dans le cas où $lv = \varphi$, on peut appliquer EXP-ADDR0F.

- $lv = \varphi$. On applique l'hypothèse de récurrence à e .
 - $e = v$. Alors on peut appliquer EXP-SET.
 - $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. On conclut avec $C = \varphi \leftarrow \bullet$.
 - $\langle e, m \rangle \rightarrow \Omega$. Idem.
- $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$. On conclut avec $C = \bullet \leftarrow e$.
- $\langle lv, m \rangle \rightarrow \Omega$. Idem.

- $i = \text{PASS}$. Alors on peut appliquer EXP-CALL-RETURN-UNIT.
- $i = \text{RETURN}(v)$. Idem avec EXP-CALL-RETURN.
- $\langle i, m_1 \rangle \rightarrow \langle i', m_2 \rangle$.
Alors on peut appliquer EXP-CALL-CTX.
- $\langle i, m \rangle \rightarrow \Omega$. On peut alors appliquer EXP-CALL-ERR.

Puisque $\Gamma \vdash x : t$, on en déduit que $x : t \in \Gamma$ (par le lemme d'inversion), donc que $x \in \text{dom}(\Gamma)$. Comme $\Gamma \models m$, $x \in \{\text{Varname}(a) / a \in \text{Visible}(m)\}$, donc $\exists a \in \text{Visible}(m). x = \text{Varname}(a)$ et $\text{Lookup}(x, m)$ ne renvoie donc pas d'erreur.

LV-DEREF : Appliquons l'hypothèse de récurrence à lv (vue en tant qu'expression).

- $lv = v$. Puisque $\Gamma \vdash v : t*$, on déduit du lemme 5.4 que $v = \varphi$ ou $v = \text{NULL}$. Dans le premier cas, la règle PHI-DEREF s'applique. Dans le second, puisque $\langle * \text{NULL}, m \rangle \rightarrow \Omega_{ptr}$, on a $\langle e, m \rangle \rightarrow \Omega_{ptr}$.
- $\langle lv, m \rangle \rightarrow \langle e', m' \rangle$. De CTX avec $C = *\bullet$, on obtient $\langle e, m \rangle \rightarrow \langle *e', m' \rangle$.
- $\langle lv, m \rangle \rightarrow \Omega$. En appliquant EVAL-ERR avec $C = *\bullet$, on obtient $\langle e, m \rangle \rightarrow \Omega$.

LV-INDEX : De même, on applique l'hypothèse de récurrence à lv .

- $lv = v$. Comme $\Gamma \vdash v : t[]$, on déduit du lemme 5.4 que $v = [\nu_1; \dots; \nu_p]$. Appliquons l'hypothèse de récurrence à e .
 - $e = v'$. Puisque $\Gamma \vdash e : \text{INT}$, on réapplique le lemme 5.4 et $v' = n$. D'après PHI-ARRAY, $\langle lv[e], m \rangle \rightarrow \langle [\nu_1; \dots; \nu_p][\widehat{n}], m \rangle$. Deux cas sont à distinguer : si $n \in [0; p-1]$, la partie droite vaut ν_{n+1} et donc $\langle lv[e], m \rangle \rightarrow \langle \nu_{n+1}, m \rangle$. Sinon elle vaut Ω_{array} et $\langle lv[e], m \rangle \rightarrow \Omega_{array}$ par EXP-ERR.
 - $\langle e, m \rangle \rightarrow \langle e', m' \rangle$. En appliquant CTX avec $C = v[\bullet]$, on en déduit
 - $\langle lv[e], m \rangle \rightarrow \langle lv[e'], m' \rangle$.
 - $\langle e, m \rangle \rightarrow \Omega$. Avec EVAL-ERR sous ce même contexte, $\langle lv[e], m \rangle \rightarrow \Omega$
- $\langle lv, m \rangle \rightarrow \langle e', m' \rangle$. On applique alors CTX avec $C = \bullet[e]$, et $\langle lv[e], m \rangle \rightarrow \langle e'[e], m' \rangle$.
- $\langle lv, m \rangle \rightarrow \Omega$. Toujours avec $C = \bullet[e]$, de EVAL-ERR il vient $\langle lv[e], m \rangle \rightarrow \Omega$.

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

- $lv = \varphi$ Alors PHI-STRUCT s'applique. Puisque $(l, t) \in S$, l'accès au champ l ne provoque pas d'erreur Ω_{field} . Donc $\langle e, m \rangle \rightarrow \langle \varphi[l], m \rangle$.
- $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$ En appliquant CTX avec $C = \bullet.l_S$, il vient $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$.
- $\langle lv, m \rangle \rightarrow \Omega$ En appliquant EVAL-ERR avec $C = \bullet.l_S$, on a $\langle lv, m \rangle \rightarrow \Omega$.

PTR-ARITH : Le schéma est similaire au cas OP-INT. Le seul cas intéressant arrive lorsque e_1 et e_2 sont des valeurs. D'après le lemme 5.4 :

- $e_1 = \text{NULL}$ ou $e_1 = \varphi$
- $e_2 = n$

D'après EXP-BINOP, $\langle e, m \rangle \rightarrow \langle e_1 \hat{\boxplus} n, m \rangle$.

On se réfère ensuite à la définition de $\hat{\boxplus}$ (page 54) : si e_1 est de la forme $\varphi[m]$, alors $e_1 \hat{\boxplus} n = \varphi[m + n]$. Donc $\langle e, m \rangle \rightarrow \langle \varphi[m + n], m \rangle$.

Dans les autres cas ($e_1 = \text{NULL}$ ou $e_1 = \varphi$ avec φ pas de la forme $\varphi'[m]$), on a $e_1 \hat{\boxplus} n = \Omega_{ptr}$. Donc d'après EXP-ERR, $\langle e, m \rangle \rightarrow \Omega_{ptr}$.

□

D.3 Préservation

(théorème 5.2)

Démonstration. On procède par induction sur la dérivation de $\langle \cdot, m \rangle \rightarrow \langle \cdot, m' \rangle$.

N.B. CTX apparaît 3 fois.

Cas $\Gamma \vdash e : t$ et $\langle e, m \rangle \rightarrow \langle e', m' \rangle$.

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

- $\Gamma \vdash e : t'$
- $\Gamma \vdash e' : t'$
- $\Gamma \vdash C(e) : t.$

$$\frac{m' = \text{Cleanup}(m)}{\langle \text{fun}(a_1, \dots, a_n) \{ \text{RETURN}(v) \} (v_1, \dots, v_n), m \rangle \rightarrow \langle v, m' \rangle} \text{ (EXP-CALL-RETURN)}$$

$$\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{PASS}\}, m \rangle \rightarrow \langle \text{PASS}, m' \rangle} \text{ (DECL-PASS)}$$

DECL-RETURN :

$$\frac{m' = \text{Cleanup}(m - x)}{\langle \text{DECL } x = v \text{ IN } \{\text{RETURN}(v')\}, m \rangle \rightarrow \langle \text{RETURN}(v'), m' \rangle} \text{ (DECL-RETURN)}$$

DECL-CTX :

$$\frac{\begin{array}{c} m' = \text{Extend}(m, x \mapsto v) \\ \langle i, m' \rangle \rightarrow \langle i', m'' \rangle \quad v' = m''[(0, x)]_A \quad m''' = m'' - x \end{array}}{\langle \text{DECL } x = v \text{ IN } \{i\}, m \rangle \rightarrow \langle \text{DECL } x = v' \text{ IN } \{i'\}, m''' \rangle} \text{ (DECL-CTX)}$$

IF-FALSE : D'après le lemme d'inversion, $\Gamma \vdash i_f$.**IF-TRUE :** D'après le lemme d'inversion, $\Gamma \vdash i_t$.**WHILE :** D'après le lemme d'inversion, $\Gamma \vdash e : t$ et $\Gamma \vdash i$. Par SEQ, on a $\Gamma \vdash i; \text{WHILE}(e)\{i\}$. Enfin par IF il vient $\Gamma \vdash \text{IF}(e)\{i; \text{WHILE}(e)\{i\}\}$.**RETURN :** Par le lemme d'inversion, $\Gamma \vdash \text{RETURN}(v)$.**EXP-CALL-CTX :**

$$\frac{\begin{array}{c} m_1 = \text{Push}(m_0, (a_1 \mapsto v_1, \dots, a_n \mapsto v_n)) \\ \langle i, m_1 \rangle \rightarrow \langle i', m_2 \rangle \quad \forall i \in [1; n], v'_i = m_2[(0, a_i)]_A \quad m_3 = \text{Pop}(m_2) \end{array}}{\langle \text{fun}(a_1, \dots, a_n)\{i\}(v_1, \dots, v_n), m_0 \rangle \rightarrow \langle \text{fun}(a_1, \dots, a_n)\{i'\}(v'_1, \dots, v'_n), m_3 \rangle} \text{ (EXP-CALL-CTX)}$$

□

D.4 Progrès pour les types qualifiés

(théorème 6.1)

Démonstration. On procède de la même manière que pour le théorème 5.1 (prouvé en annexe D.2). En fait, puisque le schéma de preuve porte sur les règles de typage, il suffit de traiter les cas supplémentaires.

TAINT : Alors $e = \diamond lv$. On applique l'hypothèse de récurrence à lv .

- $lv = \varphi$. Alors on peut appliquer **EXPR-TAINTED**.
- $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$. On conclut en utilisant **CTX** avec $C = \diamond \bullet$.
- $\langle lv, m \rangle \rightarrow \Omega$. On applique **EVAL-ERR** avec ce même C .

ADDR-OF-KERNEL : On applique l'hypothèse de récurrence à lv .

- $lv = \varphi$. Alors **EXP-ADDR-OF** s'applique.
- $\langle lv, m \rangle \rightarrow \langle lv', m' \rangle$. On applique **CTX** à $C = \& \bullet$.
- $\langle lv, m \rangle \rightarrow \Omega$. Idem avec la règle **EVAL-ERR**.

PTR-ARITH-USER : Ce cas est similaire à **OP-INT**.

GETU : On applique l'hypothèse de récurrence à e_d .

- $e_d = v_d$. On applique l'hypothèse de récurrence à e_s .
 - $e_s = v_s$.
D'après le lemme 5.4, v_s a pour forme φ_s .
On distingue la forme de φ_s :
 - $\varphi_s = \hat{\diamond} \varphi$. Alors on applique USER-GET-OK.
 - $\nexists \varphi, \varphi_s = \hat{\diamond} \varphi$. Alors idem avec USER-GET-ERR.
 - $\langle e_s, m \rangle \rightarrow \langle e'_s, m' \rangle$. Posons $C = \text{copy_from_user}(v_d, \bullet)$. On conclut avec CTX.
 - $\langle e_s, m \rangle \rightarrow \Omega$. Idem avec EVAL-ERR.
- $\langle e_d, m \rangle \rightarrow \langle e'_d, m' \rangle$. On applique CTX avec $C = \text{copy_from_user}(\bullet, e_s)$.
- $\langle e_d, m \rangle \rightarrow \Omega$. On utilise EVAL-ERR avec ce même contexte.

PUTU : Ce cas est similaire au cas GETU, en appliquant les règles USER-PUT-OK et USER-PUT-ERR.

□

D.5 Préservation pour les types qualifiés

(théorème 6.2)

ANALYSE DE BITMASKS

Par exemple, en ouvrant un fichier on précise le mode du fichier (lecture, écriture ou les deux) par les bits 1 et 2, s'il faut créer le fichier ou non s'il n'existe pas par le bit 7, s'il dans ce cas il doit être effacé par le bit 8), etc. On obtient un mode en réalisant un "ou" bit à bit entre des constantes. Ces deux utilisations du type `int` n'ont rien à voir ; il faudrait donc empêcher d'utiliser un descripteur de fichier comme un mode, et vice-versa. De même, aucun opérateur n'a de sens sur les descripteurs de fichier, mais l'opérateur `|` du "ou" bit à bit doit rester possible pour les modes.

Les fonctions suivantes permettent d'ouvrir un fichier sous UNIX.

```
int open(const char *pathname, int flags);
int open(const char *pathname, int flags, mode_t mode);
int creat(const char *pathname, mode_t mode);
```

Il est intéressant de remarquer que ces fonctions manipulent des entiers :

- la valeur de retour est un *descripteur de fichier*, qui correspond à un indice numérique dans une table interne au processus. Par exemple, 0 désigne son entrée standard, 1 sa sortie standard, et 2 son flux d'erreur standard.
- `flags` est un entier qui encode les options liées à l'ouverture du fichier.
- `mode` (en fait un *unsigned int*) encode de la même manière les permissions que doit avoir le fichier créé, le cas échéant.

Ces entiers ont une particularité en commun : ils ne sont pas utilisés comme entiers au sens arithmétique du terme. C'est-à-dire que `int` ou `unsigned int` est le bon type car il correspond à la bonne taille de données, mais les opérations comme `+` ou `/` n'ont pas de sens sur ces valeurs.

On décrit ici une technique de typage pour détecter et interdire ces mauvaises utilisations.

E.1 Ajouts

On étend les types avec un type `BITS`.

Type $t ::= \dots$

| BITS

Entier utilisé comme *bitmask*

Au niveau des valeurs, les entiers sont représentés par des valeurs entières classiques \hat{n} . C'est pourquoi on ajoute la règle suivante de compatibilité entre le type de valeurs INT et le type statique BITS.

$$\frac{}{\text{Comp}(\text{INT}, \text{BITS})} \text{ (COMP-BITS)}$$

LISTE DES FIGURES

1.1	Règles d'exécution de code et d'accès à la mémoire.	3
2.1	Mécanisme de mémoire virtuelle.	15
2.2	Implantation de la mémoire virtuelle	15
2.3	Appel de <code>gettimeofday</code>	16
2.4	Implantation de l'appel système <code>gettimeofday</code>	17
3.1	Domaine des signes	21
3.2	Quelques domaines abstraits	21
3.3	Treillis de qualificateurs	23
4.1	Fonctionnement d'une lentille	38
4.2	Fonctionnement d'une lentille indexée	39
4.3	Composition de lentilles	40
4.4	Syntaxe des expressions	42
4.5	Syntaxe des instructions	42
4.6	Syntaxe des opérateurs	43
4.7	Valeurs	44
4.8	Composantes d'un état mémoire	45
4.9	Opérations de pile	47
4.10	Nettoyage d'un cadre de pile	48
4.11	Contextes d'évaluation	52
4.12	Évaluation des <i>left-values</i>	52
4.13	Évaluation stricte ou paresseuse des <i>left-values</i>	53
4.14	Appel d'une fonction	55
4.15	Cas d'erreur lors d'une évaluation	58
4.16	Évaluation des phrases d'un programme	58
5.1	Programmes bien et mal formés	61
5.2	Types et environnements de typage	62
5.3	Typage d'une suite de phrases et d'un programme	63
5.4	Jugements d'égalité sur les types	65
5.5	Typage des phrases	67
5.6	Types sémantiques	68
5.7	Règles de typage sémantique	69
5.8	Compatibilité entre types sémantiques et statiques	69
6.1	Implantation d'un appel système qui remplit une structure par pointeur	78
6.2	Ajouts liés aux pointeurs utilisateurs (par rapport à l'évaluateur du chapitre 4)	79
6.3	Ajouts liés aux qualificateurs de types (par rapport aux figures 5.2 et 5.6)	81
6.4	Notations utilisées pour les pointeurs	82
7.1	Compilation depuis Newspeak	94
7.2	Exemple de programme C à compiler	95

7.3	Partage lors de l'unification	96
7.4	Unification par mutation de références	97
7.5	Cycle dans le graphe de types	97
7.6	Fonction principale de ptr type	98
7.7	Inférence des déclarations de variable et appels de fonction	100
7.8	Représentation des types	101
7.9	Fonction de raccourcissement des représentations de types	102
7.10	Fonction d'unification	103
7.11	Structures	104
7.12	Qualificateurs	105
7.13	Unification directe ou retardée	105
8.1	Espace d'adressage d'un processus	108
8.2	Patch résolvant le problème de pointeur utilisateur.	111
B.1	Règles d'évaluation des contextes	128
B.2	Règles d'évaluation des erreurs	128
B.3	Règles d'évaluation des left-values et expressions	129
B.4	Règles d'évaluation des instructions et phrases	130
B.5	Règles d'évaluation des extensions noyau	131
C.1	Règles de typage des constantes et variables	133
C.2	Règles de typage des opérateurs	134
C.3	Règles de typage	135
C.4	Règles de typage sémantique	135
C.5	Règles de typage des extensions noyau	136

LISTE DES DÉFINITIONS

4.1	Définition (Lentille)	37
4.2	Définition (Lentille indexée)	38
4.3	Définition (Composition de lentilles)	39
4.4	Définition (Recherche de variable)	46
4.5	Définition (Manipulations de pile)	46
4.6	Définition	55
5.1	Définition (État mémoire bien typé)	69
5.2	Définition (Hauteur d'une valeur)	70
5.3	Définition (État mémoire bien typé)	70

LISTE DES THÉORÈMES ET LEMMES

5.1	Lemme (Hauteur des chemins typés)	70
5.2	Lemme (Accès à une mémoire bien typée)	71
5.3	Lemme (Inversion)	71
5.4	Lemme (Formes canoniques)	73
5.1	Théorème (Progrès)	73
5.5	Lemme (Représentabilité)	74
5.2	Théorème (Préservation)	75
5.3	Théorème (Progrès pour les phrases)	75
5.4	Théorème (Préservation pour les phrases)	75
6.1	Théorème (Progrès pour les types qualifiés)	83
6.2	Théorème (Préservation pour les types qualifiés)	83

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/>
- [🌐⁹] The Rust Programming Language
<http://www.rust-lang.org/>

BIBLIOGRAPHIE

- [AB07] Andrew W. Appel and Sandrine Blazy. Separation logic for small-step Cminor (extended version). Research report 6138, INRIA, 2007. 29 pages. 7
- [ABD⁺07] Alex Aiken, Suhabe Bugrara, Isil Dillig, Thomas Dillig, Brian Hackett, and Peter Hawkins. An overview of the saturn project. In *Proceedings of the 7th ACM SIGPLAN-SIGSOFT workshop on Program analysis for software tools and engineering*, PASTE '07, pages 43–48, New York, NY, USA, 2007. ACM. 24
- [AH07] Xavier Allamigeon and Charles Hymans. Analyse statique par interprétation abstraite. In Eric Filiol, editor, *5ème Symposium sur la Sécurité des Technologies de l'Information et des Communications (SSTIC'07)*, Rennes, France, June 2007. 20
- [BA08] S. Bugrara and A. Aiken. Verifying the Safety of User Pointer Dereferences. In *Security and Privacy, 2008. SP 2008. IEEE Symposium on*, pages 325–338, 2008. 24
- [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. 22
- [BC05] Daniel P. Bovet and Marco Cesati. *Understanding the Linux Kernel, Third Edition*. O'Reilly Media, third edition edition, November 2005. 108
- [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. 20
- [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. 7
- [BDN09] Ana Bove, Peter Dybjer, and Ulf Norell. A brief overview of agda — a functional language with dependent types. In *Proceedings of the 22nd International Conference on Theorem Proving in Higher Order Logics, TPHOLs '09*, pages 73–78, Berlin, Heidelberg, 2009. Springer-Verlag. 25
- [BLS05] Mike Barnett, K. Rustan M. Leino, and Wolfram Schulte. The spec# programming system : an overview. In *Proceedings of the 2004 international conference on Construction and Analysis of Safe, Secure, and Interoperable Smart Devices, CASSIS'04*, pages 49–69, Berlin, Heidelberg, 2005. Springer-Verlag. 25
- [CC77] Patrick Cousot and Radhia Cousot. Abstract interpretation : a unified lattice model for static analysis of programs by construction or approximation of fix-points. 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. 20, 21

- [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/>). 20
- [CCF⁺05] Patrick Cousot, Radhia Cousot, Jérôme Feret, Laurent Mauborgne, Antoine Miné, David Monniaux, and Xavier Rival. The ASTREÉ analyzer. In Shmuel Sagiv, editor, *ESOP*, volume 3444 of *Lecture Notes in Computer Science*, pages 21–30. Springer, 2005. 22
- [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. 22
- [DDMPn10] Javier De Dios, Manuel Montenegro, and Ricardo Peña. Certified absence of dangling pointers in a language with explicit deallocation. In *Proceedings of the 8th international conference on Integrated formal methods*, IFM’10, pages 305–319, Berlin, Heidelberg, 2010. Springer-Verlag. 117
- [DM82] Luis Damas and Robin Milner. Principal type-schemes for functional programs. In *Proceedings of the 9th ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL ’82, pages 207–212, New York, NY, USA, 1982. ACM. 94
- [DRS00] Nurit Dor, Michael Rodeh, and Mooly Sagiv. CSSV : Towards a realistic tool for statically detecting all buffer overflows in C, 2000. 25
- [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. 94
- [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. 23, 87
- [FGM⁺07] J. Nathan Foster, Michael B. Greenwald, Jonathan T. Moore, Benjamin C. Pierce, and Alan Schmitt. Combinators for bidirectional tree transformations : A linguistic approach to the view-update problem. *ACM Trans. Program. Lang. Syst.*, 29(3), May 2007. 37
- [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. 22
- [Flo67] Robert W. Floyd. Assigning Meanings to Programs. In J. T. Schwartz, editor, *Proceedings of a Symposium on Applied Mathematics*, volume 19 of *Mathematical Aspects of Computer Science*, pages 19–31, Providence, 1967. American Mathematical Society. 24
- [FTA02] Jeffrey S. Foster, Tachio Terauchi, and Alex Aiken. Flow-sensitive type qualifiers. In *Proceedings of the ACM SIGPLAN 2002 Conference on Programming language design and implementation*, PLDI ’02, pages 1–12, New York, NY, USA, 2002. ACM. 23

- [GGTZ07] Stephane Gaubert, Eric Goubault, Ankur Taly, and Sarah Zennou. Static analysis by policy iteration on relational domains. In Rocco Nicola, editor, *Programming Languages and Systems*, volume 4421 of *Lecture Notes in Computer Science*, pages 237–252. Springer Berlin Heidelberg, 2007. 22
- [GMJ⁺02] Dan Grossman, Greg Morrisett, Trevor Jim, Michael Hicks, Yanling Wang, and James Cheney. Region-based memory management in Cyclone. *SIGPLAN Not.*, 37(5) :282–293, May 2002. 24
- [Gon05] Georges Gonthier. A computer-checked proof of the Four Colour Theorem. 2005. 26
- [Gor04] Mel Gorman. *Understanding the Linux Virtual Memory Manager*. Prentice Hall PTR, Upper Saddle River, NJ, USA, 2004. 108, 117
- [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. 22
- [Har88] Norm Hardy. The confused deputy (or why capabilities might have been invented). *ACM Operating Systems Review*, 22(4) :36–38, October 1988. 4, 17
- [HJLT05] Thomas Hallgren, Mark P. Jones, Rebekah Leslie, and Andrew Tolmach. A principled approach to operating system construction in haskell. In *Proceedings of the tenth ACM SIGPLAN international conference on Functional programming*, ICFP '05, pages 116–128, New York, NY, USA, 2005. ACM. 6
- [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. 7, 93, 94
- [HLA⁺05] G. Hunt, J. Larus, M. Abadi, M. Aiken, P. Barham, M. Fahndrich, C. Hawblitzel, O. Hodson, S. Levi, N. Murphy, B. Steensgaard, D. Tarditi, T. Wobber, and B. Zill. An Overview of the Singularity Project. Technical report, Microsoft Research, 2005. 6
- [Hoa69] C. A. R. Hoare. An axiomatic basis for computer programming. *Commun. ACM*, 12(10) :576–580, October 1969. 24
- [Int10] Intel, Santa Clara, CA, USA. *Intel® 64 and IA-32 Architectures Software Developer's Manual*, June 2010. 109
- [ISO99] ISO. The ANSI C standard (C99). Technical Report WG14 N1124, ISO/IEC, 1999. 6, 54
- [JMG⁺02] Trevor Jim, J. Greg Morrisett, Dan Grossman, Michael W. Hicks, James Cheney, and Yanling Wang. Cyclone : A safe dialect of c. In *Proceedings of the General Track of the annual conference on USENIX Annual Technical Conference, ATEC '02*, pages 275–288, Berkeley, CA, USA, 2002. USENIX Association. 24
- [JW04] Robert Johnson and David Wagner. Finding user/kernel pointer bugs with type inference. In *USENIX Security Symposium*, pages 119–134, 2004. 23, 87
- [KcS07] Oleg Kiselyov and Chung chieh Shan. Lightweight static capabilities. *Electr. Notes Theor. Comput. Sci.*, 174(7) :79–104, 2007. 22

- [Ker81] Brian W. Kernighan. Why Pascal is not my favorite programming language. Technical report, AT&T Bell Laboratories, April 1981. 5
- [KR88] Brian W. Kernighan and Dennis M. Ritchie. *The C Programming Language Second Edition*. Prentice-Hall, Inc., 1988. 6
- [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. 7
- [LBR99] Gary T. Leavens, Albert L. Baker, and Clyde Ruby. Jml : A notation for detailed design, 1999. 25
- [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. 22
- [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. 22
- [Mer03] J. Merrill. GENERIC and GIMPLE : a new tree representation for entire functions. In *GCC developers summit 2003*, pages 171–180, 2003. 7
- [Mil78] Robin Milner. A theory of type polymorphism in programming. *Journal of Computer and System Sciences*, 17(3) :348–375, December 1978. 5
- [MMR⁺13] Anil Madhavapeddy, Richard Mortier, Charalampos Rotsos, David Scott, Balraj Singh, Thomas Gazagnaire, Steven Smith, Steven Hand, and Jon Crowcroft. Unikernels : library operating systems for the cloud. In *Proceedings of the eighteenth international conference on Architectural support for programming languages and operating systems, ASPLOS '13*, pages 461–472, New York, NY, USA, 2013. ACM. 6
- [NCH⁺05] George C. Necula, Jeremy Condit, Matthew Harren, Scott McPeak, and Westley Weimer. Ccured : type-safe retrofitting of legacy software. *ACM Trans. Program. Lang. Syst.*, 27(3) :477–526, May 2005. 23
- [New00] Tim Newsham. Format string attacks. *Phrack*, 2000. 23
- [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. 7
- [NPW02] Tobias Nipkow, Lawrence C. Paulson, and Markus Wenzel. *Isabelle/HOL — A Proof Assistant for Higher-Order Logic*, volume 2283 of *LNCS*. Springer, 2002. 25
- [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. 43
- [Pie02] Benjamin C. Pierce. *Types and Programming Languages*. MIT Press, 2002. 5

- [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. 7
- [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. 20
- [Ric53] H. G. Rice. Classes of recursively enumerable sets and their decision problems. *Transactions of the American Mathematical Society*, 74(2) :pp. 358–366, 1953. 19
- [Spe05] Brad Spengler. grsecurity 2.1.0 and kernel vulnerabilities. *Linux Weekly News*, 2005. 20
- [SRH95] Mooly Sagiv, Thomas Reps, and Susan Horwitz. Precise interprocedural data-flow analysis with applications to constant propagation, 1995. 21
- [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. 7
- [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. 23
- [SY86] R E Strom and S Yemini. Tpestate : A programming language concept for enhancing software reliability. *IEEE Trans. Softw. Eng.*, 12(1) :157–171, January 1986. 23
- [Tan07] Andrew S. Tanenbaum. *Modern Operating Systems*. Prentice Hall Press, Upper Saddle River, NJ, USA, 3rd edition, 2007. 2
- [The04] The Coq Development Team. *The Coq Proof Assistant Reference Manual – Version V8.0*, April 2004. <http://coq.inria.fr>. 25
- [TJ92] Jean-Pierre Talpin and Pierre Jouvelot. Polymorphic type, region and effect inference. *Journal of Functional Programming*, 2 :245–271, 1992. 24
- [TT93] Mads Tofte and Jean-Pierre Talpin. A theory of stack allocation in polymorphically typed languages. Technical report, 1993. 24
- [TT94] Mads Tofte and Jean-Pierre Talpin. Implementation of the typed call-by-value λ -calculus using a stack of regions. In *Proceedings of the 21st ACM SIGPLAN-SIGACT symposium on Principles of programming languages*, POPL '94, pages 188–201, New York, NY, USA, 1994. ACM. 24
- [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. 22
- [vL11] Twan van Laarhoven. Lenses : viewing and updating data structures in Haskell. <http://www.twanvl.nl/files/lenses-talk-2011-05-17.pdf>, May 2011. 37

Résumé

Les noyaux de système d'exploitation manipulent des données fournies par les programmes utilisateur via certaines interfaces, comme les appels système. Si celles-ci sont manipulées sans prendre une attention particulière, une faille de sécurité connue sous le nom de *Confused Deputy Problem* peut amener à des fuites de données confidentielles ou l'élévation de privilège d'un attaquant.

Le but de cette thèse est d'utiliser des techniques de typage statique afin de détecter les manipulations dangereuses de pointeurs contrôlés par l'espace utilisateur.

De nombreux systèmes d'exploitation, dont le noyau Linux sur lequel nous nous basons, sont écrits dans le langage C. Afin de pouvoir appliquer les techniques de typage statique à ce langage, on commence par isoler un sous-langage sûr nommé SAFESPEAK, pour lequel on définit une sémantique opérationnelle et un système de types. Les propriétés classiques de sûreté du typage sont établies.

Dans un second temps on étend SAFESPEAK pour intégrer la notion de valeur provenant de l'espace utilisateur. On montre que l'extension triviale du système de types brise la sûreté du typage, et que l'on peut rétablir celle-ci en donnant un type particulier aux pointeurs dont la valeur est contrôlée par l'espace utilisateur, forçant leur déréférencement à se faire au sein d'un ensemble limité de fonctions sûres. Un cas d'étude est déroulé pour montrer que cette technique permet de détecter un bug qui a frappé un pilote de carte graphique AMD dans le noyau Linux.

Ces travaux ouvrent des perspectives sur deux points. D'une part, il est possible d'appliquer cette analyse à de plus grandes parties du noyau Linux. Les diverses interfaces qui exposent des données contrôlées par l'espace utilisateur peuvent bénéficier de cette vérification statique. Ceci est d'autant plus vrai que les pilotes de périphériques dont elles sont issues sont en général moins testées que le reste du code du noyau. D'autre part, cette technique revient à introduire des types abstraits opaques, séparant leur représentation de leur interface. Ce concept étant absent de C, de nombreuses erreurs de programmation notamment liées aux manipulations d'entiers peuvent également être vérifiées à l'aide de cette méthode.

Abstract

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Donec gravida, diam ullamcorper volutpat varius, est dui blandit nibh, et posuere nisi tortor ut lectus. Aliquam in ipsum at nisl luctus pulvinar. Nulla pulvinar accumsan interdum. Aliquam varius est et ligula commodo varius quis gravida justo. Ut sit amet eros orci, at iaculis metus. Cras placerat condimentum magna ut dictum. Mauris id sollicitudin quam. Ut pretium aliquet lacus, et elementum odio blandit quis.

Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Morbi tincidunt, justo sed venenatis imperdiet, diam dui blandit nisl, et sodales urna mauris sit amet felis. Sed congue sodales sollicitudin. Integer imperdiet tellus et ligula consequat quis commodo metus posuere. Donec id vehicula mi. Nullam id tincidunt orci. Nunc luctus ligula ut urna laoreet consequat. In hac habitasse platea dictumst. Sed imperdiet ultricies diam id cursus. Quisque metus orci, tempor et dapibus eget, ullamcorper in augue. Nullam non urna metus, id accumsan lacus. Duis erat arcu, porttitor non laoreet vitae, ultrices vitae tortor. Pellentesque interdum faucibus tellus eu semper. Nulla ac aliquet mauris.

In in metus purus. Nunc ullamcorper posuere turpis, vel laoreet tortor posuere at. Ut euismod arcu sit amet dui varius vel varius neque pellentesque. Lorem ipsum dolor sit amet, consectetur adipiscing elit. Aliquam nec mi interdum ipsum eleifend tincidunt. Duis bibendum sem vitae risus facilisis mattis. Sed varius eleifend nulla vel commodo. Morbi risus metus, egestas non commodo id, ullamcorper a enim. Integer et augue et orci feugiat rhoncus a in augue.

Nulla elementum felis eu eros pretium quis porttitor felis faucibus. Vivamus a vulputate tellus. Suspendisse potenti. Morbi dapibus, quam ac consectetur hendrerit, lorem urna faucibus dui, in vestibulum tellus nisl sit amet ligula. Morbi ultricies venenatis nisi sit amet scelerisque. Nam pharetra bibendum lacinia. Suspendisse potenti. Suspendisse elementum tellus et mi hendrerit sagittis.

Ut dignissim fermentum cursus. Praesent feugiat venenatis rutrum. Cras eget nibh id quam ultrices pellentesque eget non elit. Maecenas justo lorem, fermentum vitae dictum eu, egestas eu urna. In hac habitasse platea dictumst. Cras nec arcu at purus aliquet accumsan. Nam interdum malesuada orci et sagittis