

Master de recherche en informatique,  
ISTIC — ÉNS Cachan, antenne de Bretagne  
**Rapport**

Encadrant : David Cachera  
Équipe Celtique  
Inria Rennes — Irista

# Vérification probabiliste de résultats d'analyse statique

---

Antoine Bride

Rennes, 5 juin 2012

## Résumé

Ce document traite de la vérification probabiliste de résultats d'analyse statique. Il s'agit plus précisément de mettre au point un algorithme de Monte-Carlo permettant de vérifier qu'un objet (appelé certificat) que l'on nous donne est bien le résultat d'une analyse statique. À notre connaissance, aucun algorithme de la sorte n'existe. Nous avons donc commencé par nous définir un cas précis sur lequel travailler à savoir une analyse *Points-to* sur un langage impératif simple que nous définissons. Nous avons ensuite proposé un algorithme de Monte-Carlo qui soit simple à implémenter et rapide à exécuter pour ce cas en particulier. Une grande partie de notre travail a alors consisté à démontrer que l'algorithme avait de bonnes propriétés de convergence en terme de probabilité d'erreur. Nous avons enfin implémenté cet algorithme afin d'obtenir des résultats expérimentaux. Finalement, nous proposons des pistes pour élargir le champ d'application et ne plus se restreindre à un cas précis.

**Mots-clefs :** Analyse de programmes, Analyse statique, Algorithmique probabiliste, Analyse *Points-to*, Monte-Carlo

## Table des matières

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Travaux Connexes</b>	<b>4</b>
<b>3</b>	<b>Un langage de pointeur</b>	<b>5</b>
3.1	Syntaxe du langage . . . . .	6
3.2	Domaines sémantiques du langage . . . . .	7
3.3	Sémantiques des sous-expressions du langage . . . . .	7
3.4	Sémantique du Langage . . . . .	8
<b>4</b>	<b>Sémantique abstraite du langage</b>	<b>10</b>
4.1	Abstraction par rapport à la sémantique de langage . . . . .	10
4.2	Analyse dans la sémantique abstraite . . . . .	11
<b>5</b>	<b>Algorithme de vérification de résultats de l'analyse</b>	<b>12</b>
5.1	Conception de l'algorithme . . . . .	13
5.2	Vérification probabiliste d'une inclusion . . . . .	15
5.3	Vérification de l'ensemble du résultat . . . . .	16
5.4	Indépendance des inclusions . . . . .	18
<b>6</b>	<b>Implémentation et résultats expérimentaux</b>	<b>19</b>
6.1	Implémentation . . . . .	19
6.2	Résultats expérimentaux . . . . .	20
<b>7</b>	<b>Ouvertures</b>	<b>21</b>
7.1	L'utilisation de martingales . . . . .	21
7.2	Amélioration de l'implémentation . . . . .	25
<b>8</b>	<b>Conclusion</b>	<b>26</b>

# 1 Introduction

Actuellement, quand on obtient un logiciel, on n'a presque jamais de garantie que le logiciel en question se comporte comme il faut. Le consommateur n'a alors plus qu'une seule alternative : faire confiance au producteur, ou prouver lui-même que le logiciel se comporte bien. Cependant, la plupart des consommateurs n'ont pas la patience ou pas les compétences requises en informatique par faire le second choix. Un moyen de gagner du temps est d'exiger du producteur qu'il réalise une analyse de programme sur son logiciel afin de prouver que ce dernier se comporte bien.

Les méthodes d'analyse de programme sont faites pour étudier le comportement d'un programme, en établissant des propriétés sur celui-ci. On peut séparer ces méthodes en deux grandes catégories. La première catégorie, l'analyse dynamique (ou test), consiste à faire fonctionner un programme afin de voir ce qu'il se passe. Par exemple, on peut lancer un programme avec des valeurs d'entrée arbitraires sur une machine virtuelle et regarder comment la mémoire est affectée. Cette classe d'analyse n'est pas traitée dans ce document car elle n'est pas faite pour garantir que *rien de grave ne peut arriver*.

La seconde classe d'analyse de programmes, l'analyse statique, est faite pour obtenir ce type de garantie. Elle ne nécessite pas l'exécution du programme. Le principe est d'analyser le code<sup>1</sup> d'un programme afin d'établir des propriétés sur son comportement. Un problème important dans ce type d'approche est l'indécidabilité de telles propriétés. Par exemple, la question de savoir si un programme termine ou non est indécidable. En réponse à ce problème, on définit un modèle qui surapproxime le comportement d'un programme. Le but est que les propriétés que l'on cherche à établir soient décidables dans le modèle. Si les propriétés garantissent que le modèle se comporte bien, alors le programme se comporte bien car le modèle le surapproxime. Plus précisément, établir des propriétés sur le modèle consiste souvent à obtenir des équations sur les variables du modèle et à les résoudre. Cette dernière étape est souvent coûteuse en terme de calculs car elle consiste généralement à réaliser un calcul de point fixe sur un treillis.

On peut alors imaginer qu'un producteur fournissant un logiciel java réalise une analyse statique (publique) sur le `byte code java` de celui-ci et fournisse donc, en sus du `byte code`, le résultat de l'analyse en tant que certificat. Le producteur garantirait alors le bon comportement de son logiciel sans avoir à en dévoiler le code source. L'acquéreur, lui, pourrait vérifier que le certificat est bien conforme à l'analyse statique faite sur le logiciel.

Partant toujours du principe que l'acquéreur d'un programme n'est pas patient. Il nous semblait important réaliser un algorithme qui puisse s'exécuter rapidement. Pour cette raison, nous avons souhaité réaliser un algorithme probabiliste pour vérifier le certificat. L'inconvénient de tels algorithmes est qu'ils peuvent fournir un résultat erroné (Algorithmes de Monte-Carlo) ou être beaucoup plus longs à terminer que prévu (Algorithmes de Las Vegas). De plus, deux propriétés des algorithmes probabilistes nous ont semblé intéressantes.

La première de ces propriétés est qu'un algorithme probabiliste est souvent, en soit, très simple. Tout la difficulté repose dans la preuve que l'algorithme se comporte bien. Rajeev Motwani et Prabhakar Raghavan expliquent, dans leur ouvrage *Randomized Algorithms* [MR95], que les algorithmes déterministes ont une certaine tendance à toujours être plus compliqués pour réduire leur complexité et à jouer sur des propriétés mathématiques avancées qu'il faut comprendre lors de l'implémentation. Au contraire, les algorithmes probabilistes ont l'avantage de rester simple. Le principe est que dans un programme, il y a souvent un choix à faire qui détermine le succès du programme ou son temps d'exécution. Par exemple, pour l'algorithme de tri rapide, le point dur est le choix du pivot. Les algorithmes probabilistes proposent alors souvent de faire quelque chose de très peu compliqué : choisir au hasard en comptant sur le fait qu'il n'y a que peu de choix vraiment

---

<sup>1</sup>Code source, code assembleur, `byte code java`...

mauvais. Toute la difficulté est alors de prouver qu’effectivement, choisir au hasard donne de bon résultats. Ceci est transparent pour qui souhaite implémenter l’algorithme.

Le deuxième avantage de cette sorte de paradigme est d’éviter les « mauvaises » entrées. Pour reprendre l’exemple du tri rapide, lorsque l’on connaît la manière déterministe dont le pivot est choisi, on peut toujours créer une entrée qui le fera s’exécuter en un temps proche du temps quadratique. Si le pivot est choisit au hasard, en revanche, il est impossible de créer une entrée qui garantira un temps d’exécution quadratique. Le choix aléatoire du pivot peut bien sûr être malchanceux et sélectionner systématiquement le plus grand élément du tableau comme pivot. Cependant, cela ne dépend plus du tableau. Là où l’on compte sur le fait que les entrées sont bien réparties dans le cas déterministe, on ne compte plus que sur le fait que le générateur aléatoire, que l’on contrôle bien plus, donne un bon aléas dans le cas probabiliste.

Plus particulièrement nous nous sommes intéressés à l’élaboration d’un algorithme de Monte-Carlo. Il y a plusieurs raisons pour lesquelles nous avons choisit de faire cela. Premièrement, un tel algorithme permet de choisir à quel point la vérification sera rapide<sup>2</sup>. Deuxièmement, un tel algorithme n’a pas nécessairement besoin de connaître toutes propriétés données par le certificat, seulement celles qu’il vérifie. La quantité de propriétés étant éventuellement grande, cela permettrait de gagner du temps en ne demandant que celles que l’on souhaite vérifier au producteur.

À notre connaissance, aucun algorithme de ce type n’existe à ce jour. La section 2 porte donc sur l’algorithmique probabiliste en général, un algorithme probabiliste dans le domaine de l’analyse de programme et des analyses statiques proches de celle que nous considérons. De plus, nous spécifions formellement le problème que nous nous posons. D’abord nous définissons le langage que nous considérons (section 3), ensuite l’analyse qui est faite sur le langage (section 4). Une fois cela fait, nous pouvons définir notre algorithme (section 5) et discuter de son implémentation et de ses résultats (section 6). Enfin nous parlerons de travaux ultérieurs pouvant être réalisés (section 7) et nous concluons (section 8).

## 2 Travaux Connexes

En terme de travaux en ce qui concerne l’algorithmique probabiliste en analyse de programme, on peut citer notamment ceux de Sumit Gulwani et la *Random Interpretation* [Gul05, GN03, GN04a, GN04b, GN05]. Son but est de cumuler les avantages de l’analyse dynamique (ne capturer que le comportement du programme) et de l’analyse statique (capturer l’intégralité du comportement du programme). Pour cela Gulwani propose de faire tourner le programme sur un jeu d’entrées comme en analyse dynamique, mais de s’intéresser à tous les branchements possibles (notamment lors de conditionnelles ou de boucles) pour capturer le maximum de comportements possibles. Afin de toujours garder un seul état des variables à chaque endroit du programme, lorsque deux branches se réunissent, les états courants de chaque branche sont fusionnés grâce à une fonction aléatoire. Il s’agit en réalité de réaliser une somme pondérée des états de poids aléatoire. À savoir, si une variable  $x$  à la valeur  $v_1$  dans une branche et  $v_2$  dans un autre, lors de la réunion des deux branches, on tire un poids  $w$  aléatoirement et  $x$  reçoit la valeur  $wv_1 + (1 - w)v_2$ . L’intérêt, étant donné que  $w$  est tiré une fois pour toutes les variables lors de la fusion de deux branches, est que les relations linéaires<sup>3</sup> sont conservées puisque la somme est elle même linéaire. Gulwani montre dans sa thèse [Gul05] qu’il parvient à obtenir un algorithme pro-

---

<sup>2</sup>En revanche, plus la vérification est rapide, plus l’algorithme a de chances de se tromper.

<sup>3</sup>C’est à dire les relations du type  $a_1x_1 + \dots + a_nx_n = b$  où les  $x_i$  sont des variables du programmes et les  $a_i$  ainsi que  $b$  sont des constantes.

babiliste qui capture toutes les relations linéaires d'un programme et qui ne capture que celles-là avec une probabilité élevée. Il montre notamment que sa méthode peut s'adapter dans le cas de l'arithmétique affine [GN03], dans l'abstraction de fonctions [GN04a] et de l'utilisation de fonctions déclarées par l'utilisateur [GN05]. Une des principales forces de ses travaux est que l'algorithme qu'il propose est rapide et, en soit, très simple. Seule la preuve que l'algorithme se comporte bien est difficile.

Une remarque que l'on peut faire sur la sous-catégorie des algorithmes probabilistes qui nous intéresse, à savoir les algorithmes de Monte-Carlo, est qu'ils sont souvent utilisés pour répondre à une question fermée. Dans notre cas, est-ce que oui ou non, un certificat correspond à l'analyse d'un programme. Une propriété pratique apparaît si l'on peut montrer que l'algorithme ne se trompe jamais quand il donne l'une de ses réponses. Dans notre cas, imaginons que notre algorithme consiste à vérifier que le certificat est conforme à un sous-ensemble des propriétés données par l'analyse, et réponde oui si c'est le cas, non sinon. Notre algorithme répond non uniquement si il a trouvé au moins une propriété que le certificat ne vérifie pas et ne peut donc pas se tromper dans ce cas. La remarque, faite par Motwani *et al.* [MR95], est alors la suivante : si l'algorithme à une probabilité  $p < 1$  de répondre oui à tort, on lance l'algorithme  $n$  fois de suite, l'algorithme a alors une probabilité  $p^n$  de répondre  $n$  fois oui à tort (et si il répond non, une fois, on est fixé sur le résultat). Autrement dit, nous venons d'obtenir presque gratuitement la décroissance exponentielle de la probabilité d'erreur de notre algorithme contre une croissance linéaire de son temps d'exécution. Nous n'avons eu à définir formellement ni l'algorithme, ni les propriétés, ni même l'analyse.

Afin, d'obtenir une quantification précise de cette probabilité d'erreur de l'algorithme, il nous est tout de même nécessaire de définir formellement notre cadre. Nous avons choisi de travailler sur un langage de pointeur, c'est pourquoi nous allons maintenant discuter de quelques manières d'aborder l'analyse *points-to* avant de commencer les définitions formelles en section 3. Une analyse *Points-to* a pour but d'établir vers quelle adresse peut pointer chaque variable. Elle se focalise donc sur toutes les manipulations d'adresse d'un programme.

Parmi d'autres travaux, nous pouvons citer ceux de Rupesh Nasre et Ramaswamy Govindarajan qui proposent de modéliser les équations venant d'une analyse *points-to* par des systèmes d'équations linéaires [NG11]. L'inconvénient est cette méthode est que traduire une équation de l'analyse *points-to* à l'algèbre linéaire a un coût. L'avantage est que la résolution de systèmes d'équation linéaire est un problème ancien auquel beaucoup de solutions efficaces ont été proposées. On peut citer aussi les travaux de John Whaley et Monica S. Lam qui ont proposé d'utiliser des diagrammes de décision binaire pour résoudre les équations d'une analyse *points-to* [WL04]. Ils ont ensuite continué leurs travaux avec l'aide de Dzintars Avots et Michael Carbin pour réaliser un compilateur de `dataLog` (un sous-langage de `prolog`) qui est dédié à l'analyse statique en général et à l'analyse *points-to* en particulier [WACL05]. Ce compilateur nous a été utile. En effet, nous avons pu l'utiliser afin d'obtenir de vrais résultats de notre analyse statique. Cela nous a permis de modifier légèrement ces résultats pour tester le comportement de l'algorithme que nous avons implémenté.

### 3 Un langage de pointeur

Maintenant que d'autres approches d'analyse *points-to* ont été présentées, nous proposons la nôtre. Dans cette section, on définit donc le langage de pointeurs sur lequel une analyse sera faite. On commence par donner la syntaxe du langage (section 3.1), puis sa sémantique (sections 3.2, 3.3 et 3.4).

$\epsilon \in \mathbf{ExprVar}$ $\epsilon := \begin{array}{ l} x & x \in \mathbf{Var} \\ *x & x \in \mathbf{Var} \end{array}$ <p>(a) Expressions de variables</p> $e \in \mathbf{ExprVal}$ $e := \begin{array}{ l} x & x \in \mathbf{Var} \\ *x & x \in \mathbf{Var} \\ \&x & x \in \mathbf{Var} \end{array}$ <p>(b) Expressions de valeurs</p>	$b \in \mathbf{Bool}$ $b := \begin{array}{ l} \text{true} \\ \text{false} \\ \text{not } b & b \in \mathbf{Bool} \\ b_1 \text{ o } b_2 & b_1, b_2 \in \mathbf{Bool}, \text{ o} \in \{\wedge, \vee\} \\ \epsilon_1 == \epsilon_2 & \epsilon_1, \epsilon_2 \in \mathbf{ExprVar} \\ e_1 = e_2 & e_1, e_2 \in \mathbf{ExprVal} \end{array}$ <p>(c) Expressions booléennes</p>
---	---

FIG. 1 – Définition des différentes expressions présentes dans le langage

### 3.1 Syntaxe du langage

La syntaxe du langage de pointeur que l'on souhaite définir est celle d'un langage impératif standard à ceci près qu'elle considère que les variables contiennent des adresses et non des entiers. On définit d'abord les domaines syntaxiques, puis les syntaxes des sous-expressions : expressions de variables, d'adresses, et booléennes. Enfin, on considère la syntaxe d'un programme.

Comme évoqué précédemment, le langage utilisera des variables ; on a donc besoin d'un ensemble de variables **Var**. On définit de plus un ensemble d'étiquettes **Label** qui permettra d'identifier les expressions entre elles.

On considère maintenant les sous-expressions. Pour commencer, on définit l'ensemble des expressions de variables : **ExprVar** (figure 1a). On choisit de pouvoir exprimer une variable soit par la variable elle-même, soit par le déréférencement d'une autre variable à l'aide de l'opérateur  $*$ .

Ensuite, on définit l'ensemble des expressions de valeurs : **ExprVal** (figure 1b). Une valeur peut être exprimée de trois manières : à l'aide d'une variable, par le truchement d'un déréférencement (opérateur  $*$ ), ou en référençant une variable avec l'opérateur  $\&$ .

Enfin, on définit l'ensemble des expressions booléennes : **Bool** (figure 1c). Cet ensemble est défini inductivement de manière usuelle à la différence que l'on considère deux égalités. La première est le test d'égalité de variables :  $\epsilon_1 == \epsilon_2$ ,  $\epsilon_1, \epsilon_2 \in \mathbf{ExprVar}$ . Dans la section 3.3, on définira la sémantique du test  $==$  de sorte que si  $\epsilon_1 == \epsilon_2$  est évaluée à *true*, cela signifie que modifier la valeur de  $\epsilon_1$  revient à modifier celle de  $\epsilon_2$  (et réciproquement). La seconde égalité est le test d'égalité de valeurs :  $e_1 = e_2$ ,  $e_1, e_2 \in \mathbf{ExprVal}$ .

La syntaxe des sous-expressions ayant été définie, on considère maintenant le langage lui-même. Le langage, **Statement**, possède une syntaxe basique de langage impératif (figure 2). Il contient les instructions suivantes : **skip**, l'assignation d'une valeur à une variable, la séquence, un branchement conditionnel et une boucle. Les instructions sont étiquetées.

La syntaxe du langage ayant été définie, il faut maintenant donner un sens à toutes ces expressions. Dans les sous-sections suivantes, on définit donc la sémantique du programme. On considère en premier lieu les domaines sémantiques 3.2, ensuite les sémantiques des sous-expressions 3.3 et enfin la sémantique du programme lui-même 3.4.

$S \in \mathbf{Statement}, l \in \mathbf{Label}$	
$S :=$	$l[\text{skip}]$
	$l[\epsilon := e]$ $\epsilon \in \mathbf{ExprVar} \quad e \in \mathbf{ExprVal}$
	$S_1; S_2$ $S_1, S_2 \in \mathbf{Statement}$
	$\text{if } l[b] \text{ then } S_1 \text{ else } S_2 \text{ endif}$ $b \in \mathbf{Bool} \quad S_1, S_2 \in \mathbf{Statement}$
	$\text{while } l[b] \text{ do } S \text{ end}$ $b \in \mathbf{Bool} \quad S \in \mathbf{Statement}$

FIG. 2 – Syntaxe du langage

## 3.2 Domaines sémantiques du langage

Pour commencer, on définit donc les domaines sémantiques à savoir les ensembles dans lesquels les sémantiques des expressions et du programme prendront leurs arguments et leurs valeurs.

Premièrement, trois ensembles sont définis : **Add**, l'ensemble des adresses utilisées, **Value** l'ensemble des valeurs que peuvent contenir les zones mémoires (correspondant à des adresses) et  $\{\text{tt}, \text{ff}\}$  qui correspondent respectivement à vrai et à faux. Étant donné que l'on considère un langage de pointeur, on pose **Value** := **Add**. Cette égalité étant spécifique au cas présent et afin d'éviter toute confusion, nous utilisons tout de même les noms **Add** pour les adresses en tant que telles et **Value** pour les valeurs. Lorsque l'égalité entre ces deux ensembles sera nécessaire, nous l'utiliserons de manière explicite.

Dans le paragraphe précédent nous avons implicitement évoqué l'ensemble des états : **State** = **Add**  $\rightarrow$  **Value**  $\cup$   $\{\perp\}$ . Un état ( $\sigma \in \mathbf{State}$ ) associe à chaque adresse une valeur (que contient donc la zone mémoire correspondant à l'adresse) ou  $\perp$  si la zone mémoire correspondant à l'adresse n'a pas été initialisée.

Afin d'obtenir une sémantique, il est nécessaire de lier les éléments syntaxiques (**Var**) et sémantiques (**Add** et **Value**). Pour cela, on définit l'ensemble des environnements **Env** = **Var**  $\rightarrow$  **Add**  $\cup$   $\{\perp\}$  qui à chaque variable, associe l'adresse à laquelle elle correspond (ou  $\perp$  si la variable n'a pas d'adresse). Les environnements sont noté  $\rho$ .

Il est désormais possible de définir les sémantiques des sous-expressions.

## 3.3 Sémantiques des sous-expressions du langage

On définit d'abord la sémantique des expressions de variables, puis celle des expressions de valeurs, et enfin celle des expressions booléennes.

Pour commencer, on s'intéresse donc à la sémantique des expressions de variables, à savoir  $\mathcal{E}_a$  (figure 3a). Deux cas peuvent se présenter :  $x$  et  $*x$  (où  $x \in \mathbf{Var}$ ). Lorsque l'on utilise la variable  $x \in \mathbf{Var}$ , on utilise en fait l'adresse qui y correspond :  $\rho x$ . Lorsque l'on déréférence la variable  $x \in \mathbf{Var}$ , on s'intéresse à la valeur que « contient  $x$  » à savoir la valeur se trouvant à l'adresse correspondant à  $x$  :  $\sigma(\rho x)$ . Dans ce dernier cas, il faut que la valeur soit une adresse, c'est le cas car **Value** = **Add**.

Ensuite, on définit la sémantique des expressions de valeurs :  $\mathcal{E}_v$  (figure 3b). Trois cas peuvent échoir :  $x$ ,  $*x$ ,  $\&x$  et (où  $x \in \mathbf{Var}$ ). La valeur qui correspond à une variable  $x$  est celle que contient son adresse, à savoir  $\sigma(\rho x)$ . La valeur que l'on souhaite obtenir avec un déréférencement  $*x$  est la valeur se trouvant à l'adresse<sup>4</sup> que « contient »  $x$ . La sémantique de  $*x$  est donc  $\sigma(\sigma(\rho x))$  (pour  $x \in \mathbf{Var}$ ). Enfin, lorsque l'on référence une variable  $x$  avec  $\&x$ , on s'intéresse simplement à l'adresse de cette variable :  $\rho x$ .

Finalement, on utilise une sémantique des expressions booléennes usuelle,  $\mathcal{E}_b$  (figure 3c). Comme annoncé dans la section 3.1, la sémantique de  $\epsilon_1 == \epsilon_2$ ,  $\epsilon_1, \epsilon_2 \in \mathbf{ExprVar}$ ,

<sup>4</sup>C'est en réalité une valeur mais **Value** = **Add**.

$$\begin{array}{ll}
\mathcal{E}_a : \mathbf{ExprVar} \rightarrow \mathbf{State} \rightarrow \mathbf{Env} \rightarrow \mathbf{Add} & \mathcal{E}_v : \mathbf{ExprVal} \rightarrow \mathbf{State} \rightarrow \mathbf{Env} \rightarrow \mathbf{Value} \\
\mathcal{E}_a[x] \sigma \rho = \rho x & \mathcal{E}_v[x] \sigma \rho = \sigma(\rho x) \\
\mathcal{E}_a[*x] \sigma \rho = \sigma(\rho x) & \mathcal{E}_v[*x] \sigma \rho = \sigma(\sigma(\rho x)) \\
& \mathcal{E}_v[\&x] \sigma \rho = \rho x
\end{array}$$

(a) Expressions de variables

(b) Expressions de valeurs

$$\begin{array}{l}
\mathcal{E}_b : \mathbf{Bool} \rightarrow \mathbf{State} \rightarrow \mathbf{Env} \rightarrow \{\mathbf{tt}, \mathbf{ff}\} \\
\mathcal{E}_b[\mathbf{true}] \sigma \rho = \mathbf{tt} \\
\mathcal{E}_b[\mathbf{false}] \sigma \rho = \mathbf{ff} \\
\mathcal{E}_b[\epsilon_1 == \epsilon_2] \sigma \rho = \begin{cases} \mathbf{tt} & \text{si } (\mathcal{E}_a[\epsilon_1] \sigma \rho) = (\mathcal{E}_a[\epsilon_2] \sigma \rho) \\ \mathbf{ff} & \text{sinon} \end{cases} \\
\mathcal{E}_b[e_1 = e_2] \sigma \rho = \begin{cases} \mathbf{tt} & \text{si } (\mathcal{E}_v[e_1] \sigma \rho) = (\mathcal{E}_v[e_2] \sigma \rho) \\ \mathbf{ff} & \text{sinon} \end{cases} \\
\mathcal{E}_b[b_1 \circ b_2] \sigma \rho = (\mathcal{E}_b[b_1] \sigma \rho) \llbracket \circ \rrbracket (\mathcal{E}_b[b_2] \sigma \rho)
\end{array}$$

(c) Expressions booléennes

FIG. 3 – Sémantique des différentes expressions présentes dans le langage

est évaluée à  $\mathbf{tt}$  lorsque  $\mathcal{E}_a[\epsilon_1] \sigma \rho = \mathcal{E}_a[\epsilon_2] \sigma \rho$  et à  $\mathbf{ff}$  sinon. Donc, si on met une nouvelle valeur dans l'adresse correspondant à  $\epsilon_1$ , on change la valeur dans l'adresse correspondant à  $\epsilon_2$  puisqu'il s'agit de la même.

Les sémantiques des sous-expressions étant définies, on peut désormais s'intéresser à la sémantique d'un programme  $S \in \mathbf{Statement}$ .

### 3.4 Sémantique du Langage

Notre but est de vérifier qu'un programme n'atteint jamais un état « dangereux ». C'est pour cela que l'on s'intéresse à la sémantique collectrice. Il s'agit en fait de collecter tous les états possiblement atteignables lors de l'exécution d'un programme.

Pour réaliser une telle sémantique, on procède en deux étapes : d'abord on construit un graphe de flot de contrôle à partir de  $S$ , puis on « navigue » sur un tel graphe en respectant les conditions données par les arrêtes pour obtenir la sémantique du programme. La raison d'un tel choix est que cela permet d'obtenir une séparation claire entre les éléments structurant d'un programme (tel que la boucle) et les éléments affectant la mémoire. En effet, les premiers sont capturés dans la structure du graphe tandis que les seconds sont gérés par la façon dont le graphe est parcouru.

On s'intéresse à la première étape. On souhaite donc générer un graphe de flot de contrôle  $\mathit{cfg} S$  que l'on définit par un quadruplet  $(\mathbf{Vertex}, e, f, \mathbf{Edge})$  où  $\mathbf{Vertex}$  est l'ensemble des nœuds du graphe,  $e, f \in \mathbf{Vertex}$  sont respectivement l'entrée et la sortie du graphe, et  $\mathbf{Edge}$  l'ensemble des arêtes du graphe. On pose que les nœuds sont les étiquettes du programme, donc  $\mathbf{Vertex} := \mathbf{Label}$ . Deux instructions qui partagent la même étiquette ne sont donc pas distinguées. En ce qui concerne les arêtes, on définit d'abord  $\mathbf{Inst} = \{\mathbf{skip}, \epsilon := e, \mathbf{assert} b\}$ . Cet ensemble comprend les « instructions » atomiques d'un programme, celles ayant un effet ou des conditions sur l'état courant. Le principe est de lier les arêtes à des instructions de sorte que l'on ne prendra les arêtes qu'en respectant l'effet des instructions sur la mémoire.

Les nœuds ont été définis,  $e$  est simplement la première étiquette rencontrée dans  $S$  (appelée  $\mathit{entry} S$ ),  $f$  est choisie arbitrairement ; seules les arêtes sont toujours à définir. On



$$\begin{aligned}
\text{cfg}_l : \mathbf{Statement} &\rightarrow \mathcal{P}(\mathbf{Label} \times \mathbf{Inst} \times \mathbf{Label}) \\
\text{cfg}_l ({}^l[\epsilon := e]) &= \{(l, \epsilon := e, l')\} \\
\text{cfg}_l ({}^l[\text{skip}]) &= \{(l, \text{skip}, l')\} \\
\text{cfg}_l (\text{if } {}^l[b] \text{ then } S_1 \text{ else } S_2 \text{ endif}) &= \{(l, \text{assert } b, \text{entry } S_1)\} \cup \text{cfg}_l S_1 \\
&\quad \cup \{(l, \text{assert } (\text{not } b), \text{entry } S_2)\} \cup \text{cfg}_l S_2 \\
\text{cfg}_l (\text{while } {}^l[b] \text{ do } S \text{ end}) &= \{(l, \text{assert } b, \text{entry } S)\} \cup \text{cfg}_l S \\
&\quad \cup \{(l, \text{assert } (\text{not } b), l')\} \\
\text{cfg}_l (S_1; S_2) &= \text{cfg}_{\text{entry } S_2} S_1 \cup \text{cfg}_l S_2
\end{aligned}$$

FIG. 4 – Génération des arrêtes du graphe de flot de contrôle d'un programme

définit les arêtes comme des triplets de  $\mathbf{Vertex} \times \mathbf{Inst} \times \mathbf{Vertex}$ . Intuitivement, on souhaite que l'arête  $(v_1, i, v_2)$  corresponde au fait de pouvoir aller de l'étiquette  $v_1$  à l'étiquette  $v_2$  en effectuant l'instruction  $i$ . On génère les arêtes de  $\text{cfg}_l S$  inductivement en fonction des cas pouvant échoir à l'aide de fonction  $\text{cfg}_l$  (figure 4) où  $l$  en indice est une étiquette. Le principe est de garder en mémoire où l'on souhaite se rendre (avec l'indice) et de construire les arêtes permettant de se rendre de l'étiquette courante (donnée par le programme en argument) vers l'étiquette cible (indice). Par exemple  $\text{cfg}_l (\text{while } {}^l[b] \text{ do } S \text{ end})$  est la réunion de trois ensembles :

- $\{(l, \text{assert } b, \text{entry } S)\}$  car si  $b$  est vérifié on se branche sur le corps de la boucle :  $S$ ,
- $\text{cfg}_l S$  car lorsqu'on exécute le corps de la boucle et on revient au test,
- $\{(l, \text{assert } (\text{not } b), l')\}$  car si  $b$  n'est pas vérifié ( $\text{assert } (\text{not } b)$ ) on sort de la boucle.

Au final le graphe de flot de contrôle d'un programme  $S \in \mathbf{Statement}$  est le suivant :  $\text{cfg}_l S = (\mathbf{Label}, \text{entry } S, f, \text{cfg}_l S)$ .

Il faut maintenant donner une sémantique à ce graphe de flot de contrôle. Pour cela, on établit des règles quant à la manière de parcourir le graphe à l'aide d'une fonction correspondant à chaque instruction que peut contenir une arête :

$$\begin{array}{c}
\begin{array}{l}
\begin{array}{l}
\frac{}{(\sigma, \rho) \xrightarrow{\text{skip}} (\sigma, \rho)} \\
\frac{\mathcal{E}_b[b] = \text{tt } \sigma \rho}{(\sigma, \rho) \xrightarrow{\text{assert } b} (\sigma, \rho)} \\
\frac{a = \mathcal{E}_a[\epsilon] \sigma \rho \quad v = \mathcal{E}_v[e] \sigma \rho}{(\sigma, \rho) \xrightarrow{\epsilon := e} (\sigma[x \mapsto a], \rho)}
\end{array} \\
\frac{}{(\sigma, \rho) \xrightarrow{\text{skip}} (\sigma, \rho)} \quad \frac{\mathcal{E}_b[b] = \text{tt } \sigma \rho}{(\sigma, \rho) \xrightarrow{\text{assert } b} (\sigma, \rho)} \quad \frac{a = \mathcal{E}_a[\epsilon] \sigma \rho \quad v = \mathcal{E}_v[e] \sigma \rho}{(\sigma, \rho) \xrightarrow{\epsilon := e} (\sigma[x \mapsto a], \rho)}
\end{array}
\end{array}$$

Le sens que l'on donne est qu'une arête avec l'instruction  $\text{skip}$  ne change rien, qu'une arête ayant l'instruction  $\text{assert } b$  ne change rien mais ne peut être empruntée que si  $b$  est vérifiée, et qu'une arête avec l'instruction  $\epsilon := e$  modifie l'état selon les évaluations de  $\epsilon$  et  $e$  lorsqu'elle est empruntée.

On parcourt alors le graphe ainsi : Si l'on est dans l'état  $\sigma_1$  et l'environnement  $\rho_1$  au nœud  $l_1$ , on peut arriver dans l'état  $\sigma_2$  et l'environnement  $\rho_2$  au nœud  $l_2$  si et seulement si l'arête  $(l_1, i, l_2)$  existe et l'on a  $(\sigma_1, \rho_1) \xrightarrow{i} (\sigma_2, \rho_2)$ . Cela est défini formellement par la fonction suivante :

$$\rightarrow_{\text{cfg}} : (\mathbf{Label} \times \mathbf{State} \times \mathbf{Env}) \rightarrow (\mathbf{Label} \times \mathbf{State} \times \mathbf{Env})$$

$$\frac{(k_1, i, k_2) \in \text{cfg}_f S \quad \sigma_1 \xrightarrow{i} \sigma_2}{(k_1, \sigma_1, \rho) \rightarrow_{\text{cfg}} (k_2, \sigma_2, \rho)}$$

Maintenant que l'on sait quels états peuvent être atteints, on peut définir la sémantique collectrice en chaque étiquette d'un programme  $P$  :

$$\forall k \in \mathbf{Label}, \llbracket P \rrbracket_k^{\text{col}} = \{\sigma \mid \exists (l, \sigma_{(0,l)}, \rho) \rightarrow_{\text{cfg}}^* (k, \sigma, \rho)\}$$

Où  $\sigma_{(0,k)}$  est l'état initial associé à chaque étiquette et est une donnée du problème. On peut par exemple n'avoir rien d'initialisé au début :  $\forall l \in \mathbf{label}, \forall a \in \mathbf{Add}, \sigma_{(0,l)} a = \perp$ . On note  $\llbracket P \rrbracket^{\text{col}} \in \mathcal{P}(\mathbf{State})^{\#\mathbf{Label}}$  la sémantique collectrice (globale) dont la  $k^{\text{ième}}$  composante est la sémantique collectrice de l'étiquette  $k$  :  $\llbracket P \rrbracket_k^{\text{col}}$ . On remarque que  $(\mathcal{P}(\mathbf{State})^{\#\mathbf{Label}}, \cup, \cap, \subseteq)$  est un treillis fini.

Il devient alors intéressant d'écrire  $\llbracket P \rrbracket^{\text{col}}$  comme solution du système d'équation suivant :

$$\forall k \in \mathbf{Label}, X_k = X_k^{\text{init}} \cup \bigcup_{(k',i,k) \in P} \llbracket i \rrbracket X_{k'} \quad \llbracket i \rrbracket : X \mapsto \{\sigma_2 \mid \exists \sigma_1 \in X, \sigma_1 \xrightarrow{i} \sigma_2\} \quad (1)$$

Où  $X_k^{\text{init}}$  correspond à  $\sigma_{(0,k)}$ . En effet, le système d'équations ci-dessus peut être vu comme une fonction croissante dans le treillis auquel appartient  $\llbracket P \rrbracket^{\text{col}}$ . Le fait que le treillis soit fini implique qu'en itérant successivement la fonction sur son résultat (en partant des états initiaux), on atteindra son plus petit point fixe en un nombre d'étapes fini. Autrement dit, on dispose d'une manière de résoudre le système d'équation 1.

## 4 Sémantique abstraite du langage

Étant donné qu'il est difficile d'exprimer simplement les propriétés dans la sémantique collectrice, nous souhaitons nous en abstraire. Pour cela, nous utilisons une sémantique abstraite. Il s'agit d'approximer la sémantique réelle du programme de façon à travailler dans des ensembles (et des treillis) plus simples. L'inconvénient est que l'on perd en précision. L'intérêt est que l'on gagne en complexité. Nous allons donc définir une telle sémantique abstraite puis montrer qu'elle a de bonnes propriétés par rapport à la sémantique collectrice. Ensuite, nous allons analyser le système d'équations abstrait pour définir précisément à quelle question devra répondre l'algorithme que nous souhaitons réaliser.

### 4.1 Abstraction par rapport à la sémantique de langage

Afin de réaliser la sémantique abstraite, on commence par définir ses éléments constitutifs : les états abstraits notés  $\varphi$ . Ces états abstraits associent aux adresses l'ensemble des valeurs qu'elles peuvent contenir. Ainsi, pour chaque étiquette, on pourra considérer un seul état abstrait plutôt qu'un ensemble d'états :

$$\begin{aligned} \varphi : \mathbf{Add} &\rightarrow \mathcal{P}(\mathbf{Value}) \\ a &\mapsto \bigcup_{\sigma \in X} \{\sigma a\} \end{aligned}$$

La sémantique abstraite que nous souhaitons définir est alors la sémantique suivante :  $\llbracket P \rrbracket^{\#} \in (\mathbf{Add} \rightarrow \mathcal{P}(\mathbf{Value}))^{\#\mathbf{Label}}$  dont la  $k^{\text{ième}}$  composante serait l'état abstrait de l'étiquette  $k$ . Elle sera plus simple que la sémantique collectrice puisqu'elle ne comportera qu'un état abstrait dans chaque composante et non un ensemble d'états. On diminue donc le taux d'indirection afin d'obtenir un élément atomique de la sémantique (à savoir un état pour la sémantique collectrice et un état abstrait pour la sémantique abstraite).

Afin de pouvoir atteindre ce but, nous allons munir l'ensemble des sémantiques abstraites  $(\mathbf{Add} \rightarrow \mathcal{P}(\mathbf{Value}))^{\#\mathbf{Label}}$  d'une structure de treillis :  $((\mathbf{Add} \rightarrow \mathcal{P}(\mathbf{Value}))^{\#\mathbf{Label}}, \sqcup, \sqcap, \sqsubseteq)$ . L'objectif est d'obtenir, au final, une insertion de Galois entre  $(\llbracket P \rrbracket^{\text{col}}, \cup, \cap, \subseteq)$  et  $(\llbracket P \rrbracket^{\#}, \sqcup, \sqcap, \sqsubseteq)$ . En effet, cela apportera comme propriété que la solution au système d'équation 1 dans  $(\mathbf{Add} \rightarrow \mathcal{P}(\mathbf{Value}))^{\#\mathbf{Label}}$  est une surapproximation de la solution dans  $\llbracket P \rrbracket^{\text{col}}$ . Donc, si la solution  $\llbracket P \rrbracket^{\#}$  n'atteint pas d'état interdit, la solution  $\llbracket P \rrbracket^{\text{col}}$  non plus.

Pour doter  $(\mathbf{Add} \rightarrow \mathcal{P}(\mathbf{Value}))^{\#\mathbf{Label}}$  d'une structure de treillis, on définit d'abord la borne supérieure ainsi :  $\varphi \sqcup \varphi' = (a \mapsto (\varphi a) \cup (\varphi' a))$ . Autrement dit, la borne supérieure de

deux états abstraits  $\varphi_1$  et  $\varphi_2$  est l'état abstrait qui à une adresse associe toutes les valeurs qui sont associées à  $\varphi_1$  ou à  $\varphi_2$  et seulement celles-là. On définit ensuite la borne inférieure comme ceci :  $\varphi \sqcap \varphi' = (a \mapsto (\varphi a) \cap (\varphi' a))$ . Ce qui signifie que la borne inférieure de deux états abstraits  $\varphi_1$  et  $\varphi_2$  est l'état abstrait qui a une adresse associe toutes les valeurs qui sont associées à  $\varphi_1$  et à  $\varphi_2$  et uniquement celles-là. Enfin, on utilise la relation d'ordre suivante :  $\varphi \sqsubseteq \varphi'$  si et seulement si  $(\forall a, (\varphi a) \subseteq (\varphi' a))$ . Pour  $\sigma \in \mathbf{State}$ , on notera  $\sigma \sqsubseteq \varphi$  lorsque  $(a \mapsto \{\sigma a\}) \sqsubseteq \varphi$  par abus de notation.

Pour montrer qu'il s'agit bien d'une insertion de Galois on définit explicitement la fonction  $\alpha : \llbracket P \rrbracket^{\text{col}} \rightarrow \llbracket P \rrbracket^{\#}$  et l'injection  $\gamma : \llbracket P \rrbracket^{\#} \hookrightarrow \llbracket P \rrbracket^{\text{col}}$ . On pose  $\alpha : X_i \mapsto (a \mapsto \{\sigma a \mid \sigma \in X_i\})$  noté  $\varphi_i$ . On voit ici que l'on agglomère les états, en perdant un peu de précision. En effet, si l'on considère  $\mathbf{Add} = \{1, 2\}$ , on ne distingue plus  $X = \{\sigma_1 : (1 \mapsto 1, 2 \mapsto 2), \sigma_2 : (1 \mapsto 2, 2 \mapsto 1)\}$  de  $Y = \{\sigma_3 : (1 \mapsto 1, 2 \mapsto 1), \sigma_4 : (1 \mapsto 2, 2 \mapsto 2)\}$  car  $\alpha X = \alpha Y = \varphi : (1 \mapsto \{1, 2\}, 2 \mapsto \{1, 2\})$ . On pose, ensuite,  $\gamma : \varphi_i \mapsto \{\sigma \mid \forall a \in \mathbf{Add}, (\sigma a) \in (\varphi_i a)\}$ . Il s'agit de considérer tout les états qui peuvent contribuer à obtenir l'état abstrait via  $\alpha$ . Pour continuer l'exemple précédent :  $\gamma \varphi = \{\sigma_1, \sigma_2, \sigma_3, \sigma_4\}$ . On remarque que  $\gamma$  est bien injective puisque bijective sur son image d'application réciproque  $\alpha$ .

Il faut maintenant montrer que ces fonctions définissent bien une insertion de Galois, donc que  $(\alpha X_i) \sqsubseteq \varphi_j \Leftrightarrow X_i \subseteq (\gamma \varphi_j)$ . Pour cela, nous montrons la double implication :

$$\begin{aligned}
\bullet (\alpha X_i) \sqsubseteq \varphi_j &\Rightarrow X_i \subseteq (\gamma \varphi_j) \\
(\alpha X_i) \sqsubseteq \varphi_j &\Rightarrow \forall a \in \mathbf{Add}, (\alpha X_i a) \subseteq (\varphi_j a) \\
&\Rightarrow \forall \sigma, (\forall a \in \mathbf{Add}, (\sigma a) \in (\alpha X_i a)) \\
&\Rightarrow (\forall a \in \mathbf{Add}, (\sigma a) \in (\varphi_j a)) \\
&\Rightarrow \gamma(\alpha X_i) \subseteq \gamma \varphi_j \\
\text{Or, par définition, } \forall \sigma \in X_i, \forall a \in \mathbf{Add}, (\sigma a) &\in \alpha X_i a \\
\text{donc } X_i &\subseteq \gamma(\alpha X_i) \subseteq \gamma \varphi_j
\end{aligned}$$

□

$$\begin{aligned}
\bullet X_i \subseteq (\gamma \varphi_j) &\Rightarrow (\alpha X_i) \sqsubseteq \varphi_j \\
X_i \subseteq (\gamma \varphi_j) &\Rightarrow (\alpha X_i) \sqsubseteq \alpha(\gamma \varphi_j) = \varphi_j
\end{aligned}$$

□

Ceci montre que l'on a un insertion de Galois entre  $(\llbracket P \rrbracket^{\text{col}}, \cup, \cap, \subseteq)$  et  $(\llbracket P \rrbracket^{\#}, \sqcup, \sqcap, \sqsubseteq)$ . Il est donc pertinent de calculer  $\llbracket P \rrbracket^{\#}$  dans la mesure où l'on sait que si elle ne possède pas d'état abstrait dangereux alors  $\llbracket P \rrbracket^{\text{col}}$  ne possède pas d'état dangereux non plus.

## 4.2 Analyse dans la sémantique abstraite

Afin de calculer  $\llbracket P \rrbracket^{\#}$ , on calcule de l'image du système d'équation 1 par  $\alpha$ . Pour cela, on utilise le fait que  $\alpha(\mathbf{A} \cup \mathbf{B}) = (\alpha \mathbf{A}) \sqcup (\alpha \mathbf{B})$ . o, obtient :

$$\forall k \in \mathbf{Label}, \varphi_k = \varphi_k^{\text{init}} \sqcup \bigsqcup_{(k', i, k) \in P} \llbracket i \rrbracket^{\#} \varphi_{k'} \quad \text{où } \llbracket i \rrbracket^{\#} : \varphi \mapsto (a \mapsto \bigcup_{\substack{\sigma_1 \sqsubseteq \varphi \\ \sigma_1 \xrightarrow{i} \sigma_2}} (\sigma_2 a))$$

De plus, comme on ne peut plus calculer qu'une surapproximation de la sémantique collectrice, seul un sens de l'égalité nous intéresse, à savoir que le membre de droite est supérieur au membre de gauche. En effet, si une sémantique vérifie une telle inégalité, cela signifie que l'on a une surapproximation de la sémantique abstraite et donc de la sémantique collectrice. Cela implique que de la même manière que si la sémantique abstraite n'atteint pas d'état dangereux alors la sémantique collectrice non plus, si une sémantique vérifiant

une telle inégalité n'atteint aucun état dangereux, c'est aussi le cas pour la sémantique collectrice. Au final on cherche à résoudre :

$$\forall k \in \mathbf{Label}, \varphi_k \sqsupseteq \varphi_k^{\text{init}} \sqcup \bigsqcup_{(k',i,k) \in P} \llbracket i \rrbracket^\# \varphi_{k'} \quad \text{où } \llbracket i \rrbracket^\# : \varphi \mapsto (a \mapsto \bigcup_{\substack{\sigma_1 \sqsubseteq \varphi \\ \sigma_1 \xrightarrow{i} \sigma_2}} (\sigma_2 a)) \quad (2)$$

De la même manière que pour la sémantique collectrice, le système d'équations peut être vu comme une fonction croissante dans un treillis fini. On peut donc calculer une solution en itérant la fonction (correspondant au système) jusqu'à atteindre un point fixe.

On détail maintenant à quoi correspond le système d'équations en observant ce que les fonctions  $(\llbracket i \rrbracket^\# \varphi_{k'}) \sqsubseteq \varphi_k$  signifient en fonction de l'instruction  $i$  :

$$\begin{aligned} i = \text{skip} & \quad \varphi_{k'} \sqsubseteq \varphi_k \\ i = \text{assert } b & \quad \forall \sigma \sqsubseteq \varphi_{k'}, \mathcal{E}_b[b] \sigma \rho \Rightarrow \sigma \sqsubseteq \varphi_k \\ i = \epsilon := e & \quad \forall \sigma \sqsubseteq \varphi_{k'}, \begin{cases} \rho (\mathcal{E}_v[\epsilon] \sigma \rho) & \mapsto \{ \mathcal{E}_a[e] \sigma \rho \} \\ a \neq (\mathcal{E}_v[\epsilon] \sigma \rho) & \mapsto \varphi_{k'} a \end{cases} \sqsubseteq \varphi_k \end{aligned}$$

Afin de simplifier ces équations nous avons choisi de faire deux hypothèses : on part du principe que rien n'est initialisé au début et on prend le cas d'une analyse insensible au flot. La première hypothèse signifie que  $\forall k, \forall a \in \mathbf{Add}, \varphi_k^{\text{init}} a = \perp$ , autrement dit  $\varphi_k^{\text{init}}$  est élément neutre pour la borne supérieure dans le treillis, donc peut être enlevé des équations. La seconde hypothèse signifie que  $\mathbf{Label}$  est réduit à un seul élément :  $\mathbf{Label} = \{.\}$ . Les équations deviennent :

$$\varphi. \sqsupseteq \varphi. \cup \bigcup_{(.,i.,.) \in P} \llbracket i \rrbracket^\# \varphi.$$

Des cas selon  $i$  sont devenus triviaux. En effet, pour  $i = \text{skip}$ , l'équation signifie maintenant  $\varphi. \sqsubseteq \varphi.$  ce qui est une tautologie. De plus, pour  $i = \text{assert } b$ , l'équation devient  $\forall \sigma \sqsubseteq \varphi., \mathcal{E}_b[b] \sigma \rho \Rightarrow \sigma \sqsubseteq \varphi.$  Autrement dit, des éléments plus petits que  $\varphi.$  sont plus petits que  $\varphi.$  Ceci est à nouveau une tautologie.

Le seul cas restant est  $i = \epsilon := e$ . On s'intéresse aux cas selon  $\epsilon$  et  $e$  ( $x, y \in \mathbf{Var}$ ) :

$x := y$	$\varphi. (\rho y) \sqsubseteq \varphi. (\rho x)$
$x := *y$	$\forall v \in \mathbf{Add}, v \in \varphi. (\rho y) \Rightarrow \varphi. v \sqsubseteq \varphi. (\rho x)$
$x := \&y$	$\rho y \in \varphi. (\rho x)$
$*x := y$	$\forall v \in \mathbf{Add}, v \in \varphi. (\rho x) \Rightarrow \varphi. (\rho y) \sqsubseteq \varphi. v$
$*x := \&y$	$\forall v \in \mathbf{Add}, v \in \varphi. (\rho x) \Rightarrow \rho y \in \varphi. v$
$*x := *y$	$\forall v_1, v_2 \in \mathbf{Add}, (v_1 \in \varphi. (\rho x) \wedge v_2 \in \varphi. (\rho y)) \Rightarrow \varphi. v_2 \sqsubseteq \varphi. v_1$

On remarque que tous les quantificateurs universels s'appliquent à des ensembles finis (e.g. du type  $\varphi. (\rho x)$  où  $x \in \mathbf{Var}$ ). Notons  $a_1, \dots, a_A$  les éléments de  $\mathbf{Add}$  (où  $A = \#\mathbf{Add}$  le nombre d'éléments de  $\mathbf{Add}$ ). Notons  $v_{i,1}, \dots, v_{i,V_i}$  les éléments de  $\varphi. a_i$  (où  $V_i = \#(\varphi. a_i)$ ). On peut voir les équations obtenues dans la deuxième colonne de la table 1. Enfin, on peut se ramener à un test d'inclusion pour chaque cas où réalisant des unions et intersections idoines comme le montre la troisième colonne de la table 1. Ce sont ces équations que nous souhaitons vérifier à l'aide d'un algorithme probabiliste.

## 5 Algorithme de vérification de résultats de l'analyse

Nous pouvons définir la problématique précisément : on est capable de générer les équations correspondant à un programme (table 1), l'environnement  $\rho$  est connu, on nous

TAB. 1 – simplifications successives des équations à vérifier.

$x := y$	$\varphi.(\rho y) \subseteq \varphi.(\rho x)$	$\varphi.(\rho y) \subseteq \varphi.(\rho x)$
$x := *y$	Soit $a_i = \rho y$ : $\varphi.v_{i,1} \subseteq \varphi.(\rho x)$ $\vdots$ $\varphi.v_{i,V_i} \subseteq \varphi.(\rho x)$	Soit $a_i = \rho y$ : $\bigcup_{k=1}^{V_i} \varphi.v_{i,k} \subseteq \varphi.(\rho x)$
$x := \&y$	$\rho y \in \varphi.(\rho x)$	$\rho y \in \varphi.(\rho x)$
$*x := y$	Soit $a_i = \rho x$ : $\varphi.(\rho y) \subseteq \varphi.v_{i,1}$ $\vdots$ $\varphi.(\rho y) \subseteq \varphi.v_{i,V_i}$	Soit $a_i = \rho x$ : $\varphi.(\rho y) \subseteq \bigcap_{k=1}^{V_i} \varphi.v_{i,k}$
$*x := \&y$	Soit $a_i = \rho x$ : $\rho y \subseteq \varphi.v_{i,1}$ $\vdots$ $\rho y \subseteq \varphi.v_{i,V_i}$	Soit $a_i = \rho x$ : $\rho y \subseteq \bigcap_{k=1}^{V_i} \varphi.v_{i,k}$
$*x := *y$	Soient $a_i = \rho x, a_j = \rho y$ : $\varphi.v_{j,1} \subseteq \varphi.v_{i,1} \cdots \varphi.v_{j,V_j} \subseteq \varphi.v_{i,1}$ $\vdots \quad \ddots \quad \vdots$ $\varphi.v_{j,1} \subseteq \varphi.v_{i,V_i} \cdots \varphi.v_{j,V_j} \subseteq \varphi.v_{i,V_i}$	Soient $a_i = \rho x, a_j = \rho y$ : $\bigcup_{k=1}^{V_j} \varphi.v_{j,k} \subseteq \bigcap_{k=1}^{V_i} \varphi.v_{i,k}$

fournit un état abstrait  $\varphi$ . (le certificat), et on souhaite vérifier que l'état abstrait vérifie les équations. Nous concevons pour le résoudre un algorithme de Monte-Carlo. Pour cela, nous commençons par discuter des choix de conception qui peuvent être faits (section 5.1). Nous calculons ensuite une valeur théorique à la probabilité d'erreur de l'algorithme (sections 5.2 et 5.3). Enfin nous montrons qu'une hypothèse faite lors de du calcul de la valeur théorique n'est pas nécessaire (section 5.4).

## 5.1 Conception de l'algorithme

Nous souhaitons réaliser un algorithme probabiliste. Il ne faut donc pas perdre de vue que l'on cherche à réaliser un algorithme simple et efficace à implémenter. Étant données, les équations de la table 1, on voit que l'on peut réaliser au moins trois actions de façon probabiliste. Premièrement, on peut vérifier partiellement si une inclusion donnée est vérifiée, ce sera traité section 5.2. Deuxièmement, on peut ne vérifier qu'une partie des inclusions, on s'intéresse à cela en section 5.3. Troisièmement, on peut calculer approximativement les réunions et intersections, nous ne traiterons pas de cela ici. Enfin, afin de garder un algorithme le plus simple possible et de limiter l'espace nécessaire, nous décidons de ne pas garder en mémoire ce qui a déjà été vérifié. L'algorithme 1 que nous obtenons a deux paramètres le nombre d'inclusions testées  $N$  et le nombre d'appartenances testées par inclusion  $L$ .

Si l'algorithme trouve une appartenance fautive, il s'arrête et rejette le certificat. S'il termine sans trouver d'appartenance fautive, il accepte le certificat. On remarque que d'après nos choix de conception, l'algorithme peut tester plusieurs fois la même inclusion ou la même appartenance. On remarque de plus que notre algorithme de Monte-Carlo ne peut pas se tromper quand il annonce qu'un certificat est faux. En effet, cela signifie que l'algo-

**Input** : Ensemble d'équations (du type table 1), certificat à vérifier.

**Output** : Oui si le certificat vérifie les équations, non sinon.

```

1 Réaliser les unions et intersections dans les équations (cf. table 1)
2 for n ← 1 to N do
3   Choisir aléatoirement une inclusion à vérifier  $\mathbf{A}_n \subseteq \mathbf{B}_n$ 
4   for l ← 1 to N do
5     Choisir aléatoirement un élément  $a_l \in \mathbf{A}_n$ 
6     if  $a_l \notin \mathbf{B}_n$  then
7       terminer l'algorithme, répondre non
8 terminer l'algorithme, répondre oui

```

**Algorithme 1**: Notre algorithme.

rithme a trouvé une inclusion fausse.

Nous nous intéressons à quatre propriétés sur notre algorithme : sa probabilité d'erreur, le nombre de bit aléatoires qu'il génère (car cela a un coût), sa complexité spatiale et sa complexité temporelle. La probabilité d'erreur théorique va occuper toute la suite de cette section (5.2, 5.3, 5.4) dans laquelle on obtient une formule théorique dans le cas où les ensembles sont générés aléatoirement. Le nombre de bit aléatoires générés, la complexité spatiale et la complexité temporelles dépendent de l'implémentation, ils seront donc évalués dans la section 6. On note cependant qu'un algorithme testant  $N$  inclusions et  $L$  appartenances par inclusion réalise en tout  $N \times L$  tests d'appartenance. L'évaluation de la complexité temporelle consiste donc principalement à estimer le temps nécessaire à la réalisation d'un test d'appartenance.

À ce propos, on remarque que par rapport à un algorithme testant  $N$  inclusions et  $L$  appartenances par inclusion, l'algorithme testant  $N \times L$  inclusions et 1 appartenance par inclusion réalise autant de tests d'appartenance mais fait plus jouer l'aléa car il est moins restrictif sur la répartition des appartenances testées. Intuitivement, on peut donc penser qu'il est plus intéressant de maximiser le premier paramètre et de garder le second à 1. Cette intuition sera confirmée par les résultats théoriques et pratiques.

Puisque l'on va calculer des probabilités, on définit ici l'espace des événements  $\Omega$ . On rappelle que les ensembles  $\mathbf{A}$  et  $\mathbf{B}$  sont des sous-ensembles de **Value**. Soient,  $v_1, \dots, v_{\#\mathbf{Value}}$  les éléments de **Value**, on pose  $\Omega_1 = \{v_k \in \mathbf{E} \mid k \in \llbracket 1, \#\mathbf{Value} \rrbracket \text{ et } \mathbf{E} \in \mathcal{P}(\mathbf{Value})\}$ . De plus, soient  $e_1, \dots, e_{\#\mathbf{E}}$  les éléments d'un ensemble  $\mathbf{E}$ ; on pose  $\Omega_{\mathbf{E}} = \{\text{"on choisit } e_k \text{ au } l\text{-ième tirage"} \mid k \in \llbracket 1, \#\mathbf{E} \rrbracket \text{ et } l \in \llbracket 1, L \rrbracket\}$ . Enfin, soient  $I_1, \dots, I_N$  les inclusions à vérifier, on pose  $\Omega_2 = \{\text{"on teste l'inclusion } I_k \text{ à l'étape } n \mid k \in \llbracket 1, N \rrbracket \text{ et } n \in \llbracket 1, N \rrbracket\}$ . Au final, on a :

$$\Omega = \Omega_1 \times \Omega_2 \times \prod_{\mathbf{E} \in \mathcal{P}(\mathbf{Value})} \Omega_{\mathbf{E}}$$

On considère volontairement un espace d'évènements plus grand que nécessaire<sup>5</sup> pour simplifier la définition.

Dans les sections suivantes, étant donné que nous cherchons à établir une probabilité, nous essayons à chaque fois de nous ramener à des problèmes connus, notamment de tirages.

---

<sup>5</sup>Notamment en ce qui concerne les  $\Omega_{\mathbf{E}}$ .

FIG. 5 – Inclusion de **A** dans **B**

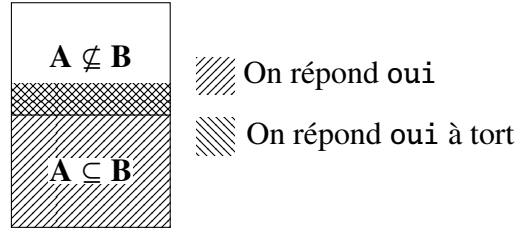
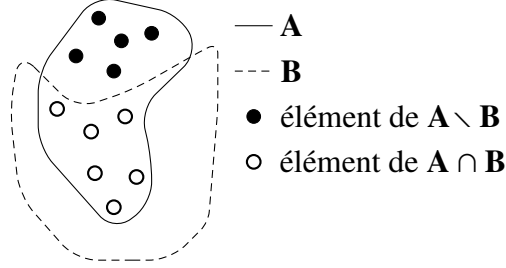


FIG. 6 – **A** est une urne, ses éléments des boules.



## 5.2 Vérification probabiliste d'une inclusion

Nous commençons par établir la probabilité d'accepter une inclusion à tort en la vérifiant partiellement.

Plus précisément, soient **A** et **B** deux ensembles. On souhaite vérifier que  $\mathbf{A} \subseteq \mathbf{B}$ . Pour cela, on choisit  $L$  fois aléatoirement un élément de **A** et on teste s'il appartient à **B**. Si, lors d'un essai où l'on prend un élément  $a \in \mathbf{A}$ , on a  $a \notin \mathbf{B}$  alors on répond non. Sinon — c'est-à-dire si à chaque essai, on prend un élément  $a \in \mathbf{A}$  tel que  $a \in \mathbf{B}$  — on répond oui.

On remarque que si  $\mathbf{A} \subseteq \mathbf{B}$  alors, on répond systématiquement oui. Autrement dit, on ne peut pas répondre non à tort. On cherche donc à évaluer la probabilité de répondre oui à tort. Autrement la probabilité de l'intersection de deux événements : "on répond oui" et " $\mathbf{A} \not\subseteq \mathbf{B}$ " (voir figure 5). On calcule donc :

$$\begin{aligned} \mathcal{P}_r(\text{"oui à tort"}) &= \mathcal{P}_r(\text{"oui"} \cap \mathbf{A} \not\subseteq \mathbf{B}) && \text{(figure 5)} \\ &= \sum_{i=1}^{\#\mathbf{A}} \mathcal{P}_r(\text{"oui"} \cap \#\mathbf{A} \setminus \mathbf{B} = i) && \text{(probabilités totales)} \\ &= \sum_{i=1}^{\#\mathbf{A}} \mathcal{P}_r(\text{"oui"} | \#\mathbf{A} \setminus \mathbf{B} = i) \mathcal{P}_r((\#\mathbf{A} \setminus \mathbf{B}) = i) && \text{(probabilité conditionnelle)} \end{aligned}$$

On va maintenant calculer une valeur explicite pour chaque membre du produit dans la somme :  $\mathcal{P}_r(\text{"oui"} | (\#\mathbf{A} \setminus \mathbf{B}) = i)$  et  $\mathcal{P}_r((\#\mathbf{A} \setminus \mathbf{B}) = i)$ .

Pour le premier membre, nous montrons que le fait de tester  $L$  fois si un élément de **A** appartient à **B** sachant que  $i$  éléments de **A** n'appartiennent pas à **B** est équivalent à un tirage avec remise. Dans la figure 6, **A** est vu comme une urne, ses éléments sont vus comme des boules. La boule est noire si l'élément n'est pas dans **B** et blanche l'élément est dans **B**. De cette façon, on voit que tirer  $L$  fois un élément de **A** qui est aussi un élément de **B**, revient à tirer  $L$  fois une boule blanche dans une urne contenant  $\#\mathbf{A}$  boules dont  $\#\mathbf{A} \setminus \mathbf{B}$  boules noires et  $\#\mathbf{A} \cap \mathbf{B}$  boules blanches. Ce type de tirage suit la loi de probabilité binomiale :

$$\mathcal{P}_r(\text{"oui"} | \#\mathbf{A} \setminus \mathbf{B} = i) = \left( \frac{\#\mathbf{A} - i}{\#\mathbf{A}} \right)^L$$

Pour le second membre, on peut voir sur la figure 7 que pour déterminer  $\mathcal{P}_r(\#\mathbf{A} \setminus \mathbf{B} = i)$ , on peut se ramener à un tirage sans remise où **Value** est l'urne,  $\#\mathbf{B}$  le nombre de boules

FIG. 7 – **Add** est une urne, **B** l'ensemble des boules blanches, **A** le tirage mis dans une urne pour réaliser le tirage de la figure 6 .

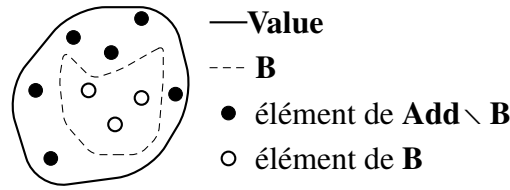
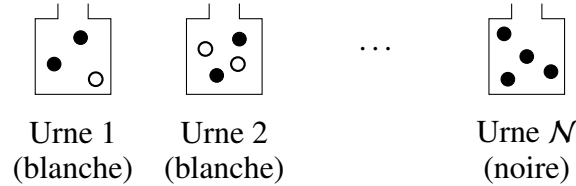


FIG. 8 – On choisit une urne aléatoirement, puis on effectue un tirage avec remise dessus.



blanches,  $\#A$  le nombre de boules tirées, et  $i$  le nombre de boules noires tirées. Ce type de tirage suit la loi de probabilité hypergéométrique :

$$\mathcal{P}_r(\#A \setminus B = i) = \frac{C_{\#Add-\#B}^i C_{\#B}^{\#A-i}}{C_{\#Add}^{\#A}}$$

Au final, on a :

$$\mathcal{P}_r(\text{“oui à tort”}) = \sum_{i=1}^{\#A} \left( \frac{\#A - i}{\#A} \right)^L \frac{C_{\#Value-\#B}^i C_{\#B}^{\#A-i}}{C_{\#Add}^{\#A}}$$

Nous sommes donc en mesure d'estimer la probabilité de croire qu'une inclusion est vraie à tort. Nous allons maintenant calculer la probabilité de croire que l'ensemble d'un certificat est vrai à tort.

### 5.3 Vérification de l'ensemble du résultat

On souhaite tester  $N$  fois une inclusion au hasard parmi les  $N$  inclusions proposées. On voit sur la figure 8 que cela revient à choisir au hasard  $N$  fois une urne parmi  $N$ . Les urnes contiennent des boules noires ou blanches, et sont numérotés. On ne connaît pas le nombre d'urnes ne contenant que des boules noires (correspondant aux inclusions vraies) et que l'on appelle « urnes noires ». On connaît pas non plus le nombre d'urnes contenant au moins une boule blanche (correspondant aux inclusions fausses) et que l'on appelle « urnes blanches ». En revanche, on sait calculer la probabilité d'une urne  $i$  d'être noire. De plus, on connaît la probabilité de ne tirer que des boules noires d'une urne blanche (cela a été calculé en section 5.2).

On prend comme hypothèse que les  $N$  tirages sont indépendants. Nous verrons dans la section 5.4 que cette hypothèse est fautive en toute généralité mais permet de simplifier la formule. De plus nous verrons dans la section 6 que les résultats expérimentaux sont proches des résultats théoriques faisant cette hypothèse.

Afin de calculer la probabilité qu'à l'algorithme de se tromper on calcule d'abord la probabilité de juger l'inclusion tirée vraie. L'hypothèse d'indépendance permet alors de déduire rapidement la probabilité de juger toutes les inclusions tirées vraies. En terme de tirage, cela signifie que l'on va d'abord calculer la probabilité de considérer une urne choisie noire, puis en déduire la probabilité de croire toutes les urnes noires.



On commence donc en calculant la probabilité  $P$  qu'une urnes tirée soit vue noire :

$$\begin{aligned}
P &= \mathcal{P}_r(\text{“urne vue noire”}) \\
&= \sum_{k=1}^N \mathcal{P}_r(\text{“urne crue noire”} | \text{“urne n}^\circ k \text{ tirée”}) \mathcal{P}_r(\text{“urne n}^\circ k \text{ tirée”}) \\
&\quad \text{d'après la formule de probabilités totales} \\
&= \sum_{k=1}^N \mathcal{P}_r(\text{“urne noire”} \cup \text{“urne blanche crue noire”} | \text{“urne n}^\circ k \text{ tirée”}) \frac{1}{N} \\
&\quad \text{car } \mathcal{P}_r(\text{“urne n}^\circ k \text{ tirée”}) = \frac{1}{N} \\
&= \frac{1}{N} \sum_{k=1}^N \mathcal{P}_r(\text{“urne noire”} | \text{“urne n}^\circ k \text{”}) + \mathcal{P}_r(\text{“urne blanche crue noire”} | \text{“urne n}^\circ k \text{”}) \\
&\quad \text{car les deux évènements sont disjoints}
\end{aligned}$$

Or, « urne noire » signifie que l'inclusion est vraie donc si l'urne  $k$  correspond à l'équation  $\mathbf{A}_k \subseteq \mathbf{B}_k$ , on a  $\mathcal{P}_r(\text{“urne noire”} | \text{“urne n}^\circ k \text{”}) = \mathcal{P}_r(\mathbf{A}_k \setminus \mathbf{B}_k = 0)$ . C'est un cas particulier du problème de la figure 6 pour  $i = 0$ . On a donc  $\mathcal{P}_r(\text{“urne noire”} | \text{“urne n}^\circ k \text{”}) = \frac{C_{\#\mathbf{B}_k}^{\#\mathbf{A}_k}}{C_{\#\mathbf{Add}}^{\#\mathbf{A}_k}}$ . De plus, « croire une urne blanche noire » correspond à accepter une inclusion à tort, cela a déjà été calculé dans la section 6, on a donc :

$$\mathcal{P}_r(\text{“urne blanche crue noire”} | \text{“urne n}^\circ k \text{”}) = \sum_{i=1}^{\#\mathbf{A}_k} \left( \frac{\#\mathbf{A}_k - i}{\#\mathbf{A}_k} \right)^L \frac{C_{\#\mathbf{Add} - \#\mathbf{B}_k}^i C_{\#\mathbf{B}_k}^{\#\mathbf{A}_k - i}}{C_{\#\mathbf{Add}}^{\#\mathbf{A}_k}}$$

Enfin comme  $\frac{C_{\#\mathbf{B}_k}^{\#\mathbf{A}_k}}{C_{\#\mathbf{Add}}^{\#\mathbf{A}_k}}$  est un cas particulier de la somme pour  $i = 0$ , on peut regrouper les deux cas. On a au final :

$$P = \mathcal{P}_r(\text{“urne crue noire”}) = \frac{1}{N} \sum_{k=1}^N \sum_{i=0}^{\#\mathbf{A}_k} \left( \frac{\#\mathbf{A}_k - i}{\#\mathbf{A}_k} \right)^L \frac{C_{\#\mathbf{Add} - \#\mathbf{B}_k}^i C_{\#\mathbf{B}_k}^{\#\mathbf{A}_k - i}}{C_{\#\mathbf{Add}}^{\#\mathbf{A}_k}}$$

Sous l'hypothèse que les tirages sont indépendants, la probabilité d'annoncer toutes les inclusions vraies est le produit des probabilités de croire l'urne choisie noire  $N$  fois. Autrement dit :

$$\mathcal{P}_r(\text{“annoncer inclusions vraies”}) = \left( \frac{1}{N} \sum_{k=1}^N \sum_{i=0}^{\#\mathbf{A}_k} \left( \frac{\#\mathbf{A}_k - i}{\#\mathbf{A}_k} \right)^L \frac{C_{\#\mathbf{Add} - \#\mathbf{B}_k}^i C_{\#\mathbf{B}_k}^{\#\mathbf{A}_k - i}}{C_{\#\mathbf{Add}}^{\#\mathbf{A}_k}} \right)^N$$

On souhaite connaître la probabilité d'annoncer toutes inclusions vraies à tort (puisqu'annoncer toutes inclusions vraies à raison est un comportement souhaité). Or d'après la loi des probabilités totales :

$$\mathcal{P}_r(\text{“annoncer toutes inclusions vraies”}) = \mathcal{P}_r(\text{“annoncer inclusions vraies à tort”}) + \mathcal{P}_r(\text{“annoncer inclusions vraies à raison”})$$

d'où

$$\mathcal{P}_r(\text{“annoncer inclusions vraies à tort”}) = \mathcal{P}_r(\text{“annoncer inclusions vraies”}) - \mathcal{P}_r(\text{“annoncer inclusions vraies à raison”})$$

La probabilité d'annoncer toutes les inclusions vraies a été calculée. La probabilité d'annoncer toutes les inclusions vraies à raison est simplement la probabilité que toutes les inclusions soient effectivement vraies. En effet, lorsque toutes les inclusions sont vraies,

l'algorithme annonce nécessairement qu'elles le sont. On a donc :

$$\begin{aligned}\mathcal{P}_r(\text{"annoncer à raison"}) &= \mathcal{P}_r(\text{"toutes les inclusions sont vraies"}) \\ &= \prod_{k=1}^N \mathcal{P}_r(\text{"urne noire"} | \text{"urne n° k"}) \\ &= \prod_{k=1}^N \frac{C_{\#B_k}^{\#A_k}}{C_{\#Add}^{\#A_k}}\end{aligned}$$

Au final on obtient la valeur théorique suivante :

$$\mathcal{P}_r(\text{"annoncer vrai à tort"}) = \left( \frac{1}{N} \sum_{k=1}^N \sum_{i=0}^{\#A_k} \left( \frac{\#A_k - i}{\#A_k} \right)^L \frac{C_{\#Add - \#B_k}^i C_{\#B_k}^{\#A_k - i}}{C_{\#Add}^{\#A_k}} \right)^N - \prod_{k=1}^N \frac{C_{\#B_k}^{\#A_k}}{C_{\#Add}^{\#A_k}} \quad (3)$$

On constate que comme on l'a expliqué intuitivement, le nombre d'inclusions testées  $N$  agit de manière plus globale que le nombre d'appartenances testées par inclusion  $L$ . Il est donc plus intéressant de prendre  $L = 1$  et de maximiser  $N$ .

## 5.4 Indépendance des inclusions

L'un des reproches qui peut être fait au résultat précédent est qu'il est obtenu sous l'hypothèse que les inclusions sont indépendantes. Cette hypothèse est fautive en général. En effet, soit les inclusions à tester sont  $\mathbf{A} \subseteq \mathbf{B}$ ,  $\mathbf{B} \subseteq \mathbf{C}$  et  $\mathbf{A} \subseteq \mathbf{C}$ . Clairement,  $\mathcal{P}_r(\mathbf{A} \subseteq \mathbf{C} | \mathbf{A} \subseteq \mathbf{B} \ \& \ \mathbf{B} \subseteq \mathbf{C}) = 1$ .

Nous allons donc essayer de calculer la probabilité qu'une inclusion soit vraie sachant qu'une autre l'est. Formellement, soient  $\mathbf{A}_1, \mathbf{B}_1, \mathbf{A}_2, \mathbf{B}_2$  des ensembles. On souhaite calculer  $\mathcal{P}_r(\mathbf{A}_2 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1)$ . Pour cela, on commence par séparer le problème en deux problèmes plus petits en partitionnant l'événement  $\mathbf{A}_2 \subseteq \mathbf{B}_2$  :

$$\mathcal{P}_r(\mathbf{A}_2 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1) = \mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1)$$

Or, par application des probabilités conditionnelles, on peut séparer le second membre en un produit de deux facteurs :

$$\mathcal{P}_r(\mathbf{A}_2 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1) = \mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2) \mathcal{P}_r(\mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1) \quad (4)$$

On va maintenant simplifier le premier facteur de l'équation 4, à savoir :  $\mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2)$ . On remarque d'abord que les événements  $\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2$  et  $\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1$  sont équivalents. On a donc :

$$\mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2) = \mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1 | \mathbf{A}_1 \subseteq \mathbf{B}_1 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2)$$

Or, on pose l'hypothèse que les événements  $\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1$  et  $\mathbf{A}_1 \subseteq \mathbf{B}_1$  sont indépendants. Cela vient du fait que  $(\mathbf{A}_2 \setminus \mathbf{A}_1) \cap \mathbf{A}_1 = (\mathbf{B}_2 \setminus \mathbf{A}_1) \cap \mathbf{A}_1 = \emptyset$ . Autrement dit,  $\mathbf{A}_1$  n'a rien à voir avec ces ensembles ; avoir une propriété sur  $\mathbf{A}_1$  ne nous donne donc pas d'information pertinente. De même, les événements  $\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1$  et  $\mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2$  sont supposés indépendants. Finalement, on a simplement :

$$\mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1 \ \& \ \mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2) = \mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1)$$

Après simplification, l'équation 4 devient :

$$\mathcal{P}_r(\mathbf{A}_2 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1) = \mathcal{P}_r(\mathbf{A}_2 \setminus \mathbf{A}_1 \subseteq \mathbf{B}_2 \setminus \mathbf{A}_1) \mathcal{P}_r(\mathbf{A}_2 \cap \mathbf{A}_1 \subseteq \mathbf{B}_2 | \mathbf{A}_1 \subseteq \mathbf{B}_1) \quad (5)$$

Afin de pouvoir continuer les calculs, on observe que chacun des tests d'inclusion équivaut à la synthèse d'un ou plusieurs tests d'appartenance. La propriété «  $\mathbf{A} \subseteq \mathbf{B}$  » est équivalente à la propriété « soient  $a_1, \dots, a_m$  les éléments de  $\mathbf{A}$ , on a  $a_1 \in \mathbf{B}, \dots, a_m \in \mathbf{B}$  ». Considérons donc les  $\mathcal{M}$  propriétés d'appartenance correspondant aux  $\mathcal{N}$  propriétés d'inclusion<sup>6</sup>. La probabilité conditionnelle ci-dessus s'applique : une propriété d'appartenance n'étant qu'une propriété d'inclusion dont l'ensemble inclus est réduit à un élément. Dans le cas particulier des appartenances, elle devient ceci :

$$\mathcal{P}_r(a_2 \in \mathbf{B}_2 | a_1 \in \mathbf{B}_1) = \mathcal{P}_r(\{a_2\} \setminus \{a_1\} \subseteq \mathbf{B}_2 \setminus \{a_1\}) \mathcal{P}_r(\{a_2\} \cap \{a_1\} \subseteq \mathbf{B}_2 | a_1 \in \mathbf{B}_1)$$

On remarque que l'un des deux facteurs est nécessairement trivial. En effet, soit  $\{a_2\} \setminus \{a_1\} = \emptyset$ , soit  $\{a_2\} \cap \{a_1\} = \emptyset$ .

On va maintenant procéder à une disjonction de cas :

- Si  $\mathbf{B}_1 = \mathbf{B}_2$ , deux cas se présentent :
- Si  $a_1 = a_2$  alors le problème est trivial :

$$\mathcal{P}_r(a_2 \in \mathbf{B}_2 | a_1 \in \mathbf{B}_1) = \mathcal{P}_r(a_1 \in \mathbf{B}_1 | a_1 \in \mathbf{B}_1) = 1$$

- Si  $a_1 \neq a_2$  :

$$\mathcal{P}_r(a_2 \in \mathbf{B}_2 | a_1 \in \mathbf{B}_1) = \mathcal{P}_r(\{a_2\} \subseteq \mathbf{B}_2 \setminus \{a_1\}) \mathcal{P}_r(\emptyset \subseteq \mathbf{B}_2 | a_1 \in \mathbf{B}_1) = \mathcal{P}_r(a_2 \in \mathbf{B}_2 \setminus \{a_1\})$$

- si  $\mathbf{B}_1 \neq \mathbf{B}_2$ , comme aucune hypothèse n'est faite entre  $\mathbf{B}_1$  et  $\mathbf{B}_2$ , on suppose qu'il y a indépendance :

$$\mathcal{P}_r(a_2 \in \mathbf{B}_2 | a_1 \in \mathbf{B}_1) = \mathcal{P}_r(a_2 \in \mathbf{B}_2)$$

Au final, nous pouvons prendre en compte la dépendance entre les différentes propriétés. Il est possible de calculer une valeur théorique semblable affinant le résultat de l'équation 3 en prenant en compte les dépendances. Pour cela, tout comme dans cette section, il faut séparer les propriétés d'inclusion en propriétés « atomiques » d'appartenance. Il faut alors estimer la probabilité de tester chaque sous ensemble  $\mathbf{E}$  d'appartenance (tirage). On peut ensuite établir la probabilité que l'ensemble des appartenances est vérifié par le certificat sachant que les appartenances de  $\mathbf{E}$  sont vérifiées (grâce aux résultats de cette section). Enfin, il suffit de regrouper.

Nous ne le faisons pas ici car l'équation 3 possède une erreur faible en pratique. Néanmoins, il nous a paru important de s'assurer que l'on peut alléger les hypothèses si cela est nécessaire pour de futurs travaux.

## 6 Implémentation et résultats expérimentaux

Dans la section 5, on a défini un algorithme de Monte-Carlo de vérification de certificats. Nous avons réalisé une implémentation de cet algorithme. Dans la section 6.1, nous discutons des choix d'implémentation et notamment ce que cela implique en termes de complexité temporelle et spatiale. Dans la section 6.2 nous donnons quelques résultats expérimentaux de notre implémentation.

### 6.1 Implémentation

Nous avons implémenté notre algorithme en caml. Dans une première approche nous avons choisi de représenter les valeurs (**Value**) par des entiers et les ensembles par des tableaux triés. De cette façon, nous obtenons une complexité spatiale dans le cas le pire

---

<sup>6</sup>Numéros d'indice choisis arbitrairement.

en  $O(\mathcal{KL})$  où  $\mathcal{K}$  est le nombre d'ensembles et  $\mathcal{L}$  la taille maximale des ensembles. La complexité temporelle dans le cas le pire est elle en  $O(NL(\log_2 \mathcal{L}))$  puisque nous réalisons une recherche dichotomique. Enfin le nombre de bits aléatoires générés est de  $\log_2 \mathcal{K}$  pour choisir aléatoirement un ensemble et au pire  $\log_2 \mathcal{L}$  pour choisir aléatoirement un élément dans un ensemble. Le nombre de bits aléatoires générés est donc de  $N \log_2 \mathcal{K} + NL \log_2 \mathcal{L}$  dans le cas le pire. On remarque que si  $L = 1$  on obtient  $N(\log_2 \mathcal{KL})$  bit générés.

Pour résumer, comme  $\mathcal{K}$  correspond aux inclusions qui correspondent aux instructions et  $\mathcal{L}$  correspond aux valeurs qui correspondent aux adresses donc aux variables. On a une complexité spatiale linéaire en le nombre d'instructions et en le nombre de variables. La complexité temporelle est logarithmique en le nombre de variables et linéaire en chacun des paramètres de l'algorithme. Enfin, on génère un nombre de bits logarithmique en le nombre d'inclusion et de variables (réunis) et linéaire en chacun des paramètres  $L$  et  $N$  de l'algorithme.

## 6.2 Résultats expérimentaux

Lorsque l'algorithme a été implémenté, nous avons réalisé des tests en deux temps. D'abord, nous avons testé l'algorithme avec des certificats générés aléatoirement que nous avons comparés avec nos valeurs théoriques (équation 3). Ensuite, nous avons essayé d'affiner le comportement du producteur. En effet, des certificats générés aléatoirement, correspondent à un producteur qui ne réaliserait pas l'analyse de son logiciel. Un producteur pourrait, au contraire réaliser l'analyse, s'apercevoir qu'elle ne garantit pas que son logiciel se comporte bien et donc modifier le résultat de l'analyse pour camoufler cela. Quel que soit le cas considéré les équations ont toujours été générées aléatoirement et n'ont jamais correspondu à un vrai programme.

Dans un premier temps, pour des certificats réalisés aléatoirement, nous avons réalisé des tests pour des tailles de l'ensemble des valeurs (**Add = Value**) allant de 1 à 11 (valeurs qui peuvent correspondre à une fonction d'un programme), un nombre d'inclusions vérifiées allant de 1 à 11 et un nombre d'appartenances vérifiées par inclusion allant de 1 à 11. À chaque triplet de valeur, nous avons lancé 1 000 000 de fois l'algorithme et compté le nombre de fois qu'il obtenait un résultat erroné.

De ces tests ressortent trois résultats principaux. Le premier est que la probabilité d'erreur est bien exponentiellement décroissante (l'échelle des ordonnées de la figure 9 est logarithmique).

Le deuxième résultat (figure 9a) est que comme prévu, le nombre de test d'appartenances vérifiées par inclusion  $L$  a une influence moins grande que le nombre d'inclusions vérifiées  $N$ . En effet, la figure 9a compare les courbes pour  $L = 1$  et  $L = 11$ . On voit alors, par exemple que pour  $L = 1$  et  $N = 9$  (9 tests d'appartenance au total), le taux d'erreur est de l'ordre de 1 pour 1 000 000, alors que pour  $L = 11$  et  $N = 3$  (33 tests d'appartenance au total), le taux d'erreur est de l'ordre de 1 pour 1 000.

Le troisième résultat (figure 9b) est que bien que faisant une hypothèse erronée en général, la formule théorique fournit une valeur proche des résultats expérimentaux.

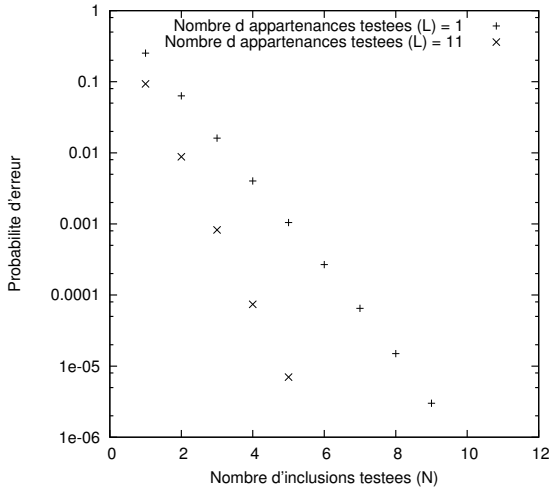
Dans un deuxième temps, pour un vrai certificat légèrement modifié, nous avons établi un certain nombre de manières dont la modification peut être faite. D'abord, des valeurs arbitraires peuvent être retirées d'ensembles arbitraires, ensuite certaines valeurs peuvent être tabou et retirées de tous les ensembles auxquels elles appartiennent.

Dans le premier cas, (figure 10) nous avons réalisé des tests pour avec un ensembles des valeurs (**Add = Value**) de taille  $100^7$ . Nous avons réalisé des tests pour un nombre d'erreurs allant de 1 à 101 avec un pas de 10 et un nombre d'inclusions vérifiées ( $N$ ) allant

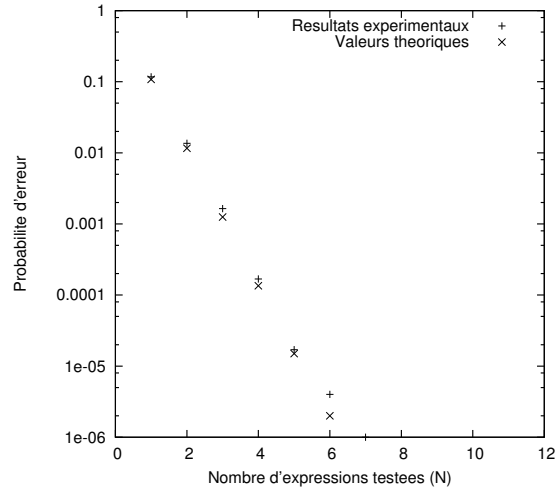
<sup>7</sup>Ceci car si l'ensemble des valeurs est trop petit, 1 erreur est déjà important

FIG. 9 – Résultats expérimentaux pour des certificats générés aléatoirement.

(a) Comparaison de l'influence du nombre d'inclusions testées et du nombre d'appartenances testées par inclusions ( $\#Add = 11$ ).



(b) Comparaison des résultats expérimentaux et des résultats théoriques ( $\#Add = 11$ ). Nombre d'appartenances testées ( $L = 5$ ).



de 1 à 101 avec un pas de 10. Le nombre d'appartenances vérifiées par inclusion ( $L$ ) était fixé à 1. Nous avons lancé 1000 fois l'algorithme pour chaque paire de paramètre.

Nous avons observé que pour une seule erreur la convergence est lente (figure 10). En effet avec 101 tests d'inclusion, le taux d'acceptation erronée est de l'ordre de 0.5. Cependant, il ne faut pas perdre de vue que le certificat est bien plus grand que l'ensemble des valeurs. En réalisant 101 tests d'inclusion, on réalise donc un nombre de tests très inférieur au nombre d'appartenance total à vérifier. On n'a pas fait plus de tests car les temps de calcul devenaient trop long pour lancer assez souvent l'algorithme. Afin de résoudre ce problème nous proposons dans les ouvertures (section 7) une implémentation ayant une meilleure complexité temporelle.

## 7 Ouvertures

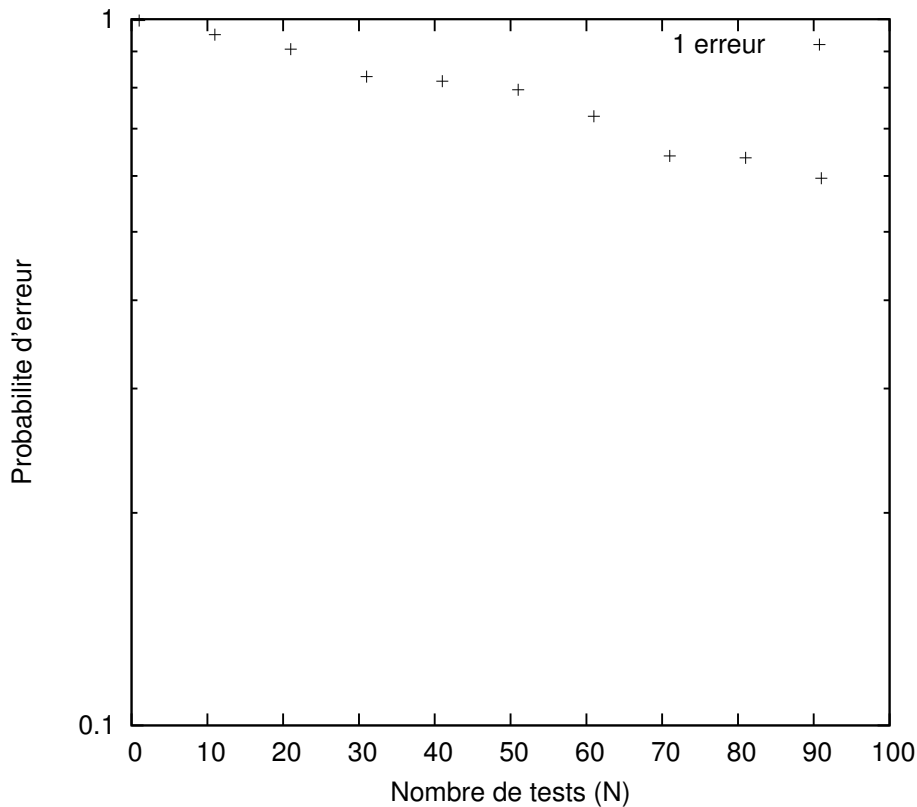
Notre travail ayant été présenté, nous donnons des pistes qui peuvent être suivies pour le généraliser ou l'améliorer. La section 7.1, traite des martingales, un outil de théorie de probabilité puissant et qui pourrait permettre de donner un cadre théorique général à la vérification probabiliste de certificat (plutôt que de traiter l'exemple d'une analyse spécifique). La section 7.2 propose, quant à elle, des améliorations possibles de l'implémentation de l'algorithme.

### 7.1 L'utilisation de martingales

Dans l'optique de généraliser le problème, nous allons montrer que l'on peut modéliser ledit problème par une surmartingale. L'algorithme choisi pour répondre au problème correspond alors à la *stratégie choisie pour optimiser la vitesse de convergence de la surmartingale*. L'intérêt d'une telle modélisation est d'avoir un cadre indépendant d'un certain nombre d'hypothèses et donc de pouvoir éventuellement obtenir des formules qui en sont indépendantes. On pourra alors savoir sous quelles hypothèses un algorithme est meilleur qu'un autre.

D'abord on explique brièvement ce qu'est une martingale et quelle est la théorie sous-

FIG. 10 – Résultats expérimentaux pour des certificats presque vrais.



jacente (section 7.1.1). Pour une définition complète, on pourra se référer aux sections 10 et 11 du cours IPPA ([Gal]). Pour une définition plus intuitive, on pourra lire la section 4.4 de *Randomized Algorithms* ([MR95]). Ensuite, on montre que le problème peut se modéliser à l'aide d'une martingale (section 7.1.2).

### 7.1.1 Définition d'une martingale

Nous commençons donc par définir ce sur quoi toute la théorie des martingales se base : le conditionnement et les filtrations. Pour cela on doit définir formellement l'espace des probabilités. Nous considérons donc un univers  $\Omega$ , et une tribu  $\mathcal{F}$  sur  $\Omega$  sur laquelle sera définie la mesure de probabilité  $\mathcal{P}$ . Soit  $X$  une variable aléatoire sur  $\Omega$  et  $F$  une sous-tribu de  $\mathcal{F}$ , on va définir l'espérance conditionnelle de  $X$  par rapport à  $F$ . Intuitivement, il s'agit de la variable aléatoire qui est la meilleure approximation de  $X$  sur  $F$ . Formellement, elle est définie ainsi :

$$\mathcal{E}(X|F) : \mathbf{F} \in F \mapsto \begin{cases} \frac{\mathcal{E}(X\mathbf{1}_{\mathbf{F}})}{\mathcal{P}_r(\mathbf{F})} & \text{si } \mathcal{P}_r(\mathbf{F}) \neq 0 \\ 0 & \text{sinon} \end{cases}$$

Où  $\mathbf{1}_{\mathbf{F}}$  est la fonction qui vaut 1 sur  $\mathbf{F}$  et 0 ailleurs. Une propriété importante de l'espérance conditionnelle est qu'elle est linéaire ( $\mathcal{E}(aX + bY|F) = a\mathcal{E}(X|F) + b\mathcal{E}(Y|F)$ ). De plus  $\mathcal{E}(X|F)$ , étant définie comme une variable aléatoire sur  $F$  est évidemment,  $F$ -mesurable<sup>8</sup>. De même, si  $X$  est  $F$ -mesurable, sa meilleure approximation sur  $F$  est elle-même puisque  $X$  peut-être définie sur  $F$ ; autrement dit,  $\mathcal{E}(X|F) = X$ .

Par exemple considérons un lancé de dé équilibré. On a  $\Omega = \{1, 2, 3, 4, 5, 6\}$ ,  $\mathcal{F} = \mathcal{P}(\Omega)$  et  $\mathcal{P}_r : \mathbf{F} \in \mathcal{F} \mapsto (\#\mathbf{F})/6$ . Soit  $X$  la variable aléatoire représentant la valeur du dé et  $F$  la

<sup>8</sup>Pour rappel : soit  $F$  une tribu, une variable  $X$  est  $F$ -mesurable si et seulement si elle peut être définie comme une variable aléatoire sur  $F$ .

tribu engendrée par {pair, impair}. On a  $\mathcal{E}(X|F)(\text{pair}) = 4$ , c'est à dire l'espérance de  $X$  sur les valeurs paires, et  $\mathcal{E}(X|F)(\text{impair}) = 3$ , c'est à dire l'espérance de  $X$  sur les valeurs impaires.

On définit ensuite une filtration d'une tribu comme une approximation de plus en plus précise de celle-ci. Formellement, une filtration de  $\mathcal{F}$  notée  $(F_n)_{n \in \mathbf{N}}$  est une suite de sous-tribus de  $\mathcal{F}$  telle que  $\forall n \in \mathbf{N}, F_n \subseteq F_{n+1}$ .

Par exemple, considérons  $\Omega = \{ \text{population française} \}$  et  $\mathcal{F} = \mathcal{P}(\Omega)$ . Une filtration peut être la suite de tribus engendrée par  $\Omega$  ( $F_0$ ), { hommes, femmes } ( $F_1$ ), { hommes de moins de 15 ans, hommes de plus de 15 ans, femmes de moins de 15 ans, femmes de plus de 15 ans } ( $F_2$ ). En conditionnant une variable aléatoire par une filtration, on obtient alors une suite variables aléatoires de plus en plus précises. Par exemple soit  $Y$  la variable aléatoire qui donne la masse d'un français choisi au hasard. On a alors  $\mathcal{E}(Y|F_0)$  qui traduit la masse moyenne des Français,  $\mathcal{E}(Y|F_1)$  qui correspond à une variable donnant soit la masse moyenne des hommes en France soit la masse moyenne des femmes en France, et  $\mathcal{E}(Y|F_2)$  qui ajoute la distinction entre les plus et moins de 15 ans.

Nous pouvons maintenant définir ce que sont les martingales. Le principe d'une martingale est de modéliser l'évolution d'une variable aléatoire « à mesure que nos connaissances augmentent ». On pose donc une filtration  $(F_n)_{n \in \mathbf{N}}$  qui correspond à l'augmentation des connaissances. En effet, plus  $n$  grandit, plus  $F_n$  grandit et est fine, plus on peut faire des distinctions entre les événements. On pose alors une suite de variables aléatoires  $(X_n)_{n \in \mathbf{N}}$  qui est la valeur que l'on souhaite estimer. On exige que pour tout  $n$ ,  $X_n$  est  $F_n$  mesurable. Cela signifie que l'on connaît la valeur de  $X_n$ , si l'on a connaissance de  $F_n$ . Une surmartingale (qui le type de martingale qui nous intéresse) possède de plus la propriété suivante :  $\forall n \in \mathbf{N}, X_n \geq \mathcal{E}(X_{n+1} | F_n)$ . Cela revient à dire qu'*en moyenne* la suite  $(X_n)_{n \in \mathbf{N}}$  décroît.

Nous allons prendre l'exemple de l'estimation qu'à notre algorithme probabiliste de se tromper.

### 7.1.2 Modéliser le comportement de l'algorithme par une surmartingale

Pour modéliser le comportement de l'algorithme par une surmartingale nous définissons d'abord l'univers, puis la filtration, pour enfin définir la suite de variables aléatoires.

Afin de définir  $\Omega$  nous considérons à nouveau les  $\mathcal{M}$  appartenances, notées  $A_1, \dots, A_{\mathcal{M}}$ , que l'algorithme peut vérifier. De plus, on partage l'algorithme de paramètres  $N$  et  $L$  en  $N \times L$  étapes atomiques. Chacune de ces étapes correspond au test d'une appartenance. On pose alors :

$$\Omega = \prod_{\substack{i \in \llbracket 1, \mathcal{M} \rrbracket \\ j \in \llbracket 1, N \times L \rrbracket}} \{ \text{appartenance } i \text{ vraie testée étape } j \} \times \{ \text{appartenance } i \text{ fausse testée étape } j \}$$

Intuitivement on peut calculer sur  $\Omega$  que chaque appartenance soit fausse ou non ainsi que les probabilités de tester chaque appartenance à chaque étape de l'algorithme.

Nous souhaitons modéliser le fait qu'au fur et à mesure que l'algorithme se déroule, nous connaissons les appartenances qui ont été testées ainsi que si elles sont vraies ou non. On définit alors  $(F_n)_{n \in \mathbf{N}}$  inductivement. On pose  $F_0 = \{\Omega, \emptyset\}$  c'est la tribu qui n'apporte aucun information sur  $X$ . La tribu  $F_{n+1}$ , par définition, est ensuite engendrée par l'ensemble suivant :

$$F_n \cup \bigcup_{i \in \llbracket 1, \mathcal{M} \rrbracket} \{ \text{appartenance } i \text{ vraie testée étape } n+1 \} \cup \{ \text{appartenance } i \text{ fausse testée étape } n+1 \}$$

On définit enfin  $(X_n)_{n \in \mathbf{N}}$  comme la probabilité de répondre oui à tort après l'étape  $n$ . On a bien la propriété que  $\forall n, X_n$  est  $F_n$ -mesurable car cela revient à dire que si l'on sait

tout de ce qui a été testé jusqu'à l'étape  $n$ , on sait calculer la probabilité de se tromper après l'étape  $n$ . On remarque que  $X_0$  correspond simplement à la probabilité de recevoir un certificat vrai.

Pour montrer que l'on a un surmartingale, on va montrer que la suite  $(X_n)_{n \in \mathbb{N}}$  décroît :

$$\forall n, \forall \mathbf{F} \in F_n, X_n(\mathbf{F}) \geq (x \mapsto X_{n+1}(\mathbf{F} \wedge \{x\})) \quad (6)$$

Intuitivement, cette formule dit que la valeur de  $X_n$  quand le résultat des  $n$  premiers lancés  $\mathbf{F}$  est connu est plus grande que toutes les valeurs possibles de  $X_{n+1}$  quand le résultat du  $(n+1)$ -ième lancé  $x$  sera connu. Cette propriété est plus forte que la propriété des surmartingales. En effet, si la suite décroît, *a fortiori* elle décroît *en moyenne*. Formellement, quand en prenant les espérances on obtient  $\forall n, \forall \mathbf{F} \in F_n, X_n \geq \mathcal{E}(X_{n+1} | \mathbf{F})$ . Autrement dit,  $\forall n \in \mathbb{N}, X_n \geq \mathcal{E}(X_{n+1} | F_n)$ .

Pour montrer que la propriété 6 est vraie, on considère sans perte de généralité que c'est l'appartenance  $A_1$  qui est testée à l'étape  $n+1$ <sup>9</sup>. Nous procédons alors à une disjonction de cas sur  $A_1$ .

**premier cas**  $A_1$  est fausse. Dans ce cas on répond non. On ne peut donc pas répondre oui à tort :  $X_{n+1}(\mathbf{F} \wedge \{x\}) = 0 \leq X_n(\mathbf{F})$  car  $X_n \in [0, 1]$  (évaluation de probabilité).

**second cas**  $A_1$  est vraie. Dans ce cas nous procédons à une disjonction de cas sur  $\mathbf{F}$ . Si il existe «  $A_i$  est fausse » dans  $\mathbf{F}$  alors  $X_n = X_{n+1} = 0$ . Ce cas n'est pas très intéressant en pratique car l'algorithme s'est arrêté et a répondu non. Sinon,  $\mathbf{F}$  est de la forme  $\{A_{f_1} \text{ est vraie}, \dots, A_{f_n} \text{ est vraie}\}$ . On va maintenant expliciter  $X_n$  et  $X_{n+1}$ .

La variable aléatoire  $X_n$  représente la probabilité de dire oui à tort après  $n$  étapes. On a alors  $\mathcal{P}_r(X_n | \mathbf{F}) = 1 - \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F})$ . En effet comme  $\mathbf{F}$  ne contient que des appartenances vraies, on répond oui après  $n$  étapes, donc la probabilité de dire oui à tort est la probabilité que toutes les appartenances ne soient pas vérifiées.

De la même manière, la variable aléatoire  $X_{n+1}$  représentant la probabilité de dire oui à tort après  $n+1$  étapes, on a  $\mathcal{P}_r(X_{n+1} | \mathbf{F}) = 1 - \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F} \wedge \{A_1 \text{ est vraie}\})$ .

Au final l'inégalité  $\mathcal{P}_r(X_n | \mathbf{F}) \geq \mathcal{P}_r(X_{n+1} | \mathbf{F})$  est équivalente à  $\mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F}) \leq \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F} \wedge \{A_1 \text{ est vraie}\})$ . Or, soit  $P = \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F})$ , on a :

$$\begin{aligned} P &= \mathcal{P}_r(A_1 \text{ vraie} \& \dots \& A_M \text{ vraie} | \mathbf{F}) \\ &= \mathcal{P}_r(A_1 \text{ vraie} | \mathbf{F}) \mathcal{P}_r(A_2 \text{ vraie} \& \dots \& A_M \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\}) \\ &\quad (\text{probabilités conditionnelles}) \\ &\leq \mathcal{P}_r(A_2 \text{ vraie} \& \dots \& A_M \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\}) \\ &\quad \text{car } \mathcal{P}_r(A_1 \text{ vraie} | \mathbf{F}) \leq 1 \\ &= \mathcal{P}_r(A_1 \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\}) \mathcal{P}_r(A_2 \text{ vraie} \& \dots \& A_M \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\} \wedge \{A_1 \text{ vraie}\}) \\ &\quad \text{car } \mathcal{P}_r(A_1 \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\}) = 1 \\ &= \mathcal{P}_r(A_1 \text{ vraie} \& \dots \& A_M \text{ vraie} | \mathbf{F} \wedge \{A_1 \text{ vraie}\}) \\ &\quad (\text{probabilités conditionnelles}) \\ &= \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F} \wedge \{A_1 \text{ est vraie}\}) \end{aligned}$$

Finalement, on a bien  $\mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F}) \leq \mathcal{P}_r(\text{toutes appartenances vraies} | \mathbf{F} \wedge \{A_1 \text{ est vraie}\})$ . □

<sup>9</sup>En effet, les appartenances sont numérotés arbitrairement.



En conclusion, on peut représenter la probabilité qu'a notre algorithme de se tromper par une surmartingale. Nous remarquons que nous n'avons pas utilisé d'hypothèse sur la dépendance entre les relations d'appartenance. De plus, notre preuve ne repose que sur le fait que l'algorithme teste des appartenances. La surmartingale que nous avons définie est donc plus générale que notre algorithme. D'abord, elle peut modéliser un algorithme différent de 1. Ensuite, on peut parfaitement la transposer à une analyse sensible au flot (ou **Label**  $\neq \{.\}$ , cf. section 4). Pour finir, un algorithme de vérification d'un certificat censé être une sémantique collectrice (section 3.4) peut aussi être modélisée par cette même surmartingale. En effet cela consiste à tester que des ensembles d'ensembles sont inclus les uns dans les autres. Ces tests peuvent être modélisés de manière atomique par plusieurs tests d'appartenance.

Un intérêt des martingales, outre un cadre général est l'existence d'outils mathématiques permettant de borner (voire calculer) la vitesse de convergence. Dans le cas de notre martingale, il est évident que  $(X_n)_{n \in \mathbb{N}}$  converge vers 0 après un nombre infini d'étapes (car il y a un nombre fini d'appartenance à vérifier). On peut cependant se demander à partir de quelle étape on a  $X_n \leq 0.05$  par exemple. La théorie des martingales propose un outil pour cela : le temps d'arrêt. Intuitivement, le temps d'arrêt est une étape auquel on peut choisir de s'arrêter. Autrement dit, quand on atteint le temps d'arrêt, on le sait. Par exemple, « la dernière étape à laquelle on teste  $A_1$  » n'est pas un temps d'arrêt car au moment où l'on atteint cette étape, on ne sait pas que l'on ne testera plus  $A_1$ . En revanche, « l'étape à partir de laquelle  $X_n \leq 0.05$  » est un temps d'arrêt car  $X_n$  décroît. On peut définir la variable aléatoire qui donne le numéro de l'étape ou le temps d'arrêt advient. La théorie des martingales fournit des outils pour calculer la loi de cette variable aléatoire.

Maintenant que nous avons une possible amélioration de l'approche théorique, nous proposons une manière d'améliorer l'implémentation.

## 7.2 Amélioration de l'implémentation

Nous proposons ici deux améliorations de l'implémentation possibles, elles sont indépendantes. La première vise à améliorer le temps d'exécution de notre algorithme. La seconde consiste à faire en sorte de pouvoir générer des équations qui correspondent à des vrais programmes.

En ce qui concerne la complexité temporelle, nous pensons l'améliorer en implémentant les ensembles sous forme tableaux de bits. Le bit  $n$  du tableau correspondrait à la valeur  $n$ , un bit à 0 signifierait que la valeur n'appartient pas à l'ensemble, un bit à 1 signifierait que la valeur appartient à l'ensemble. L'avantage de cette implémentation est que la complexité spatiale reste en  $O(KL)^{10}$  alors que la complexité temporelle passe en  $O(NL)$  puisque les accès aux tableaux se feront en temps constant. La difficulté de cette implémentation est de prendre efficacement un élément aléatoire dans un ensemble (étape 5 de l'algorithme, cf. 5.1).

À propos de la génération d'équations à partir de vrais programmes, nous étudions la possibilité de nous interfacer à `sawja`<sup>11</sup>. Cela nous permettrait alors d'analyser du `byte code java` ce qui est précisément le contexte dans lequel nous nous plaçons en section 1.

<sup>10</sup>La constante est à ce propos réduite puisque l'on n'a plus besoin que d'un bit par valeur plutôt que la taille d'un entier (32 bits le plus souvent).

<sup>11</sup><http://sawja.inria.fr/>

## 8 Conclusion

Pour conclure, nous avons, à notre connaissance, réalisé le premier algorithme de vérification d'un certificat par rapport à une analyse statique. Nous sommes, pour ce faire, placé dans un sous-contexte précis que nous maîtrisons afin de pouvoir acquérir une meilleure connaissance de la problématique. De plus, nous avons implémenté notre solution pour ne pas perdre vue l'objectif applicatif que nous avons. L'implémentation a été testée et les premiers résultats sont encourageants bien que du travail puisse encore être fait en termes de complexité temporelle et de pertinences des équations testées.

De plus, nous sommes maintenant en mesure de proposer des axes de recherche plus généraux afin de ne plus avoir s'astreindre au sous-problème que nous nous sommes posé. Nous estimons que l'utilisation de martingales, notamment, permet d'aller dans ce sens.

## Références

- [Gal] Jean-François Le Gall. Intégration, probabilités, processus aléatoires. [www.math.u-psud.fr/~jflgall/IPPA2.pdf](http://www.math.u-psud.fr/~jflgall/IPPA2.pdf).
- [GN03] Sumit Gulwani and Georges C. Necula. Discovering affine equalities using random interpretation. In *30<sup>th</sup> Annual Symposium ACM on Principles of Programming Languages*, pages 74–84, 2003.
- [GN04a] Sumit Gulwani and Georges C. Necula. Global value numbering using random interpretation. In *31<sup>st</sup> Annual Symposium ACM on Principles of Programming Languages*, pages 342–352, 2004.
- [GN04b] Sumit Gulwani and Georges C. Necula. A polynomial-time algorithm for global value numbering. In *11<sup>th</sup> Static Analysis Symposium*, volume 3148 of *LNCS*, pages 212–227. Springer, 2004.
- [GN05] Sumit Gulwani and Georges C. Necula. Precise interprocedural analysis using random interpretation. In *32<sup>nd</sup> Annual Symposium ACM on Principles of Programming Languages*, 2005.
- [Gul05] Sumit Gulwani. *Program analysis using random interpretation*. PhD thesis, University of California, Berkeley, 2005.
- [MR95] Rajeev Motwani and Prabhakar Raghavan. *Randomized algorithms*. Cambridge university press, 1995.
- [NG11] Rupesh Nasre and Ramaswamy Govindarajan. Points-to analysis as a system of linear equations. In *Lecture Notes in Computer Science*, volume 6337, pages 422–438. Springer, 2011.
- [WACL05] John Whaley, Dzintars Avots, Michael Carbin, and Monica S. Lam. Using datalog with binary decision diagrams for program analysis. In *Lecture Notes in Computer Science*, volume 3780, pages 98–118. Springer, 2005.
- [WL04] John Whaley and Monica S. Lam. Cloning-based context sensitive pointer alias analysis using binary decision diagrams. In *conference on Programming language Design and Implementation*, 2004.