

Coccinelle

un outil de reconnaissance
de motifs et de transformation de code C

Julien Brunel

Office National d'Études et de Recherches Aérospatiales
Département Traitement de l'Information et Modélisation

5 mai 2009

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking
- 6 Conclusion

Outline

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking
- 6 Conclusion

Coccinelle

Coccinelle

- Reconnaissance de motifs dans du code C
- Transformation de code
- Basé sur une extension de CTL (model checking sur le graphe de flot de contrôle)

Utilisation

- À l'origine, pour pallier le problème de l'*évolution collatérale* dans le noyau linux (modification d'API)
- Deuxième utilisation : détection et correction de bugs

Outline

- 1 Introduction
- 2 Langage de Coccinelle**
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking
- 6 Conclusion

SmPL : 'Semantic Patch Language'

- Comme le langage de patch (Unix), mais **générique**.
- On s'abstrait des détails inutiles
 - Espaces, indentations, commentaires
 - Choix du nom des variables (**metavariab**les)
 - Autres variations dans la manière de coder («**isomorphismes**»)
- On raisonne sur les chemins du **Graphe de Flot de Contrôle**

Un **patch sémantique** peut modifier un grand nombre de fichiers.

Structure et syntaxe d'un patch sémantique

```
@regle_1@
```

```
//Déclaration de meta-variables;
```

```
expression x;
```

```
@@
```

```
//Description d'un motif
```

```
x = kmalloc(...)
```

```
...
```

```
- kfree(x)
```

```
@regle_2@
```

```
//Déclaration de meta-variables;
```

```
expression y != regle_1.x;
```

```
@@
```

```
//Description d'un motif
```

Recherche de fuites de mémoire

```
@leaks exists@
```

```
expression  $x$ ,  $E$ ,  $F$ ;
```

```
statement  $S$ ;
```

```
@@
```

```
*  $x = \text{kmalloc}(\dots)$ 
```

```
...
```

```
if ( $x == \text{NULL}$ )  $S$ 
```

```
... when  $!= \text{kfree}(x)$ 
```

```
    when  $!= (x = E \mid E = x)$ 
```

```
* return  $F$ ;
```


Recherche de fuites de mémoire

```
@leaks exists@
```

```
expression  $x$ ,  $E$ ;
```

```
statement  $S$ ;
```

```
constant  $C$ ;
```

```
@@
```

```
*  $x = \text{kmalloc}(\dots)$ 
```

```
...
```

```
if ( $x == \text{NULL}$ )  $S$ 
```

```
... when  $!= \text{kfree}(x)$ 
```

```
    when  $!= (x = E \mid E = x)$ 
```

```
* return (  $-C \mid \text{NULL}$  );
```

Recherche de fonctions qui retournent NULL

```

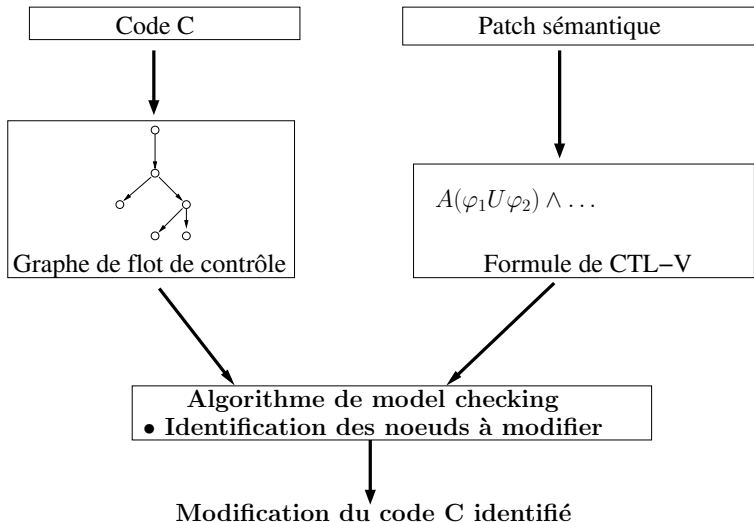
@returns_null exists@
identifiant f;
expression E, E1;
@@
f(...) {
    ...
    (
        return NULL;
    |
        E = NULL;
        ... when != E = E1
        return E;
    )
}

```

Correction des appels à `usb_submit_urb`

```
@@
expression lock, flags;
expression urb;
@@
    spin_lock_irqsave(lock, flags) ;
    <...
- usb_submit_urb(urb)
+ usb_submit_urb(urb, GFP_ATOMIC)
    ...>
    spin_unlock_irqrestore(lock, flags) ;
@@
expression urb;
@@
- usb_submit_urb(urb)
+ usb_submit_urb(urb, GFP_KERNEL)
```

Comment ça marche ?



Outline

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle**
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking
- 6 Conclusion

Exigences pour la sémantique

- 1 gérer les méta-variables
- 2 permettre aux méta-variables de prendre différentes valeurs dans des chemins différents du GFC
- 3 enregistrer les différentes valeurs des méta-variables
- 4 enregistrer les endroits où les transformations doivent être effectuées

CTL translation and model checking

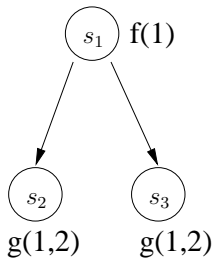
Semantic patch

$f(\dots);$

- $g(\dots);$

CTL representation

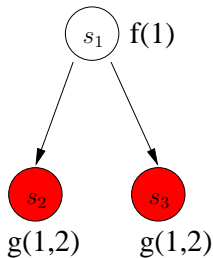
$f(\dots) \wedge AX g(\dots)$



CTL translation and model checking

Semantic patch

$f(\dots);$
 ~~$g(\dots);$~~



CTL representation

$f(\dots) \wedge AX g(\dots)$

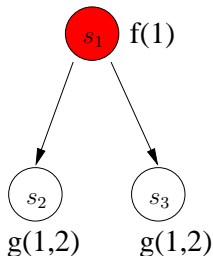
Model checking algorithm

$$SAT(g(\dots)) = \{s_2, s_3\}$$

CTL translation and model checking

Semantic patch

$f(\dots);$
 ~~$g(\dots);$~~



CTL representation

$f(\dots) \wedge AX g(\dots)$

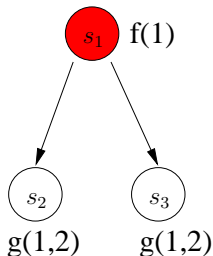
Model checking algorithm

$$\begin{aligned}
 \text{SAT}(g(\dots)) &= \{s_2, s_3\} \\
 \text{SAT}(AX g(\dots)) &= \{s_1\}
 \end{aligned}$$

CTL translation and model checking

Semantic patch

$f(\dots);$
 $- g(\dots);$



CTL representation

$f(\dots) \wedge AX g(\dots)$

Model checking algorithm

$$\begin{aligned}
 \text{SAT}(g(\dots)) &= \{s_2, s_3\} \\
 \text{SAT}(AX g(\dots)) &= \{s_1\} \\
 \text{SAT}(f(\dots) \wedge AX g(\dots)) &= \{s_1\}
 \end{aligned}$$

Adding an environment (CTL-FV)

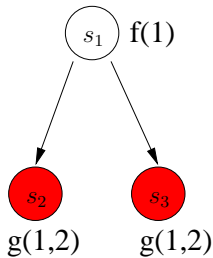
Semantic patch

$f(x);$

- $g(x, y);$

CTL representation

$$f(x) \wedge AX g(x, y)$$



Model checking algorithm

$$\text{SAT}(g(x, y)) = \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 2])\}$$

Adding an environment (CTL-FV)

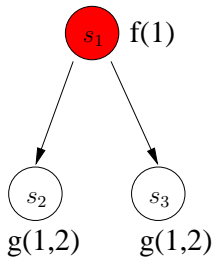
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$$f(x) \wedge AX g(x, y)$$



Model checking algorithm

$$\begin{aligned} \text{SAT}(g(x, y)) &= \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 2])\} \\ \text{SAT}(AX g(x, y)) &= \{(s_1, [x \mapsto 1, y \mapsto 2])\} \end{aligned}$$

Adding an environment (CTL-FV)

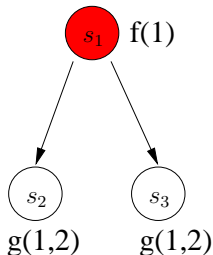
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$$f(x) \wedge AX g(x, y)$$



Model checking algorithm

$$\begin{aligned}
 \text{SAT}(g(x, y)) &= \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 2])\} \\
 \text{SAT}(AX g(x, y)) &= \{(s_1, [x \mapsto 1, y \mapsto 2])\} \\
 \text{SAT}(f(x) \wedge AX g(x, y)) &= \{(s_1, [x \mapsto 1, y \mapsto 2])\}
 \end{aligned}$$

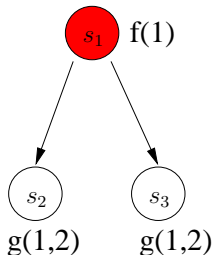
Adding an environment (CTL-FV)

Semantic patch

 $f(x);$ $- g(x, y);$

CTL representation

$$f(x) \wedge AX g(x, y)$$



Model checking algorithm

$$\text{SAT}(g(x, y)) = \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 2])\}$$

$$\text{SAT}(AX g(x, y)) = \{(s_1, [x \mapsto 1, y \mapsto 2])\}$$

$$\text{SAT}(f(x) \wedge AX g(x, y)) = \{(s_1, [x \mapsto 1, y \mapsto 2])\}$$

Problem : y has to be the same everywhere.

Example

Semantic patch

```
@leaks exists@
expression  $x$ ,  $E$ ;
statement  $S$ ;
constant  $C$ ;
@@
```

```
*  $x = \text{kmalloc}(\dots)$ 
...
if ( $x == \text{NULL}$ )  $S$ 
... when  $!= \text{kfree}(x)$ 
    when  $!= (x = E \mid E = x)$ 
* return (  $-C \mid \text{NULL}$  );
```

C code

```
a = kmalloc (...);
if (!a) {...}
if (x) return -1;
if (y) return -2;
kfree (a);
```

Adding existential quantification (CTL-V)

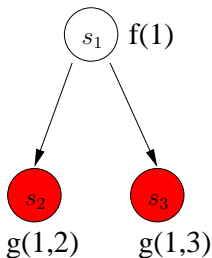
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$\exists x.f(x) \wedge AX \exists y.g(x, y)$



Model checking algorithm

$SAT(g(x, y)) = \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 3])\}$

Adding existential quantification (CTL-V)

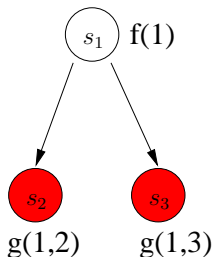
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$\exists x.f(x) \wedge AX \exists y.g(x, y)$



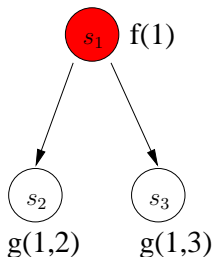
Model checking algorithm

$$\begin{aligned}
 \text{SAT}(g(x, y)) &= \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 3])\} \\
 \text{SAT}(\exists y.g(x, y)) &= \{(s_2, [x \mapsto 1]); (s_3, [x \mapsto 1])\}
 \end{aligned}$$

Adding existential quantification (CTL-V)

Semantic patch

- $f(x)$;
- $g(x, y)$;



CTL representation

$$\exists x. f(x) \wedge AX \exists y. g(x, y)$$

Model checking algorithm

$$\begin{aligned} \text{SAT}(g(x, y)) &= \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 3])\} \\ \text{SAT}(\exists y. g(x, y)) &= \{(s_2, [x \mapsto 1]); (s_3, [x \mapsto 1])\} \\ \text{SAT}(AX \exists y. g(x, y)) &= \{(s_1, [x \mapsto 1])\} \end{aligned}$$

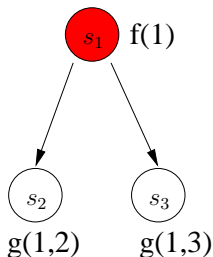
Adding existential quantification (CTL-V)

Semantic patch

- $f(x)$;
- $g(x, y)$;

CTL representation

$$\exists x. f(x) \wedge \text{AX} \exists y. g(x, y)$$



Model checking algorithm

$$\begin{aligned}
 \text{SAT}(g(x, y)) &= \{(s_2, [x \mapsto 1, y \mapsto 2]); (s_3, [x \mapsto 1, y \mapsto 3])\} \\
 \text{SAT}(\exists y. g(x, y)) &= \{(s_2, [x \mapsto 1]); (s_3, [x \mapsto 1])\} \\
 \text{SAT}(\text{AX} \exists y. g(x, y)) &= \{(s_1, [x \mapsto 1])\} \\
 \text{SAT}(f(x) \wedge \text{AX} \exists y. g(x, y)) &= \{(s_1, [x \mapsto 1])\}
 \end{aligned}$$

Adding witnesses (CTL-VW)

Goal : collect information about what and where to transform

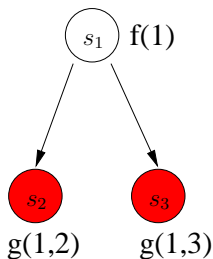
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$\exists x. f(x) \wedge AX \exists y. g(x, y)$



Model checking algorithm

$SAT(g(x, y)) = \{ (s_2, [x \mapsto 1, y \mapsto 2], ()); (s_3, [x \mapsto 1, y \mapsto 3], ()) \}$

Adding witnesses (CTL-VW)

Goal : collect information about what and where to transform

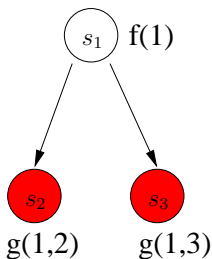
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$$\exists x. f(x) \wedge AX \exists y. g(x, y)$$



Model checking algorithm

$$\begin{aligned} \text{SAT}(g(x, y)) &= \{ (s_2, [x \mapsto 1, y \mapsto 2], ()); (s_3, [x \mapsto 1, y \mapsto 3], ()) \} \\ \text{SAT}(\exists y. g(x, y)) &= \{ (s_2, [x \mapsto 1], (\langle s_2, [y \mapsto 2] \rangle)); \\ &\quad (s_3, [x \mapsto 1], (\langle s_3, [y \mapsto 3] \rangle)) \} \end{aligned}$$

Adding witnesses (CTL-VW)

Goal : collect information about what and where to transform

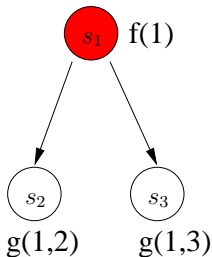
Semantic patch

$f(x)$;

- $g(x, y)$;

CTL representation

$$\exists x. f(x) \wedge AX \exists y. g(x, y)$$



Model checking algorithm

$$\begin{aligned} \text{SAT}(g(x, y)) &= \{ (s_2, [x \mapsto 1, y \mapsto 2], ()); (s_3, [x \mapsto 1, y \mapsto 3], ()) \} \\ \text{SAT}(\exists y. g(x, y)) &= \{ (s_2, [x \mapsto 1], (\langle s_2, [y \mapsto 2] \rangle)); \\ &\quad (s_3, [x \mapsto 1], (\langle s_3, [y \mapsto 3] \rangle)) \} \\ \text{SAT}(AX \exists y. g(x, y)) &= \{ (s_1, [x \mapsto 1], (\langle s_2, [y \mapsto 2] \rangle, \langle s_3, [y \mapsto 3] \rangle)) \} \end{aligned}$$

Outline

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW**
- 5 Algorithme de model checking
- 6 Conclusion

Definition (Syntaxe de CTL-VW)

$$\phi ::= p(\bar{x}) \mid \phi \wedge \phi \mid \phi \vee \phi \mid \neg \phi \mid \exists x. \phi \mid$$
$$AX\phi \mid EX\phi \mid A[\phi U \phi] \mid E[\phi U \phi]$$

Modèle $(S, \rightarrow, \text{Label})$

- S : ensemble d'états
- $\rightarrow \subseteq S \times S$ relation de transition
- $\text{Label} : S \rightarrow \mathcal{P}(\text{Atom})$ fonction de label

Environnement θ

$\theta : \text{MetaVar} \rightarrow (\text{Val} + \mathcal{P}(\text{Val}))$

$$\theta = [x \mapsto 1, y \mapsto \{3, 4\}]$$

Témoin ω , forêt Ω

$\omega \in S \times \text{MetaVar} \times (\text{Val} + \mathcal{P}(\text{Val})) \times \text{ForetTemoins}$

Ω : (multi)ensemble de témoins

$$\Omega = \{\langle s_0, y, 2, \emptyset \rangle, \langle s_3, x, 1, \{\langle s_4, y, 3, \emptyset \rangle\} \rangle\}$$

Definition (Semantique de CTL-VW)

$$\begin{array}{ll}
s \models_{\theta, \emptyset} p(\bar{x}) & \Leftrightarrow p(\theta(\bar{x})) \in \text{Label}(s) \\
s \models_{\theta, \Omega_1 \uplus \Omega_2} \phi_1 \wedge \phi_2 & \Leftrightarrow s \models_{\theta, \Omega_1} \phi_1 \wedge s \models_{\theta, \Omega_2} \phi_2 \\
s \models_{\theta, \Omega} \phi_1 \vee \phi_2 & \Leftrightarrow s \models_{\theta, \Omega} \phi_1 \vee s \models_{\theta, \Omega} \phi_2 \\
s \models_{\theta, \emptyset} \neg \phi & \Leftrightarrow \forall \Omega. s \not\models_{\theta, \Omega} \phi \\
\\
s \models_{\theta, \{\langle s, x, v, \Omega \rangle\}} \exists x. \phi & \Leftrightarrow s \models_{\theta[x \mapsto v], \Omega} \phi \\
\\
s \models_{\theta, \Omega} \text{AX} \phi & \Leftrightarrow \exists (\Omega_{s'})_{s' \in \text{next}(s)}. \\
& \bigoplus_{s' \in \text{next}(s)} \Omega_{s'} = \Omega \wedge \forall s' \in \text{next}(s). s' \models_{\theta, \Omega_{s'}} \phi \\
s \models_{\theta, \Omega} \text{EX} \phi & \Leftrightarrow \exists s' \in \text{next}(s). s' \models_{\theta, \Omega} \phi
\end{array}$$

Definition (AU and EU semantics)

$$\begin{aligned}
 \mathbf{s} \models_{\theta, \Omega} \mathbf{A}[\phi_1 \mathbf{U} \phi_2] &\Leftrightarrow \exists \Sigma \in \Pi(\mathbf{s}). \exists (\Omega_{s'})_{s' \in \bar{\Sigma}}. \\
 \forall \pi \in \Sigma. \pi[|\pi|] \models_{\theta, \Omega_{\pi[|\pi|]}} \phi_2 \wedge \forall 0 \leq j < |\pi|. \pi[j] \models_{\theta, \Omega_{\pi[j]}} \phi_1 \\
 &\wedge \biguplus_{s' \in \bar{\Sigma}} \Omega_{s'} = \Omega
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{s} \models_{\theta, \Omega} \mathbf{E}[\phi_1 \mathbf{U} \phi_2] &\Leftrightarrow \exists \Sigma \in \Pi(\mathbf{s}). \exists \pi \in \Sigma. \exists (\Omega_j)_{j \in 0..|\pi|}. \\
 \pi[|\pi|] \models_{\theta, \Omega_{|\pi|}} \phi_2 \wedge \forall 0 \leq j < |\pi|. \pi[j] \models_{\theta, \Omega_j} \phi_1 \\
 &\wedge \biguplus_{j \leq |\pi|} \Omega_j = \Omega
 \end{aligned}$$

SmPL vers CTL-VW

 $f(E1)$ \dots
 $g(E2)$ $f(E1)$ \wedge $AX(A[\dots$ \cup $g(E2)])$

SmPL vers CTL-VW

 $f(E1)$ \dots
 $g(E2)$

$$\exists E1.(f(E1)$$

$$\wedge$$

$$AX(A[\dots$$

$$\cup$$

$$\exists E2.g(E2)]))$$

SmPL vers CTL-VW

 $f(E1)$ \dots
 $g(E2)$

$$\begin{aligned} & \exists E1.(f(E1) \\ & \quad \wedge \\ & \quad AX(A[\neg(f(E1) \vee \exists E2.g(E2)) \\ & \quad \quad \cup \\ & \quad \quad \exists E2.g(E2)])) \end{aligned}$$

Outline

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking**
- 6 Conclusion

Algorithmes de model checking

Fonction SAT :

$$\text{SAT}(p(\bar{x})) = \{(s, \theta, \emptyset) \mid p(\theta(\bar{x})) \in \text{Label}(s) \wedge \text{dom}(\theta) = \text{fv}(p(\bar{x}))\}$$

$$\text{SAT}(\phi_1 \wedge \phi_2) = \text{conj}(\text{SAT}(\phi_1), \text{SAT}(\phi_2))$$

$$\text{SAT}(\phi_1 \vee \phi_2) = \text{SAT}(\phi_1) \cup \text{SAT}(\phi_2)$$

$$\text{SAT}(\neg\phi) = \text{neg}(\text{SAT}(\phi))$$

$$\text{SAT}(\exists x.\phi) = \{\text{existsone}(x, t) \mid t \in \text{SAT}(\phi)\}$$

$$\text{SAT}(\text{AX } \phi) = \text{pre}_{\forall}(\text{SAT}(\phi))$$

$$\text{SAT}(\text{EX } \phi) = \text{pre}_{\exists}(\text{SAT}(\phi))$$

$$\text{SAT}(\text{A}[\phi_1 \text{ U } \phi_2]) = \text{SAT}_{\text{AU}}(\text{SAT}(\phi_1), \text{SAT}(\phi_2))$$

$$\text{SAT}(\text{E}[\phi_1 \text{ U } \phi_2]) = \text{SAT}_{\text{EU}}(\text{SAT}(\phi_1), \text{SAT}(\phi_2))$$

Algorithme de model checking (2)

Conjonction

$$\text{conj}(T_1, T_2) = \{(s, \theta_1 \sqcap \theta_2, \Omega_1 \uplus \Omega_2) \mid (s, \theta_1, \Omega_1) \in T_1 \wedge (s, \theta_2, \Omega_2) \in T_2\}$$

Quantificateur existentiel

$$\begin{aligned} \text{existsone}(x, (s, \theta, \Omega)) = & \\ & (s, \theta - [x \mapsto v], \{\langle s, x, v, \Omega \rangle\}), \text{ if } \theta^+(x) = v \\ & (s, \theta - [x \not\mapsto V], \{\langle s, x, V, \Omega \rangle\}), \text{ if } \theta^-(x) = V \\ & (s, \theta, \{\langle s, x, \emptyset, \Omega \rangle\}), \text{ otherwise} \end{aligned}$$

Correction et complétude

On voudrait montrer

$$(s, \theta, \Omega) \in \text{SAT}(\phi) \iff s \models_{\theta, \Omega} \phi$$

Difficulté due aux «bindings» négatifs

- Un environnement $[x \mapsto 2]$ «correspond» à plusieurs environnements positifs $[x \mapsto 1], [x \mapsto 3], \dots$
- Un résultat (s, θ, Ω) de l'algorithme de model checking peut «correspondre» à plusieurs triplets pour la sémantique

Correction et complétude (2)

Relation d'ordre

On définit

- Une relation d'ordre sur les «bindings», qui induit une relation d'ordre sur les environnements et sur les témoins.

$$[x \mapsto 2] \sqsubseteq [x \mapsto \{3, 4\}] \sqsubseteq [x \mapsto \{3\}]$$

- Pour $T \subseteq S \times \text{Env} \times \text{WitForest}$, $\eta_\phi(T) =$

$$\{(s, \theta, \Omega) \mid \exists (s', \theta', \Omega') \in T. s = s' \wedge \theta \sqsubseteq \theta' \wedge \Omega \sqsubseteq \Omega'\}$$

Theorem (Correction et complétude)

$$s \models_{\theta, \Omega} \phi \iff (s, \theta, \Omega) \in \eta_\phi(\text{SAT}(\phi))$$

Outline

- 1 Introduction
- 2 Langage de Coccinelle
- 3 Intuition de la sémantique de Coccinelle
- 4 Syntaxe et sémantique de CTL-VW
- 5 Algorithme de model checking
- 6 Conclusion**

Résultats expérimentaux

Démarche de la détection et correction de bugs

- Spécification d'un motif aussi «correct» que possible («correct» = conforme à ce que l'on recherche)
- Incorrect \Rightarrow existence de *faux positifs* (fausses alarmes)
- Le taux de faux positifs est diminué en raffinant le patch sémantique

Noyau Linux

- environ 25 000 fichiers
 - 250 M de code
 - + de 8 millions de lignes de code
- \Rightarrow + de 180 patches acceptés dans le noyau

Benchmarks

Tests sur un 8 processeurs, 16GB de RAM, 3GHz

patch sémantique	temps	patch sémantique	temps
leaks	3 mn 30	categorize functions	11 mn
returns_null	4 mn 20	detect bad call of a function	2 s
usb_submit	1mn 40	correct bad call of a function	2 s
		bad pointer check	9 mn

Perspectives

Perspectives

- Utiliser coccinelle sur du code critique
- Patches et reconfiguration de code à grande échelle
- Tester le respect de standards de codage
- Étendre coccinelle à d'autres langage que C
- Gérer les flots de données

Plus d'informations

Documentation, téléchargement, exemples, etc. sur le site du projet :

`http://www.emn.fr/x-info/coccinelle/`

Participation à Coccinelle, points de contact

- Université de Copenhague (Julia Lawall `julia@diku.dk`)
- ONERA (Julien Brunel `julien.brunel@onera.fr`)
- Université de Aalborg
- École des Mines de Nantes

Références

- Détection et correction de bugs



WYSIWIB : A declarative approach to finding protocols and bugs in Linux code, DSN'09, Lisbonne, Portugal, 2009

- Sémantique de Coccinelle, model checking de CTL-VW



A foundation for flow-based program matching using temporal logic and model checking, POPL'09, Savannah, Géorgie, 2009