

Inférence de flots d'information pour ML

Vincent Simonet
INRIA Rocquencourt – Projet Cristal

Séminaire DÉMONS (L.R.I.)

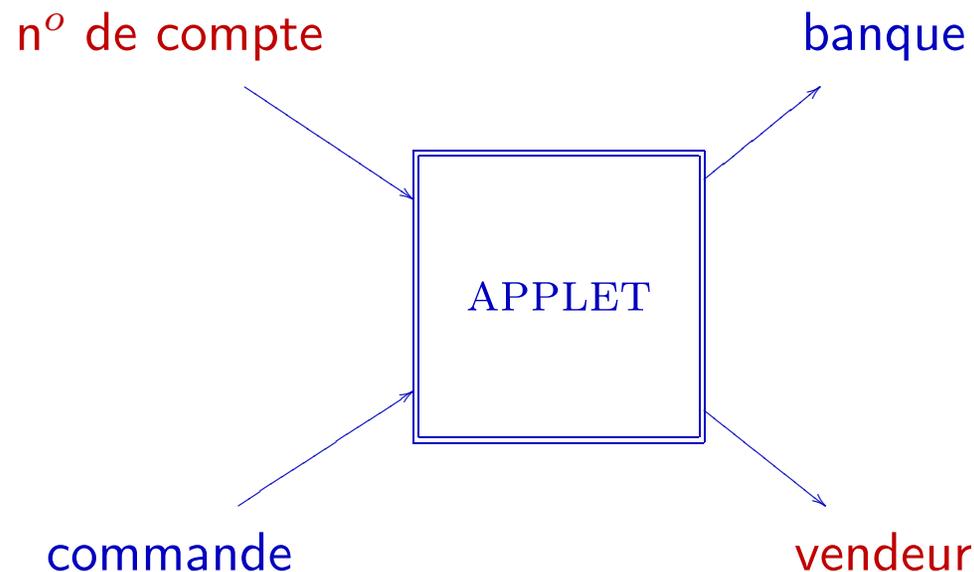
Vendredi 30 novembre 2001

Vincent.Simonet@inria.fr

<http://cristal.inria.fr/~simonet/>

Analyse de flots d'information

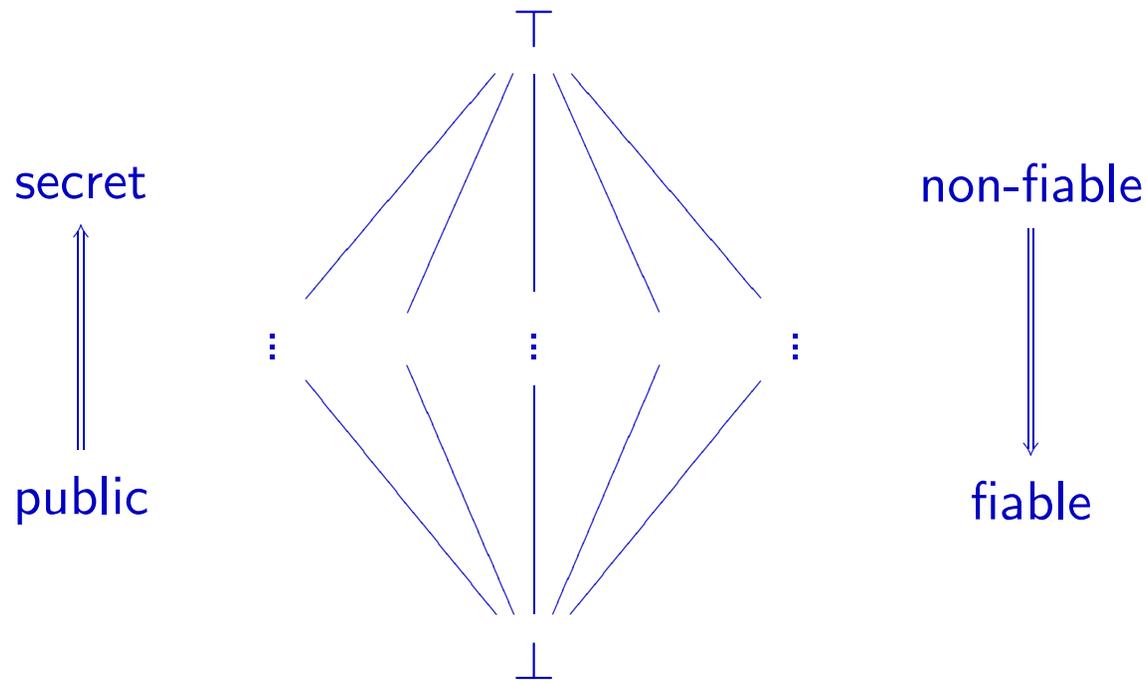
L'analyse de flots d'information consiste à analyser les dépendances entre les entrées et les sorties d'un programme de manière à vérifier qu'il satisfait certaines propriétés de confidentialité ou d'intégrité vis-à-vis des données qu'il manipule.



Analyse de flots d'information

Niveaux d'information

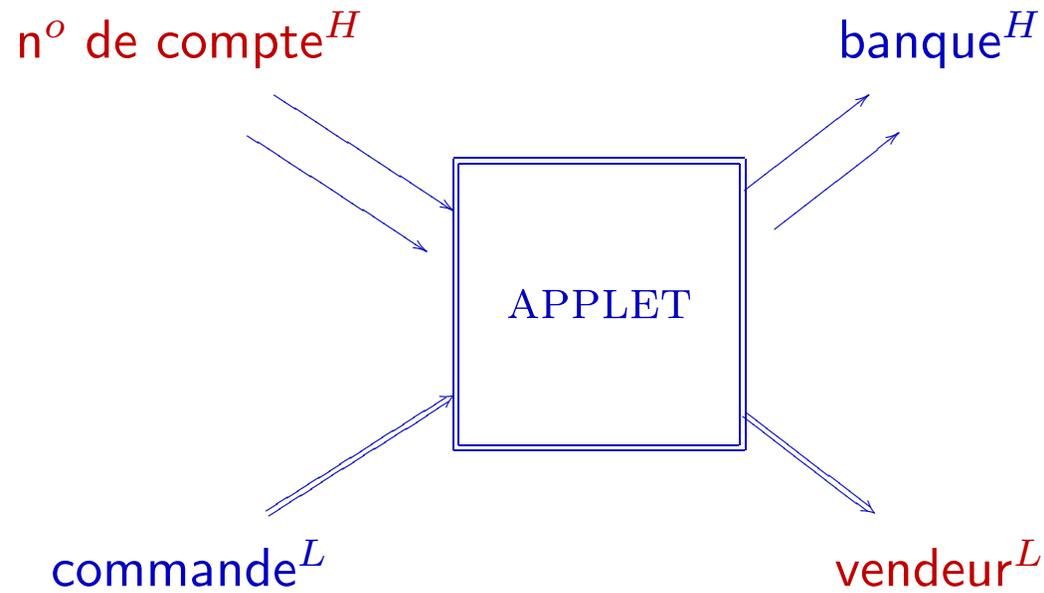
On associe à chaque donnée un **niveau d'information**, dans un treillis \mathcal{L} . Ces niveaux peuvent représenter différentes propriétés : sécurité, intégrité...



Pour la suite, on fixe $\mathcal{L} = \{L \leq H\}$.

Analyse de flots d'information

Non-interférence



L'indépendance entre une entrée et une sortie est exprimée sous forme d'une propriété de **non-interférence**.

Analyse de flots d'information

Typage

Les flots d'informations peuvent être représentés par des types annotés :

$$\text{compte}^H \times \text{commande}^L \rightarrow \text{banque}^H \times \text{vendeur}^L$$

ou plus généralement par des schémas polymorphes contraints :

$$(\forall \alpha \beta \gamma \delta) [\alpha \sqcup \beta \leq \gamma, \beta \leq \delta] \text{compte}^\alpha \times \text{commande}^\beta \rightarrow \text{banque}^\gamma \times \text{vendeur}^\delta$$

Analyse de flots d'information

Systemes existants

Dennis Volpano et Geoffrey Smith (1997)

Systeme de type sur un langage impératif simple. Limité au premier ordre et à un nombre fini de références globales.

Nevin Heintze et Jon G. Riecke SLam Calculus (1997)

Lambda-calcul avec références et processus. Le typage des références n'est pas fin. Il n'est pas énoncé ni prouvé de propriété de sûreté.

Andrew C. Myers JFlow (1999)

Analyse de flots d'information sur Java. Le système est complexe et n'est pas prouvé.

Steve Zdancewic et Andrew C. Myers (2001)

Analyse sur un langage de bas niveau avec continuations linéaires.

Le langage ML

λ -calcul en appel par valeur avec let-polymorphe

| | | |
|-----------|--------------------|-------------------------|
| x | k | $\lambda x.e$ |
| $e_1 e_2$ | let $x = v$ in e | bind $x = e_1$ in e_2 |

avec références

| | | |
|---------|--------------|------|
| ref e | $e_1 := e_2$ | $!e$ |
|---------|--------------|------|

et exceptions

| | | | |
|-----------------|-----------|--|----------------------------|
| εe | raise e | e_1 handle $\varepsilon x \succ e_2$ | e_1 handle $x \succ e_2$ |
|-----------------|-----------|--|----------------------------|

Le langage ML

Formes v -normales

$$v ::= x \mid k \mid \lambda x.e \mid \varepsilon v$$

$$e ::= v v \mid \text{ref } v \mid v := v \mid !v \mid \text{raise } v \mid \text{let } x = v \text{ in } e \mid E[v]$$

$$E ::= \text{bind } x = [] \text{ in } e \mid [] \text{ handle } \varepsilon x \succ e \mid [] \text{ handle } x \succ e$$

On peut mettre toute expression habituelle sous cette forme, à condition de choisir un ordre d'évaluation. Par exemple :

$$e_1 e_2 \Rightarrow \begin{cases} \text{bind } x_1 = e_1 \text{ in } (\text{bind } x_2 = e_2 \text{ in } x_1 x_2) & \text{éval. de gauche à droite} \\ \text{bind } x_2 = e_2 \text{ in } (\text{bind } x_1 = e_1 \text{ in } x_1 x_2) & \text{éval. de droite à gauche} \end{cases}$$

Exemples de flot d'information

Flots directs

$x := \text{not } y$

$x := (\text{if } y \text{ then false else true})$

Flots indirects

if y then $x := \text{false}$ else $x := \text{true}$

$x := \text{true}; \text{if } y \text{ then } x := \text{false} \text{ else } ()$

$x := \text{true}; (\text{if } y \text{ then raise } A \text{ else } ()) \text{ handle } _ \succ x := \text{false}$

Exemples de flot d'information

Niveau associé au « compteur de programme »

Supposons que y représente une donnée « secrète » (H) :

$$\text{if } y \text{ then } \underbrace{x := \text{false}}_{pc=H} \text{ else } \underbrace{x := \text{true}}_{pc=H}$$
$$\underbrace{x := \text{true}}_{pc=L}; \text{ if } y \text{ then } \underbrace{x := \text{false}}_{pc=H} \text{ else } ()$$
$$\text{let } f = \lambda b.(x := b) \text{ in } \underbrace{f \text{ true}}_{pc=L}; \text{ if } y \text{ then } \underbrace{f \text{ false}}_{pc=H} \text{ else } ()$$

À chaque point du programme correspond un niveau pc qui représente les informations que l'on peut apprendre en sachant que ce point a été exécuté.

Système de types

Approche semi-syntaxique

(exemples dans le cas du typage ML)

Système logique

Types bruts

e.g. int , $\text{int} \rightarrow \text{int} \dots$

Polytypes

e.g. $\{t \rightarrow t \mid t \text{ type brut}\}$

Système syntaxique

Expressions de type

e.g. int , α , $\alpha \rightarrow \alpha \dots$

Schémas de type

e.g. $\forall \alpha. \alpha \rightarrow \alpha$

On raisonne sur le système logique. Le système syntaxique est interprété dans le système logique.

Système de types

Algèbre de types

Les **niveaux** d'information ℓ, pc appartiennent à un treillis \mathcal{L} .

Les exceptions sont décrites par des **rangées** d'alternatives r :

$$\begin{aligned} a & ::= \text{Abs} \mid \text{Pre } pc \\ r & ::= \{ \varepsilon \mapsto a \}_{\varepsilon \in \mathcal{E}} \end{aligned}$$

Les types sont annotées par des **niveaux** et des **rangées** :

$$t ::= \text{int}^\ell \mid \text{unit} \mid (t \xrightarrow{pc [r]} t)^\ell \mid t \text{ ref}^\ell \mid r \text{ exn}^\ell$$

Système de types

Jugements

On distingue deux formes de jugements.

Jugements sur des valeurs

$$\Gamma \vdash v : t$$

Jugements sur des expressions

$$pc, \Gamma \vdash e : t [r]$$

Système de types

Contraintes

Contraintes de sous-typage $t_1 \leq t_2$

La relation de sous-typage étend l'ordre sur les niveaux aux types. Par exemple :

$$\text{int}^{\ell_1} \leq \text{int}^{\ell_2} \Leftrightarrow \ell_1 \leq \ell_2 \quad t_1 \text{ ref}^{\ell_1} \leq t_2 \text{ ref}^{\ell_2} \Leftrightarrow t_1 = t_2 \text{ et } \ell_1 \leq \ell_2$$

Gardes $\ell \triangleleft t$

Les gardes permettent de « marquer » un type par un niveau d'information :

$$pc \triangleleft \text{int}^{\ell} \Leftrightarrow pc \leq \ell \quad pc \triangleleft t \text{ ref}^{\ell} \Leftrightarrow pc \leq \ell$$

Système de types

Sous-typage et polymorphisme

Sous-typage et polymorphisme interviennent de manière orthogonale :

Sous-typage Permet d'augmenter librement le niveau d'une donnée (e.g. considérer *secrète* une donnée *publique* ou *non-fiable* une donnée *fiable*) :

$$\frac{\Gamma \vdash v : t \quad t \leq t'}{\Gamma \vdash v : t'}$$

Polymorphisme Nécessaire pour pouvoir utiliser la même fonction avec des entrées de niveaux différents :

$$\text{let } succ = \lambda x.(x + 1)$$

Système de types

Références

$$\frac{\text{REF} \quad \Gamma \vdash v : t \quad pc \triangleleft t}{pc, \Gamma \vdash \text{ref } v : t \text{ ref}^l [r]}$$

$$\frac{\text{DEREF} \quad \Gamma \vdash v : t' \text{ ref}^l \quad t' \leq t \quad l \triangleleft t}{pc, \Gamma \vdash !v : t [r]}$$

$$\frac{\text{ASSIGN} \quad \Gamma \vdash v_1 : t \text{ ref}^l \quad \Gamma \vdash v_2 : t \quad pc \triangleleft t \quad l \triangleleft t}{pc, \Gamma \vdash v_1 := v_2 : \text{unit} [r]}$$

Le contenu d'une référence a un niveau supérieur à égal

- au pc du point où la référence est créée,
- au pc de chaque point du programme où son contenu est susceptible d'être modifié.

Système de types

Exceptions

RAISE

$$\frac{\Gamma \vdash v : \text{typexn}(\varepsilon)}{pc, \Gamma \vdash \text{raise } (\varepsilon v) : * \quad [\varepsilon : \text{Pre } pc; \partial\text{Abs}]}$$

HANDLE

$$\frac{\begin{array}{l} pc, \Gamma \vdash e_1 : t \quad [\varepsilon : \text{Pre } pc'; r_1] \\ pc \sqcup pc', \Gamma[x \mapsto \text{typexn}(\varepsilon)] \vdash e_2 : t \quad [\varepsilon : a_2; r_2] \quad pc' \triangleleft t \end{array}}{pc, \Gamma \vdash e_1 \text{ handle } \varepsilon x \succ e_2 : t \quad [\varepsilon : a_2; r_1 \sqcup r_2]}$$

Non-interférence

On considère une expression e de type int^L avec un « trou » x marqué H :

$$(x \mapsto t) \vdash e : \text{int}^L \qquad H \triangleleft t$$

Théorème de non-interférence

$$\text{Si } \begin{cases} \vdash v_1 : t \\ \vdash v_2 : t \end{cases} \text{ et } \begin{cases} e[x \Leftarrow v_1] \rightarrow^* v'_1 \\ e[x \Leftarrow v_2] \rightarrow^* v'_2 \end{cases} \text{ alors } v'_1 = v'_2$$

Le résultat de l'évaluation de e ne dépend pas de la valeur d'entrée placée dans le trou.

Preuve de non-interférence

1. On définit une extension *ad hoc* du langage qui permet de raisonner sur les points communs et les différences entre deux programmes.
2. On vérifie que le système de types pour le langage étendu possède *subject reduction*.
3. On montre que la non-interférence pour le langage initial est une conséquence de *subject reduction*.

Preuve de non-interférence

Calcul partagé

Le calcul partagé permet de considérer simultanément les exécutions indépendantes de deux expressions en gardant trace des sous-expressions qu'elles partagent.

Syntaxe

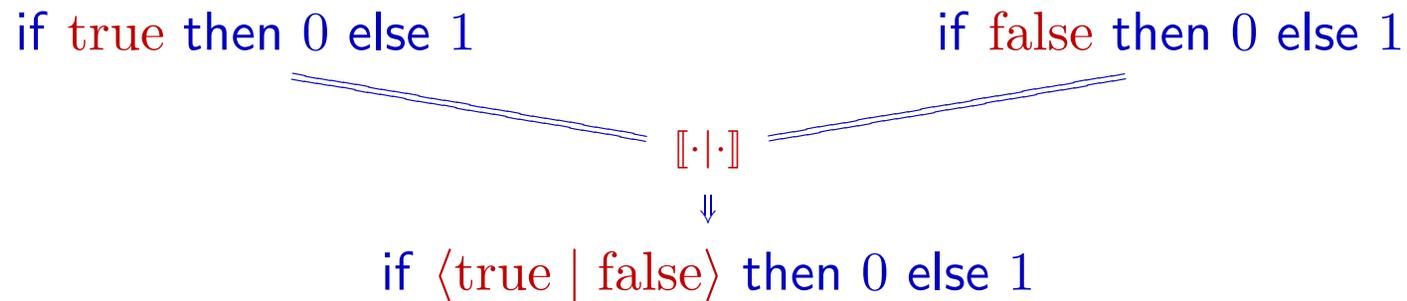
$$v ::= \dots \mid \langle v \mid v \rangle \qquad e ::= \dots \mid \langle e \mid e \rangle$$

On ne considère que les expressions où les $\langle \dots \mid \dots \rangle$ ne sont pas imbriqués.

Preuve de non-interférence

Codage

Une expression du calcul partagé permet de représenter une paire d'expressions du calcul simple en partageant certaines sous-expressions.

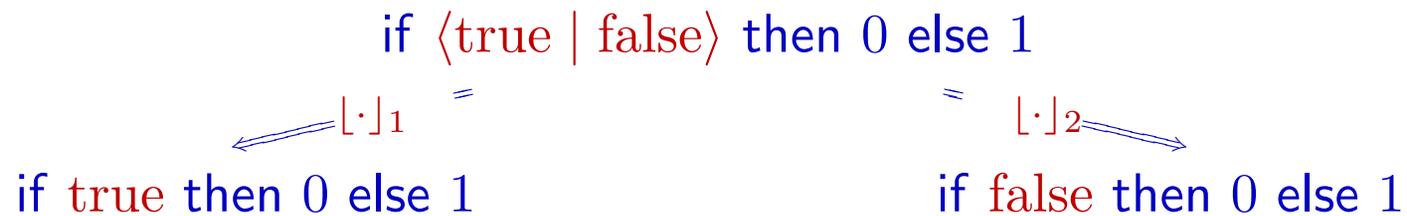


<if true then 0 else 1 | if false then 0 else 1>

Preuve de non-interférence

Projections

Les projections $[\cdot]_1$ et $[\cdot]_2$ permettent de retrouver les deux expressions simples codées par une expression du calcul partagé.



Preuve de non-interférence

Réduction du calcul partagé

On dérive la sémantique du calcul partagé de celle du calcul simple : à chaque fois que des crochets bloquent la réduction, il faut les déplacer.

$$(\lambda x.e) v \rightarrow e[x \leftarrow v] \quad (\beta)$$

$$\langle v_1 \mid v_2 \rangle v \rightarrow \langle v_1 [v]_1 \mid v_2 [v]_2 \rangle \quad (\text{lift-app})$$

Exemples

$$\begin{aligned} \langle \lambda x.x \mid \lambda x.x + 1 \rangle \mathbf{3} &\rightarrow \langle (\lambda x.x) \mathbf{3} \mid (\lambda x.x + 1) \mathbf{3} \rangle \\ &\rightarrow \langle \mathbf{3} \mid (\lambda x.x + 1) \mathbf{3} \rangle \rightarrow \langle \mathbf{3} \mid 4 \rangle \end{aligned}$$

$$\begin{aligned} \langle \lambda x.x \mid \lambda x.x + 1 \rangle \langle \mathbf{3} \mid \mathbf{2} \rangle &\rightarrow \langle (\lambda x.x) \mathbf{3} \mid (\lambda x.x + 1) \mathbf{2} \rangle \\ &\rightarrow \langle \mathbf{3} \mid (\lambda x.x + 1) \mathbf{3} \rangle \rightarrow \langle \mathbf{3} \mid \mathbf{3} \rangle \end{aligned}$$

Preuve de non-interférence

Simulation

Correction

$$\text{Si } e \rightarrow e' \quad \text{alors } \begin{cases} [e]_1 \rightarrow^= [e']_1 \\ [e]_2 \rightarrow^= [e']_2 \end{cases}$$

(calcul partagé) (calcul simple)

Complétude

$$\text{Si } \begin{cases} e_1 \rightarrow^* v_1 \\ e_2 \rightarrow^* v_2 \end{cases} \quad \text{alors } [[e_1 \mid e_2]] \rightarrow^* [[v_1 \mid v_2]]$$

(calcul simple) (calcul partagé)

Preuve de non-interférence

Typage de $\langle \dots \mid \dots \rangle$

Pour que $\langle v_1 \mid v_2 \rangle$ ait le type t , on impose que $H \triangleleft t$:

$$\frac{\text{BRACKET} \quad \Gamma \vdash v_1 : t \quad \Gamma \vdash v_2 : t \quad H \triangleleft t}{\Gamma \vdash \langle v_1 \mid v_2 \rangle : t}$$

Ainsi :

- Une valeur de type int^H peut être en entier k ou un crochet $\langle k_1 \mid k_2 \rangle$.
- Une valeur de type int^L ne peut être qu'un entier k .

Preuve de non-interférence

Subject reduction et non-interférence

Soit $(x \mapsto t) \vdash e : \text{int}^L$ avec $H \triangleleft t$.

Subject-reduction

Si $\vdash e' : \text{int}^L$ et $e' \rightarrow^* v'$ alors $\vdash v' : \text{int}^L$

$$\begin{array}{c} \uparrow \\ e' = e[x \Leftarrow v] \\ \downarrow \end{array}$$

$$\begin{array}{c} \downarrow \\ v' = k \\ \downarrow \end{array}$$

Non-interférence pour le calcul partagé

Si $\vdash v : t$ et $e[x \Leftarrow v] \rightarrow^* v'$ alors $[v']_1 = [v']_2$

Preuve de non-interférence

Non-interférence

Soit $(x \mapsto t) \vdash e : \text{int}^L$ avec $H \triangleleft t$.

Non-interférence pour le calcul partagé

Si $\vdash v : t$ et $e[x \Leftarrow v] \rightarrow^* v'$ alors $[v']_1 = [v']_2$

$$v = \langle v_1 \mid v_2 \rangle$$

$$v' = \llbracket v'_1 \mid v'_2 \rrbracket$$

Non-interférence pour le calcul simple

Si $\begin{cases} \vdash v_1 : t \\ \vdash v_2 : t \end{cases}$ et $\begin{cases} e[x \Leftarrow v_1] \rightarrow^* v'_1 \\ e[x \Leftarrow v_2] \rightarrow^* v'_2 \end{cases}$ alors $v'_1 = v'_2$

Extension du langage

On souhaite généralement étendre le langage étudié :

Pour accroître la richesse du langage Ajout de sommes, de paires. Étude d'un cas général pour des « primitives » des langages réels (opérations arithmétiques, de comparaison).

Pour mieux typer certains idiomes

$e_1 \text{ finally } e_2 \hookrightarrow \text{bind } x = (e_1 \text{ handle } y \succ e_2; \text{ raise } y) \text{ in } e_2; x$

$e_1 \text{ handle } x \succ e_2 \text{ reraise} \hookrightarrow e_1 \text{ handle } x \succ (e_2; \text{ raise } x)$

Extension du langage

Sommes et paires

$$t ::= \dots \mid (t_1 + t_2)^\ell \mid t_1 \times t_2$$

Il n'est pas nécessaire de munir le constructeur de types \times d'un niveau d'information propre. Les niveaux présents dans les types des composantes suffisent :

$$\begin{aligned} \ell \triangleleft t_1 \times t_2 &\Leftrightarrow \ell \triangleleft t_1 \wedge \ell \triangleleft t_2 \\ t_1 \times t_2 \blacktriangleleft \ell &\Leftrightarrow t_1 \blacktriangleleft \ell \wedge t_2 \blacktriangleleft \ell \end{aligned}$$

Extension du langage

Primitives

Les primitives monomorphes telles que les opérations arithmétiques n'entraînent pas de difficulté :

$$\frac{\Gamma \vdash v_1 : \text{int}^\ell \quad \Gamma \vdash v_2 : \text{int}^\ell}{pc, \Gamma \vdash v_1 + v_2 : \text{int}^\ell \quad [\partial\text{Abs}]}$$

Le typage de primitives polymorphes qui « filtrent » la structure de leurs entrées nécessite une nouvelle forme de contraintes :

$$\frac{\Gamma \vdash v_1 : t \quad \Gamma \vdash v_2 : t \quad t \blacktriangleleft \ell}{pc, \Gamma \vdash v_1 = v_2 : \text{bool}^\ell \quad [\partial\text{Abs}]} \quad \frac{\Gamma \vdash v : t \quad t \blacktriangleleft \ell}{pc, \Gamma \vdash \text{hash } v : \text{int}^\ell \quad [\partial\text{Abs}]}$$

$t \blacktriangleleft \ell$ contraint *tous* les niveaux d'information apparaissant dans t et ses sous termes à être inférieurs à ℓ .

Vers une extension du compilateur Caml

Le langage étudié permet de prendre en compte la quasi-totalité du langage Caml (à l'exception de la bibliothèque `threads`).

L'implantation d'un prototype est en cours. Elle nécessite de résoudre plusieurs problèmes liés à l'utilisation d'un système de types doté de sous-typage :

- Efficacité de l'algorithme d'inférence
- Lisibilité des types inférés
- Clarté des messages d'erreur
- ...

Vers une extension du compilateur Caml

Inférence de types

Un algorithme d'inférence comporte deux parties relativement distinctes.

Un jeu de règles d'inférence On peut dériver de manière quasi-systématique un jeu de règles d'inférence d'un système tel que le nôtre.

$$\frac{\text{REF} \quad \Gamma \vdash v : t \quad pc \triangleleft t}{pc, \Gamma \vdash \text{ref } v : t \text{ ref}^\ell [r]} \rightsquigarrow \frac{\text{INF-REF} \quad \Gamma, C \vdash v : \alpha}{\pi, \Gamma, C \cup \{\beta = \alpha \text{ ref}^\lambda, \pi \triangleleft \alpha\} \vdash \text{ref } v : \beta [\rho]}$$

Un solveur Les schémas de types comportent des ensembles de contraintes. Il faut tester leur satisfiabilité et les simplifier.

Vers une extension du compilateur Caml

Déclarations de types

```
type ('a, 'b) list =  
  []  
  | (::) of 'a * ('a, 'b) list  
level 'b
```

```
type ('a, 'b, 'c) queue = {  
  mutable in: ('a, 'b) list;  
  mutable out: ('a, 'b) list  
}  
level 'c
```

Vers une extension du compilateur Caml

Exemple : listes

```
let rec length = function
  [] -> 0
  | _ :: l -> 1 + length l
```

```
val length :  $\forall[\alpha \leq \beta]. * \text{list}^\alpha \rightarrow \text{int}^\beta$ 
```

```
let rec iter f = function
  [] -> ()
  | x :: l -> f x; iter f l
```

```
val iter :  $\forall[\sqcup \delta \leq \gamma]. (\alpha \xrightarrow{\gamma [\delta]} *)^\gamma \rightarrow \alpha \text{list}^\gamma \xrightarrow{\gamma [\delta]} \text{unit}$ 
```

Vers une extension du compilateur Caml

Exemple : queues

```
let push p elt =  
  p.in <- elt :: p.in
```

```
val push :  $\forall[\gamma \leq \beta].(\alpha, \beta) \text{ queue}^\gamma \rightarrow \alpha \xrightarrow{\gamma [*]} \text{unit}$ 
```

```
let rec pop p = match p.out with  
  hd :: tl -> p.out <- tl; hd  
| [] -> match p.in with  
  [] -> raise Empty  
  | _ -> balance p; pop p
```

```
val pop :  $\forall[\alpha \leq \alpha', \beta \triangleleft \alpha', \gamma \sqcup \pi \leq \beta].(\alpha, \beta) \text{ queue}^\gamma \xrightarrow{\pi [\text{Empty}:\beta; *]} \alpha'$ 
```

Référence

- [1] François Pottier and Vincent Simonet. [Information flow inference for ML](#). In *Proceedings of the 29th ACM Symposium on Principles of Programming Languages (POPL'02)*. <http://crystal.inria.fr/~simonet/publis/fpottier-simonet-popl02.ps.gz>.