

Théorie, conception et réalisation d'un langage de programmation adapté à XML

Alain Frisch

Soutenance de thèse - 13/12/2004

Directeur: Giuseppe Castagna

Plan

Théorie, conception, et réalisation
d'un langage de programmation adapté à XML

Plan

Théorie, conception, et réalisation
d'un langage de programmation adapté à XML

① ...XML.

Plan

Théorie, conception, et réalisation
d'un langage de programmation adapté à XML

- 1 ...XML.
- 2 ...langage de programmation adapté ...

Plan

Théorie, **conception**, et réalisation
d'un **langage de programmation adapté** à XML

- 1 ...XML.
- 2 ...langage de programmation adapté ...
- 3 ...conception ...

Plan

Théorie, conception, et réalisation d'un langage de programmation adapté à XML

- 1 ...XML.
- 2 ...langage de programmation adapté ...
- 3 ...conception ...
- 4 Théorie ...

Plan

Théorie, conception, et réalisation d'un langage de programmation adapté à XML

- 1 ...XML.
- 2 ...langage de programmation adapté ...
- 3 ...conception ...
- 4 Théorie ...
- 5 ...réalisation ...

Plan

- 1 **Introduction**
 - Contexte, motivations, objectif initial
- 2 CDuce : aperçu du langage
- 3 Contributions théoriques
 - Algèbre de types et sous-typage
 - Noyau fonctionnel, système de types
 - Aspects algorithmiques
- 4 Conclusion

XML et schémas

XML

- Un langage de **balises**.
- Pour représenter des données de nature **arborescente**.
- Représentation indépendante de l'application.

Schémas

- Chaque application définit des **contraintes** sur les documents XML qu'elle peut manipuler :
 - balises (nom, structure d'imbrication),
 - attributs (présence, valeurs autorisées),
 - texte (présence).
- Un tel ensemble de contraintes est un **schéma** XML.
- Il existe des langages pour écrire des schémas de manière formelle : DTD, XML-Schema, Relax-NG.

XML : schémas

Transformation XML

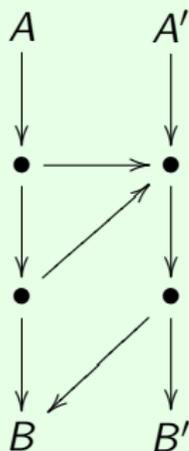
Schéma A $\xrightarrow{\text{Programme T}}$ Schéma B

Garanties

- Le programme T s'engage à fournir un document XML de schéma B si on lui donne un document XML de schéma A.
- **Peut-on le croire ?**

XML : schémas

Application complexe



Toutes les opérations effectuées sur les documents à l'intérieur de l'application sont-elles légales ?

Opérations illégales

- Consulter un attribut qui peut ne pas exister.
- Supprimer un élément qui doit être présent.

XML : schémas et types

- Typage dans les langages de programmation : interdire des opérations illégales
(`1 + "Hello"`, `customer.shutdown()`).
- Il est naturel de voir les schémas XML comme des **types** de données, à l'intérieur des applications.

XML		langages
schémas	↔	types
documents	↔	valeurs

xDuce (Hosoya, Vouillon, Pierce)

Pour le programmeur

- Un langage pour la transformation de documents XML.
- Types \simeq DTD. Valeurs = arbres XML.
- Style fonctionnel. Fonctions récursives et filtrage.

Théorie

- Types = automates d'arbres \rightsquigarrow langages réguliers d'arbres.
- Sous-typage = inclusion des langages.

Regular Expression Types for XML (ICFP 2000).

Regular Expression Pattern Matching for XML (POPL 2001).

Sous-typage

Sous-typage dans XDuce

- Définition sémantique du sous-typage. Pourtant :
 - Système de types : dépend du sous-typage.
 - Sémantique du langage : dirigée par les types (filtrage).

Paradoxe ?

Sous-typage

Sous-typage dans XDuce

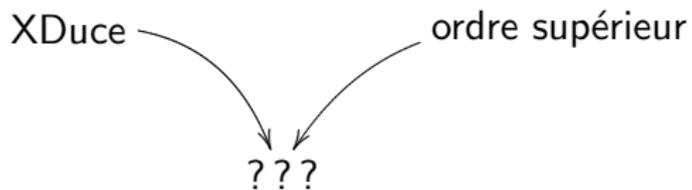
- Définition sémantique du sous-typage. Pourtant :
 - Système de types : dépend du sous-typage.
 - Sémantique du langage : dirigée par les types (filtrage).

Paradoxe ?

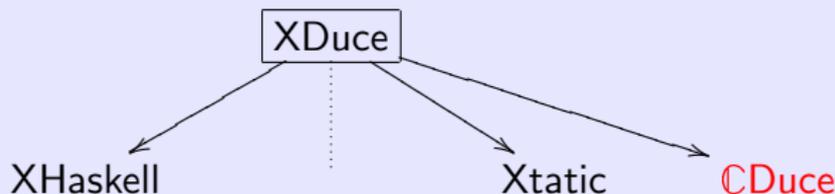
Non !

- La sémantique des types ne dépend pas du langage.

XDuce + ordre supérieur



XDuce : descendance



- XHaskell : types XDuce dans Haskell (sous-typage encodé dans les *type classes*); pas de filtrage XML, pas de mélange entre XML et fonctions, pas de sous-typage sémantique. (Zhuo Ming Lu, Sulzmann.)
- Xtatic : types XDuce dans C# (sous-typage par nom pour les objets); filtrage limité. Perte du côté structurel des types. (Pierce, Schmitt, Gapeyev, Levin.)

Objectif initial de la thèse

Étendre l'approche ensembliste de XDuce à un calcul d'ordre supérieur, qui permet de mélanger librement fonctions et documents XML, dans un cadre structurel.

Coté langage

- XDuce est un langage fonctionnel, mais sans ordre supérieur : ça manque !
- Plus spécialement pour XML : transformations paramétrables de manière modulaire, systèmes de *templates* de première classe, ...

Coté λ -calcul

Prospectif

Objectif initial de la thèse

Étendre l'approche ensembliste de X_Duce à un calcul d'ordre supérieur, qui permet de mélanger librement fonctions et documents XML, dans un cadre structurel.

Coté langage

Coté λ -calcul

- Défi intellectuel : sous-typage ensembliste avec types récurifs, combinaisons booléennes, constructeurs de types.
- Sous-typage ensembliste ($t \leq s \iff \llbracket t \rrbracket \subseteq \llbracket s \rrbracket$) :
 - Facile à expliquer. Exhiber valeur pour illustrer $t \not\leq s$.
 - Définition conceptuellement simple, indépendante d'un algorithme ou d'un système axiomatique complexe.

Prospectif

Objectif initial de la thèse

Étendre l'approche ensembliste de XDuce à un calcul d'ordre supérieur, qui permet de mélanger librement fonctions et documents XML, dans un cadre structurel.

Coté langage

Coté λ -calcul

Prospectif

- Interaction données/comportements sur le web (formulaires, javascript, applets, servlets, . . .)

Plan

- 1 Introduction
 - Contexte, motivations, objectif initial
- 2 CDuce : aperçu du langage**
- 3 Contributions théoriques
 - Algèbre de types et sous-typage
 - Noyau fonctionnel, système de types
 - Aspects algorithmiques
- 4 Conclusion

Aperçu de CDuce

Types XML

```
type Bib = [ Book* ]  
type Book = <book>[ Title Subtitle? Author+ ]  
type Title = <title>[ PCDATA ]  
type Subtitle = <title>[ PCDATA ]  
type Author = <author>[ PCDATA ]
```

Fonctions et filtrage

```
let title(Book -> String) <book>[ <title>x_* ] -> x
```

Aperçu de CDuce

Types XML

```
type Bib = [ Book* ]
type Book = <book>[ Title Subtitle? Author+ ]
type Title = <title>[ PCDATA ]
type Subtitle = <title>[ PCDATA ]
type Author = <author>[ PCDATA ]
```

Fonctions et filtrage

```
let title(Book -> String) <book>[ <title>x_* ] -> x

let author(Book -> [ Author+ ]) <book>[ (x::Author | _) * ] -> x
```

Aperçu de CDuce

Ordre supérieur

```
type FBook = Book -> String
type ABook = <book print=FBook>[ Title Subtitle? Author+ ]
type ABib = [ ABook* ]
Remarque: ABook ≤ Book   ABib ≤ Bib
```

Aperçu de CDuce

Ordre supérieur

```
type FBook = Book -> String
type ABook = <book print=FBook>[ Title Subtitle? Author+ ]
type ABib = [ ABook* ]
  Remarque: ABook ≤ Book   ABib ≤ Bib

let set(<book>c : Book)(f : FBook) : ABook = <book print=f>c
```

Aperçu de CDuce

Ordre supérieur

```
type FBook = Book -> String
type ABook = <book print=FBook>[ Title Subtitle? Author+ ]
type ABib = [ ABook* ]
  Remarque: ABook ≤ Book   ABib ≤ Bib

let set(<book>c : Book)(f : FBook) : ABook = <book print=f>c

let prepare(b : Bib) : ABib = map b with x -> set x title
```

Aperçu de C Duce

Ordre supérieur

```

type FBook = Book -> String
type ABook = <book print=FBook>[ Title Subtitle? Author+ ]
type ABib = [ ABook* ]
  Remarque: ABook ≤ Book   ABib ≤ Bib

let set(<book>c : Book)(f : FBook) : ABook = <book print=f>c

let prepare(b : Bib) : ABib = map b with x -> set x title

type Ul = <ul>[ Li+ ]
type Li = <li>[ PCDATA ]

let display (ABib -> Ul; ABook -> Li)
| <book print=f>_ & x -> <li>(f x)
| [] -> raise "Empty bibliography"
| p -> <ul>(map p with z -> display z)

```

Aperçu de CDuce

Ordre supérieur

```
let change(p : Book -> Bool)(f : FBook)(b : ABib) : ABib =  
  map b with x -> if (p x) then set x f else x
```

Aperçu de CDuce

Ordre supérieur

```
let change(p : Book -> Bool)(f : FBook)(b : ABib) : ABib =  
  map b with x -> if (p x) then set x f else x
```

```
type HasSub = <_>[ *_ Subtitle *_ ]
```

```
let change_if_sub =  
  change (fun (Book -> Bool) HasSub -> 'true | _ -> 'false)
```

```
let subtitle_first (Bib -> Bib; ABib -> ABib)  
  [ (x::HasSub | y::_)* ] -> x @ y
```

Plan

- 1 Introduction
 - Contexte, motivations, objectif initial
- 2 CDuce : aperçu du langage
- 3 Contributions théoriques**
 - Algèbre de types et sous-typage
 - Noyau fonctionnel, système de types
 - Aspects algorithmiques
- 4 Conclusion

XDuce revisité

Encodages des valeurs et des types XML

- Séquences.

$$[v_1 \dots v_n] \equiv (v_1, (\dots, (v_n, 'nil) \dots))$$

$$[t_1 * t_2?] \equiv \mu\alpha. (t_1 \times \alpha) \vee (t_2 \times 'nil) \vee 'nil$$

- Éléments XML.

$$\langle \text{tag} \rangle [v_1 \dots v_n] \equiv (\text{tag}, [v_1, \dots, v_n])$$

Des types bien connus

Types : produits, atomiques, réunions, récursifs.

Vers CDuce

Encodages des valeurs et des types XML

- Séquences.

$$[v_1 \dots v_n] \equiv (v_1, (\dots, (v_n, 'nil) \dots))$$

$$[t_1 * t_2?] \equiv \mu\alpha. (t_1 \times \alpha) \vee (t_2 \times 'nil) \vee 'nil$$

- Éléments XML.

$$\langle \text{tag} \rangle [v_1 \dots v_n] \equiv (\text{tag}, [v_1, \dots, v_n])$$

- Constructeurs : produit, **flèche**, enregistrement.
- Connecteurs booléens : réunion, intersection, différence.
- Types récursifs.

Algèbre de types

Types mal formés : à éviter

$$t = \mu\alpha.\alpha$$

$$t = \mu\alpha.\alpha \vee \alpha$$

$$t = \mu\alpha.\neg\alpha$$

Stratification

$$\mathbb{T} = (T, \tau) \quad \text{avec } \tau : T \longrightarrow F(T)$$

$$C \in F(T) := b \mid t \times t \mid t \rightarrow t \mid C \vee C \mid C \wedge C \mid \neg C \mid 0 \mid 1$$

Deux conditions :

- Décomposition finie des types
- Existence des types rékursifs

Algèbre de types : stratification

Tautologies booléennes

- Liberté pour quotienter modulo tautologies booléennes.
- $CVC \simeq C$
 $C \wedge \neg C \simeq 0$

Partage

- Liberté pour partager ou non types modulo déroulement.
- $\mu\alpha. (\alpha \times \alpha) \forall t \stackrel{?}{=} \mu\alpha. t \forall (\alpha \times (\mu\beta. \beta \times \alpha \forall t) \forall \alpha \times \alpha)$

Algèbre de types

Résultats

- Cadre catégorique : souplesse d'implémentation ; naturalité des algorithmes.
- Constructions de plusieurs algèbres, avec partage minimal, maximal, intermédiaire, et caractérisations.
- Représentation des combinaisons booléennes par arbres de décision ternaire.

Sous-typage

Définition ensembliste

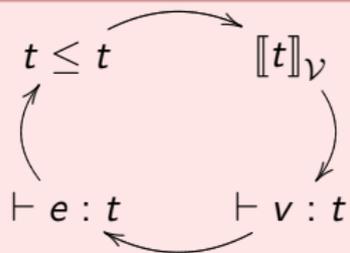
$$t \leq s \iff \llbracket t \rrbracket_{\mathcal{V}} \subseteq \llbracket s \rrbracket_{\mathcal{V}} \quad \text{où} \quad \llbracket t \rrbracket_{\mathcal{V}} = \{v \mid \vdash v : t\}$$

Sous-typage

Définition ensembliste

$$t \leq s \iff \llbracket t \rrbracket_{\mathcal{V}} \subseteq \llbracket s \rrbracket_{\mathcal{V}} \quad \text{où} \quad \llbracket t \rrbracket_{\mathcal{V}} = \{v \mid \vdash v : t\}$$

Circularité !



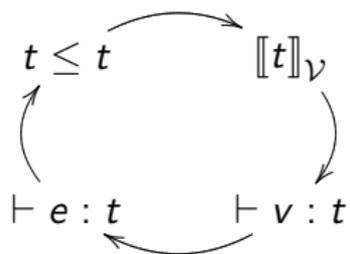
Circularité : comment la dépasser ?



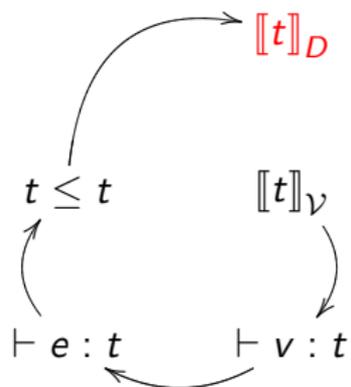
Circularité : comment la dépasser ?



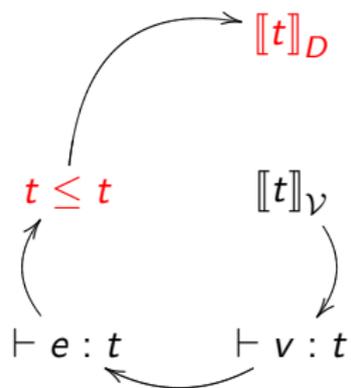
Amorçage



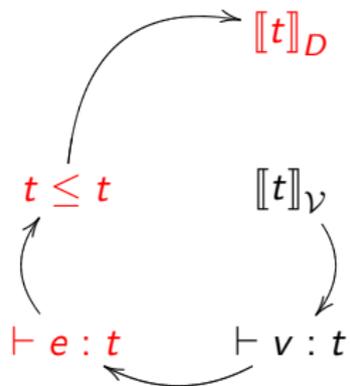
Amorçage



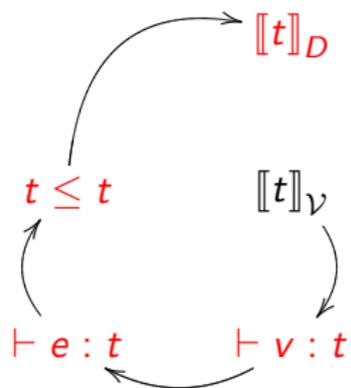
Amorçage



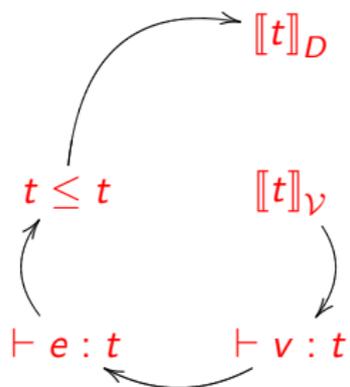
Amorçage



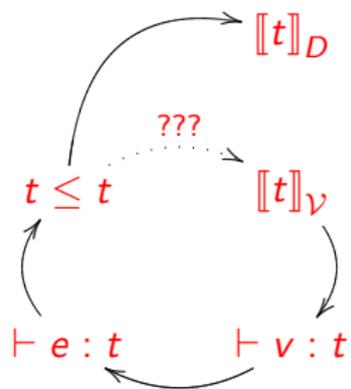
Amorçage



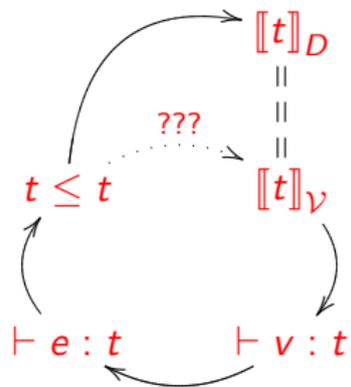
Amorçage



Amorçage



Amorçage



Modèles

Interprétation ensembliste : $\llbracket _ \rrbracket : T \rightarrow \mathcal{P}(D)$

$$\llbracket t_1 \vee t_2 \rrbracket = \llbracket t_1 \rrbracket \cup \llbracket t_2 \rrbracket$$

$$\llbracket t_1 \wedge t_2 \rrbracket = \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket$$

$$\llbracket \neg t \rrbracket = D \setminus \llbracket t \rrbracket$$

$$\llbracket 0 \rrbracket = \emptyset$$

$$\llbracket 1 \rrbracket = D$$

Modèles

Interprétation ensembliste : $\llbracket _ \rrbracket : T \rightarrow \mathcal{P}(D)$

$$\llbracket t_1 \vee t_2 \rrbracket = \llbracket t_1 \rrbracket \cup \llbracket t_2 \rrbracket$$

$$\llbracket t_1 \wedge t_2 \rrbracket = \llbracket t_1 \rrbracket \cap \llbracket t_2 \rrbracket$$

$$\llbracket \neg t \rrbracket = D \setminus \llbracket t \rrbracket$$

$$\llbracket 0 \rrbracket = \emptyset$$

$$\llbracket 1 \rrbracket = D$$

Constructeurs : $\mathbb{E}[_] : T \rightarrow \mathcal{P}(ED)$

$$\mathbb{E}D = \mathcal{C} + D \times D + D^D$$

$$\mathbb{E}\llbracket t_1 \times t_2 \rrbracket = \llbracket t_1 \rrbracket \times \llbracket t_2 \rrbracket \subseteq D \times D$$

$$\mathbb{E}\llbracket t_1 \rightarrow t_2 \rrbracket = \llbracket t_2 \rrbracket^{\llbracket t_1 \rrbracket} \subseteq D^D$$

$$\mathbb{E}\llbracket t_1 \vee t_2 \rrbracket = \mathbb{E}\llbracket t_1 \rrbracket \cup \mathbb{E}\llbracket t_2 \rrbracket \dots$$

Modèles

Peut-on avoir $D = \mathbb{E}D = \mathcal{C} + D \times D + D^D$?

Modèles

Peut-on avoir $D = \mathbb{E}D = \mathcal{C} + D \times D + D^D$?

Non !

- Pour une simple raison de cardinalité.
- On veut *in fine* que l'ensemble des valeurs soit un modèle.

Modèles

Peut-on avoir $D = \mathbb{E}D = \mathcal{C} + D \times D + D^D$?

Non !

- Pour une simple raison de cardinalité.
- On veut *in fine* que l'ensemble des valeurs soit un modèle.

Faisons « comme si »

Définition : $(D, \llbracket - \rrbracket)$ est un modèle si :

$$\forall t, s. \llbracket t \rrbracket \subseteq \llbracket s \rrbracket \iff \mathbb{E}\llbracket t \rrbracket \subseteq \mathbb{E}\llbracket s \rrbracket$$

(ou : $\forall t. \llbracket t \rrbracket = \emptyset \iff \mathbb{E}\llbracket t \rrbracket = \emptyset$)

Modèles

Résultats

- Caractérisation axiomatique des modèles en terme de simulation (notion coinductive).
- Existence d'un modèle universel $\llbracket _ \rrbracket^0$, qui induit la plus grande relation de sous-typage.
- Existence de modèles non-universels (plus difficile!).

$$t = \mu\alpha.(\mathbb{0} \rightarrow \mathbb{0}) \setminus (\alpha \rightarrow \mathbb{0}) \quad \llbracket t \rrbracket = \emptyset?$$

Raisonnements dans les modèles

- Propriétés simples : transitivité de \leq , variance des constructeurs, tautologies booléennes, égalités ensemblistes :
 $(t_1 \times s_1) \wedge (t_2 \times s_2) \simeq (t_1 \wedge t_2) \times (s_1 \wedge s_2)$
 $(t \rightarrow s_1) \wedge (t \rightarrow s_2) \simeq t \rightarrow (s_1 \wedge s_2)$
 $(t_1 \rightarrow s) \wedge (t_2 \rightarrow s) \simeq (t_1 \vee t_2) \rightarrow s$
- Très utile pour préparer l'étude algorithmique du sous-typage et pour mener l'étude méta-théorique du système de types !

Sous-typage : décomposition

$$\bigwedge_{i \in I} t_i \times s_i \leq \bigvee_{i \in J} t_i \times s_i$$
$$\iff \forall J' \subseteq J. \left(\bigwedge_{i \in I} t_i \leq \bigvee_{i \in J'} t_i \right) \text{ ou } \left(\bigwedge_{i \in I} s_i \leq \bigvee_{i \in J \setminus J'} s_i \right)$$

$$\bigwedge_{i \in I} t_i \rightarrow s_i \leq \bigvee_{i \in J} t_i \rightarrow s_i$$
$$\iff \exists j \in J. \forall I' \subseteq I. \left(t_j \leq \bigvee_{i \in I'} t_i \right) \text{ ou } \left(I' \neq I \text{ et } \bigwedge_{i \in I \setminus I'} s_i \leq s_j \right)$$

Coller au langage

$$(t_1 \rightarrow s_1) \wedge (t_2 \rightarrow s_2) \stackrel{?}{\leq} (t_1 \vee t_2) \rightarrow (s_1 \wedge s_2)$$

$$\text{int} \rightarrow \text{int} \stackrel{?}{\leq} \mathbb{1} \rightarrow \mathbb{1}$$

Résultat : variantes du système

- Avec/sans fonctions surchargées.
- Avec/sans erreur de type à l'application.

Plan

- 1 Introduction
 - Contexte, motivations, objectif initial
- 2 CDuce : aperçu du langage
- 3 Contributions théoriques**
 - Algèbre de types et sous-typage
 - **Noyau fonctionnel, système de types**
 - Aspects algorithmiques
- 4 Conclusion

Noyau fonctionnel, système de types

$$\frac{\Gamma \vdash e : t_1 \quad t_1 \leq t_2}{\Gamma \vdash e : t_2}$$

$$\overline{\Gamma \vdash c : b_c} \quad \overline{\Gamma \vdash x : \Gamma(x)}$$

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash (e_1, e_2) : t_1 \times t_2} \quad \frac{\Gamma \vdash e_1 : t \rightarrow s \quad \Gamma \vdash e_2 : t}{\Gamma \vdash e_1 e_2 : s}$$

$$\frac{\forall i = 1..n. (x : t_i), \Gamma \vdash e : s_i}{\Gamma \vdash \lambda^{t_1 \rightarrow s_1; \dots; t_n \rightarrow s_n} x. e : \bigwedge_{i=1..n} (t_i \rightarrow s_i)}$$

Noyau fonctionnel, système de types

$$\frac{\Gamma \vdash e : t_0 \quad (x : t_0 \wedge t), \Gamma \vdash e_1 : s_1 \quad (x : t_0 \setminus t), \Gamma \vdash e_2 : s_2}{\Gamma \vdash (x = e \in t ? e_1 | e_2) : s_1 \vee s_2}$$

Système de types

Résultats

- **Sûreté** pour une sémantique à petits pas.
- On obtient un **modèle** en prenant $\llbracket t \rrbracket_{\nu} = \{v \mid \vdash v : t\}$. Il est équivalent au modèle d'amorce (pour n'importe quel choix du modèle d'amorce). **La boucle est bouclée !**
- **Exactitude** locale du typage de l'application et du filtrage.
- Pas de type « principal », mais **algorithme d'inférence**.
- Itérateurs génériques, typage par déroulement des types.

Plan

- 1 Introduction
 - Contexte, motivations, objectif initial
- 2 CDuce : aperçu du langage
- 3 Contributions théoriques**
 - Algèbre de types et sous-typage
 - Noyau fonctionnel, système de types
 - **Aspects algorithmiques**
- 4 Conclusion

Calcul du sous-typage

Pour un modèle universel, le prédicat unaire $P(t) \equiv \llbracket t \rrbracket \neq \emptyset$ est un prédicat **inductif** (mais pas sur la structure de t qui est cyclique!).

Modularisation

- 1 Reformuler $P(t)$ comme plus petit point fixe d'un opérateur monotone (contraintes booléennes).
- 2 Calculer ce plus petit point fixe.

Résultats

- 1 Divers optimisations (\Leftarrow raisonnements ensemblistes).
- 2 Solveur *top-down* sans *back-tracking*, avec implémentation légère.

Compilation du filtrage

Un exemple

```
type A = <a>[ A* ]  
type B = <b>[ B* ]  
  
fun (A|B -> Int) A -> 0 | B -> 1  
≈  
fun (A|B -> Int) <a>_ -> 0 | _ -> 1
```

- Tirer partie des informations de **types statiques**
- Autoriser un **style plus déclaratif** (\rightsquigarrow robuste)

Résultat

- Cadre pour raisonner sur les stratégies de compilation.
- Critère de correction des optimisations.
- Heuristique de compilation.

Principales contributions

Aspects sémantiques

- Extension de XDuce avec l'**ordre supérieur**, en préservant l'approche ensembliste.
- Calcul de **fonctions surchargées**. Variante sans surcharge.
- Formulation originale pour les types et motifs récursifs.
- Filtrage : motifs intersection, capture plus puissante, **typage exact**.
- **Enregistrements** extensibles et sous-typage ensembliste.

Principales contributions

Aspects algorithmiques

- **Compilation optimisée** du filtrage.
- Modularisation de l'**algorithme de sous-typage** et calcul efficace.
- Implémentation efficace : **représentation des valeurs**.

Implémentation de CDuce

- \simeq 20 000 lignes de code (OCaml).
- Distribution publique depuis juin 2003.
- Implémentation assez efficace (ordre de grandeur des moteurs XSLT et XQuery pour transformations/requêtes équivalentes).
- Prise en compte de tout XML, Unicode, Namespaces.
Traduction DTD \rightsquigarrow types, et partiellement pour XML Schema.
- Quelques dizaines d'utilisateurs.
- Utilisations : sites web statiques, dynamiques (y compris en production), enseignement.
- Système d'interface typé avec OCaml :
 - Utiliser des bibliothèques OCaml existantes.
 - Appeler CDuce depuis OCaml (projets mixtes).

Thèses autour de CDuce

- Analyse de sécurité (M. Burelle)
- Langage de requêtes (C. Miachon)
- Itérateurs, chemins, ... (K. Nguyen)

Travaux futurs

Sémantique, typage

- Polymorphisme paramétrique.
- Extension de ML avec types et opérations XML.
- Inférence des types de fonctions.
- Logique plus puissante (ex : contraintes d'arité, pointeurs ID/IDREF).

Travaux futurs

Algorithmique, compilation

- Représentation des valeurs dirigée par le type et le contexte d'utilisation.
- Implémentation du partage optimal des types rékursifs modulo déroulement + tautologies booléennes.
- Modèle de coût pour la compilation du filtrage.

Merci !

<http://www.cduce.org/>