

Analyse de Données Structurées- Cours 2

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

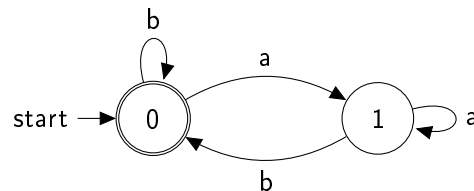
treinen@pps.univ-paris-diderot.fr

28 janvier 2015

© Ralf Treinen 2015

Exemple($\Sigma = \{a, b\}$)

- ▶ (Q, q_0, F, Δ) où
 - | $Q = \{0, 1\}$
 - | $q_0 = 0$
 - | $F = \{0\}$
 - | $\Delta(0, a) = 1 \quad \Delta(0, b) = 0$
 - | $\Delta(1, a) = 1 \quad \Delta(1, b) = 0$
- ▶ Représentation graphique:



Automates Déterministes

- ▶ Étant donné un alphabet fini Σ .
- ▶ Un automate est un tuple (Q, q_0, F, Δ) où
 - | Q est un ensemble fini d'états
 - | $q_0 \in Q$ est l'état initial
 - | $F \subseteq Q$ est l'ensemble d'états acceptants
 - | $\Delta: Q \times \Sigma \rightarrow Q$ la fonction de transition.

Exécutions sur un mot

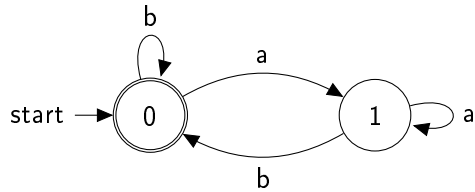
- ▶ Définition d'une fonction d'exécution $\Delta^*: Q \times \Sigma^* \rightarrow Q$ par récurrence:
 - | $\Delta^*(q, \epsilon) = q$
 - | $\Delta^*(q, xw) = \Delta^*(\Delta(q, x), w)$ où $x \in \Sigma$
- ▶ Un automate $A = (Q, q_0, F, \Delta)$ **accepte** le mot $w \in \Sigma^*$ ssi $\Delta^*(q_0, w) \in F$.
- ▶ $\mathcal{L}(A) = \{w \in \Sigma^* \mid \Delta^*(q_0, w) \in F\}$
- ▶ Un langage L est **régulier** quand il existe un automate A tel que $L = \mathcal{L}(A)$.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

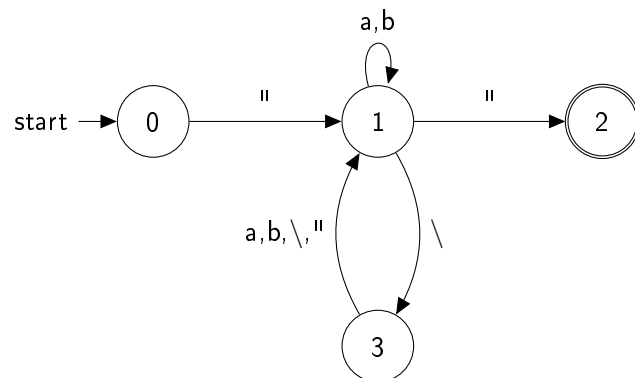
Sur l'exemple



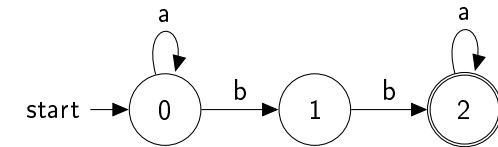
- ▶ $\Delta^*(q_0, bbab) = 0 \in F$: mot accepté
- ▶ $\Delta^*(q_0, bababa) = 1 \notin F$: mot pas accepté
- ▶ $\mathcal{L}(A) = \{\epsilon\} \cup \{wb \mid w \in \Sigma^*\}$

Exemple : un automate pour les strings de Java

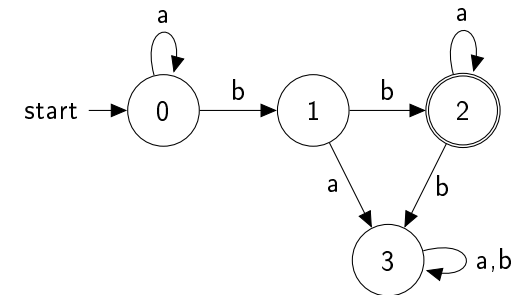
- ▶ Simplification : $\Sigma = \{a, b, \backslash, \text{"}\}$
- ▶ Expression rationnelle vue au cours 1 : $(a \mid b \mid \backslash \mid \text{"})^*$
- ▶ Automate :

Raccourci : Automate avec une fonction Δ partielle

▶



▶ Raccourci pour :

Vu dans le cours *Langages et Automates* du S3

- ▶ Automates *non*-déterministes
- ▶ Automates avec ϵ -transitions
- ▶ Algorithme de *déterminisation*
- ▶ Algorithme de *minimisation*
- ▶ Traduction d'un automate en expression régulière
- ▶ Traduction d'une expression régulière en automate
- ▶ Un langage est rationnel si et seulement s'il est régulier.

Une classe Java très simple pour un automate donné

- ▶ Représenter les états par des entiers $0, 1, \dots, nbe - 1$
- ▶ Représenter l'ensemble des états acceptants par un vecteur de booléens
- ▶ Représenter les symboles de l'alphabet par des entiers $0, \dots, nbs - 1$
- ▶ Représenter la fonction de transition par un tableau à deux dimensions.
- ▶ On utilise la valeur -1 pour l'état poubelle implicite.

Implémentation naïve II

```

    }
}

public static void main(String[] args) {
    int q=q0;
    for(int i=0; i < args[0].length(); i++) {
        q=delta[q][int_of_char(args[0].charAt(i))];
        if (q==-1) {break;}
    }
    if ( q > -1 && accepting[q] ) {
        System.out.println("accepted");
    } else {
        System.out.println("not_accepted");
    }
}

```

Implémentation naïve I

```

public class Abba {

    static int nbe=3;    /* nombre d'etats */
    static int nbs=2;    /* nombre de symboles */
    static int q0=0;     /* etat initial */
    static boolean[] accepting = {false, false, true};
    static int[][] delta = {{0,1}, {-1,2}, {2,-1}};

    public static int int_of_char(char c) {
        if (c=='a') {
            return 0;
        } else {
            return 1;
        }
    }
}

```

Implémentation naïve III

```

    }
}
}

```

Un problème en pratique

- ▶ La taille du tableau des transitions est le nombre d'états \times le nombre de symboles.
- ▶ Gourmand en mémoire quand l'alphabet est très grand (Unicode : $> 1.000.000$ caractères).
- ▶ En pratique il y a très peu de flèches qui partent d'un état, toutes étiquetées par une classe de symboles.
- ▶ Il convient de trouver une implémentation plus efficace en mémoire, à l'aide des classes de caractères.

L'algorithme naïf I

```
public class GrepNaif {  
    public static void main(String[] args) {  
        String motif = args[0];  
        String texte = args[1];  
        int r = motif.length();  
        int n = texte.length();  
        boolean found = false;  
  
        for(int i=0; i<n-r+1; i++) {  
            for(int j=0; j<r; j++) {  
                if (texte.charAt(i+j) != motif.charAt(j)) {  
                    break;  
                }  
                if (j==r-1) {  
                    found=true;  
                }  
            }  
        }  
    }  
}
```

Parenthèse : recherche d'un mot dans un texte

- ▶ On veut savoir si un mot de longueur r (le motif) apparaît dans un texte de longueur n .
- ▶ Cas d'application typique : $n \gg r$.
- ▶ L'algorithme naïf qu'on a vu en IF1/IP1 a une complexité $n * r$.

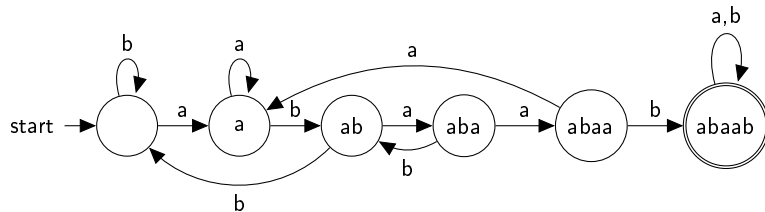
L'algorithme naïf II

```
        }  
    }  
    if (found) {  
        break;  
    }  
}  
  
if (found) {  
    System.out.println("accepted");  
} else {  
    System.out.println("not accepted");  
}  
}
```

Recherche d'un mot dans un texte : l'idée

- ▶ Les automates permettent un algorithme plus efficace !
- ▶ Idée : on construit l'automate qui correspond à l'expression régulière $\Sigma^* r \Sigma^*$.
- ▶ Cet automate s'exécute sur le texte avec une complexité de n : Chaque caractère du texte est traité une seule fois.
- ▶ Construction de l'automate : on peut utiliser la traduction générique des expressions régulières vers des automates ...
- ▶ ... ou faire une construction directe.

L'automate pour $\Sigma^* abaab \Sigma^*$



La construction directe de l'automate

- ▶ Les états sont les *préfixes* du motif r .
- ▶ Exemple : $r = abaab$.
Préfixes : $\epsilon, a, ab, aba, abaa, abaab$.
- ▶ Signification : l'état signifie le *préfixe le plus long* de r qui est un *suffixe* de la partie du texte vue.
- ▶ État initial : ϵ .
- ▶ États acceptants : $\{r\}$
- ▶ Construction de Δ ? Voir l'exemple.

Construction de l'automate

- ▶ Un algorithme efficace pour la construction de l'automate, pour un motif r quelconque ?
- ▶ C'est le célèbre algorithme de *Knuth, Morris, et Pratt*.
- ▶ Complexité linéaire dans la taille du motif.
- ▶ Conséquence : Complexité $n + r$ pour la recherche d'un motif dans un texte (au lieu de $n * r$ avec l'algorithme naïf).
- ▶ Voir un cours d'Algorithmique.

L'objectif de l'analyse lexicale

- ▶ Découper un texte d'entrée en une séquence de *lexèmes*, et les représenter par des *jetons* (*tokens* en anglais)
- ▶ À la base : Classification des lexèmes qui peuvent paraître dans un texte d'entrée, à l'aide des expressions régulières.
- ▶ La phase suivante de l'analyse (l'analyse syntaxique, voir plus tard) va travailler sur le résultat de ce découpage : il s'agit d'une *abstraction* du texte d'entrée.

Définition des catégories lexicales

Sur l'exemple des expressions arithmétiques (on utilise \ comme symbole d'échappement) :

- ▶ INT : `[0..9]+(e[0..9]+)?`
- ▶ IDENT : `[a..zA..Z][a..zA..Z0..9]*`
- ▶ PARG : `\(`
- ▶ PARD : `\)`
- ▶ MULT : `*`
- ▶ PLUS : `\+`

Exemple

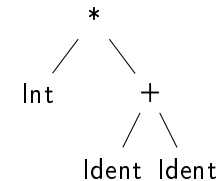
- ▶ Texte d'entrée :

(7	5	6	e	2		*		(e	5	e	7		+		v	a	l	e	u	r	2))
---	---	---	---	---	---	--	---	--	---	---	---	---	---	--	---	--	---	---	---	---	---	---	---	---	---

- ▶ Jetons :

PARG INT MULT PARG IDENT PLUS IDENT PARD PARD

- ▶ Résultat de l'analyse syntaxique (voir plus tard) :



Jetons avec arguments

- ▶ En réalité, on veut aussi garder certaines informations avec les jetons, comme la valeur d'une constante entière, ou le nom d'un identificateur.
- ▶ Certains jetons doivent donc avoir un argument :
 - └ IDENT(string)
 - └ INT(int) (c'est bien int et pas string!)
- ▶ Séquence des jetons obtenue sur l'exemple :
PARG INT(75600) MULT PARG IDENT("e5e7") PLUS IDENT("valeur2") PARD PARD

Ignorer des informations pas pertinentes

L'analyse lexicale sert aussi à faire abstraction de certaines informations dans le texte d'entrée qui ne sont pas pertinentes pour l'analyse du texte. Souvent il s'agit de :

- ▶ Les espaces : sont utiles pour indiquer la fin d'un mot. Les espaces sont utiles *pour* l'analyse lexicale, mais une fois le découpage fait on peut les oublier.
- ▶ Les commentaires : souvent l'analyse lexicale vérifie l'écriture correcte des commentaires, mais ne les représente pas dans sa sortie.

Quelle information retenir dans les jetons

- ▶ On retient dans les jetons seulement l'information qui est utile pour la suite.
- ▶ La distinction entre information utile/inutile dépend de l'application.
- ▶ Par exemple : Les commentaires peuvent être utiles à retenir pour certaines applications.
- ▶ Il peut être utile de conserver avec les jetons aussi des informations de *localisation* : nom du fichier source, numéro de ligne, numéro de colonne.

Exemple

Différents textes d'entrée qui peuvent donner la même séquence de jetons :

- ▶ `34 * (x + y)`
- ▶ `34*(x+y)`
- ▶ `34 *(x+ y)`
- ▶ `34 * (x+y) (* Ceci est un comment ire *)`

Résoudre les ambiguïtés

- ▶ L'analyse lexicale va, pour produire le jeton suivant, chercher un *préfixe* du reste du texte qui correspond à une des catégories lexicales, et construire le jeton correspondant.
- ▶ Il y a deux sources d'ambiguïtés :
 - ┆ Des préfixes de longueurs différents peuvent être reconnues
 - ┆ Les expressions régulières peuvent avoir une intersection non vide

Préfixes de longueur différentes reconnues

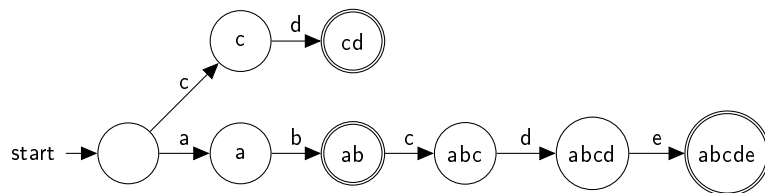
Exemple :

- ▶ Catégorie lexicale :
 - └ IDENT : [a..z]+
- ▶ Début du texte d'entrée :
xyz
- ▶ Plusieurs possibilités de découpage :
 1. IDENT("x") IDENT("y") IDENT("z")
 2. IDENT("xy") IDENT("z")
 3. IDENT("x") IDENT("yz")
 4. IDENT("xyz")
- ▶ La règle normale est : on cherche le préfixe *maximal*.

Exécution de l'automate pour le recherche d'un préfixe maximal

Exemple (artificiel) :

- ▶ Catégorie lexicales :
 - └ AB : ab
 - └ CD : cd
 - └ ABCDE : abcde
- ▶ Automate :



Plusieurs expressions régulières s'appliquent

Exemple :

- ▶ Catégories lexicales :
 - └ PUBLIC : public
 - └ IDENT : [a..z]+
- ▶ Début du texte d'entrée :
public public tion
- ▶ Plusieurs possibilités de découpage :
 1. IDENT("public") IDENT("publication")
 2. PUBLIC IDENT("publication")
- ▶ La règle normale est : à longueur égale du mot reconnu, c'est la première expression régulière qui gagne.

Exécution de l'automate pour le recherche d'un préfixe maximal

- ▶ Si on cherche le préfixe le plus *court* reconnu on doit s'arrêter dès qu'on arrive dans un état acceptant.
- ▶ Si on cherche le préfixe le plus *long* reconnu on doit continuer à lire tant que possible, et quand on passe par un état acceptant :
 - └ mémoriser l'état acceptant ;
 - └ mémoriser la position dans le mot d'entrée

Si l'automate ne peut plus continuer : remettre le pointeur de lecture dans l'entrée à la position où on a vu le dernier état acceptant.

Analyse de Données Structurées- Cours 3

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

4 février 2015

© Ralf Treinen 2015

Spécification d'une analyse lexicale

- Définition des catégories lexicales différentes: expressions régulières
- Définition de la stratégie: chercher le mot le plus long ou le plus court, comment résoudre des ambiguïtés
- Définir le type des jetons produits avec leurs arguments éventuels.

Objectif de l'analyse lexicale

- Lire le texte d'entrée, et faire un premier traitement en vue d'une simplification pour les étapes suivantes:
- *Découpage* de l'entrée en *lexèmes* (des mots élémentaires)
- *Classer* les lexèmes identifiés, création de jetons
- *Interpréter* les lexèmes quand pertinent, par exemple transformer une suite de chiffres en un entier
- *Abstraire* l'entrée: ignorer des détails non pertinents pour la suite (espaces, commentaires. `:::`)

Implémenter une analyse lexicale

- Soit *écrire un programme* (Java ou autre) à la main, basé sur un automate fini: discuté au dernier cours.
Les premiers compilateurs étaient effectivement écrits de cette façon (compilateur du langage FORTRAN, Backus et al. 1957: 18 personnes-années).
- Soit faire *engendrer* un analyseur lexicale à partir d'une spécification (ce cours).
C'est la technique utilisée pour l'écriture des compilateurs modernes.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

L'interface de l'analyseur lexicale

- ▶ On pourrait imaginer que l'analyseur lexicalement crée une liste Java avec tous les jetons créés lors de l'analyse.
- ▶ Problème: cette liste risquerait d'être très longue.
- ▶ Normalement, la phase suivante de l'analyse a seulement besoin de lire les jetons une fois dans l'ordre.
- ▶ Pour ces raisons, l'analyseur lexicalement crée les jetons un après l'autre *à la demande* : fonction qui renvoie le jeton suivant.

Le code engendré dans le cas de jflex

- ▶ Une classe pour l'analyseur lexicalement, le nom de la classe peut être défini dans la spécification (dans nos exemples: `Lexer`).
- ▶ La création d'un objet de cette classe (un analyseur) prend en argument un objet qui représente le *flux d'entrée*, par exemple un fichier, ou l'entrée standard.
- ▶ Il y a une méthode pour demander le jeton suivant. Le nom de cette méthode, et le type des jetons, peuvent être définis dans la spécification.

Différents générateurs

- ▶ Existents pour presque tous les langages de programmation.
- ▶ Le premier générateur était *lex*, publié en 1975 par Mike Lesk et Eric Schmidt. Engendré du code en C.
- ▶ Successeur: *flex*, 1987.
- ▶ Les générateurs modernes sont souvent issus de flex. Nous utilisons ici un générateur pour Java: *jflex*.
- ▶ Les générateurs pour d'autres langages de programmation sont très similaires.

La spécification

- ▶ Trois parties, séparées par des lignes `%%` :
 - └ code utilisateur
 - └ options et déclarations
 - └ règles lexicales
- ▶ Code utilisateur:
 - └ copiés intégralement au début du fichier engendré (avant la définition de la classe)
 - └ parties souvent vides (sauf commentaires, et `import...`)



Lapartie *OptionsetDéclarations*

- Options: commencent avec le symbole `%`. Parmi les options les plus importantes:
 - | `%class nom` : donne le nom de la classe engendrée.
 - | `%public` : la classe engendrée est publique
 - | `%type t` : le type de résultat de la fonction `yylex`.
 - | `%unicode` : accepte des caractères Unicode en entrée de l'analyse lexicale (recommandé)
 - | `%line` : compte les lignes pendant l'analyse lexicale (disponible en `yyline`).
 - | `%column` : compte les colonnes pendant l'analyse lexicale (disponible en `yycolumn`).
 - | `%state s` : Déclaration de l'état `s` (voir plus tard)

Lapartie *OptionsetDéclarations*

- Macros: système de définitions d'expressions régulières
 - | met les mots entre apostrophes `"` et `"`
 - | pour utiliser une expression régulière préalablement définie, par exemple `unom r : {r}`.
 - | classes de caractères, par exemple `[a-z]`
 - | quelques classes de caractères prédéfinies, par exemple `[:letter:]`, `[:digit:]`, `[:uppercase:]`, `[:lowercase:]`.

Lapartie *OptionsetDéclarations*

- Code entre `{` et `}` (peut être sur plusieurs lignes):
 - | copie au début de la classe engendrée
 - | le code a donc accès aux champs de la classe (par exemple, `yyline`, `yycolumn`)
- Code entre `%eofv 1{` et `%eofv 1}` : code exécuté quand l'analyse lexicale arrive à la fin de l'entrée (défaut: `null`).

Lapartie *RèglesLexicales*

- Séquence de *expression-régulière* `{ code-java }`
- dans le cas le plus simple, le code Java est un `return ...`
- Règles d'exécution: on cherche le lexème le plus long possible, et on applique l'action de la première expression régulière qui s'applique.



Le premier exemple

- Analyse lexicale pour des expressions arithmétiques comme vu la dernière fois.
- Petite différence au premier exemple : les entiers ne contiennent pas d'exposant.
- Définition des classes pour les Symboles (type de jetons), puis pour les jetons éventuellement avec des arguments.
- Le fichier de spécification pour `jflex`.
- Un petit programme principal pour tester.

Fichier Sym.java

```
public enum Sym {
    INT, IDENT, PARG, PARD, MULT, PLUS;

    public String str() {
        switch (this) {
            case INT: return "INT";
            case IDENT: return "IDENT";
            case PARG: return "PARG";
            case PARD: return "PARD";
            case MULT: return "MULT";
            case PLUS: return "PLUS";
            default : return "impossible";
        }
    }
}
```

Fichier Token.java I

```
class Token {
    protected Sym symbol;
    public Token(Sym s) {
        symbol = s;
    }
    public Sym symbol() {
        return symbol;
    }
    public String str() {
        return (symbol.str());
    }
}

class StringToken extends Token {
    private String value;
    public StringToken(Sym s, String s) {
        super(s);
        value = s;
    }
}
```

Fichier Token.java II

```
    }
    public String str() {
        return (symbol.str() + '(' + value + ')');
    }
}

class IntToken extends Token {
    private int value;
    public IntToken(Sym s, int i) {
        super(s);
        value = i;
    }
    public String str() {
        return (symbol.str() + '(' + value + ')');
    }
}
```



Fichier arith.flex I

```

%%
%public
%class Lexer
%unicode
%type Token

%{
    private Token token(Symtype){
        return new Token(type);
    }
    private StringToken token(Symtype, String value){
        return new StringToken(type, value);
    }
    private IntToken token(Symtype, int value){
        return new IntToken(type, value);
    }
}%

```

Fichier Test.java

```

import java.io.*;

class Test{

    public static void main(String[] args) throws Exception{
        File input = new File(args[0]);
        Reader reader = new FileReader(input);
        Lexer lexer = new Lexer(reader);
        Token t;
        do {
            t = lexer.yylex();
            if (t != null){System.out.println(t.str());}
        } while (t != null);
    }
}

```

Fichier arith.flex II

```

EspaceChar=[\n\r\f\t]
Ch=[0-9]
Le=[a-zA-Z]

%%
{Ch}+ { return token(Sym.INT,
                    Integer.parseInt(yytext()));}

{Le}({Le}{Ch}) * { return token(Sym.IDENT, yytext());}
"(" { return token(Sym.PARG);}
")" { return token(Sym.PARD);}
"*" { return token(Sym.MULT);}
"+" { return token(Sym.PLUS);}
{EspaceChar}+ {}

```

L'automate créé

- ▶ Création des classes de caractères : tous les caractères qui ne sont jamais distingués par les expressions régulières sont groupés dans la même classe.
- ▶ Les classes créées doivent être disjointes.
- ▶ Exemple : expressions régulières :
"end"
[-z] *
- ▶ Quatre classes de caractères disjointes : [e], [n], [d], [-cf-mo-z]

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



L'automate créé

- ▶ Création d'un automate non-déterministe pour l'union de toutes les expressions régulières.
- ▶ Déterminiser l'automate (et éliminer les `-` transitions).
- ▶ Minimiser l'automate.
- ▶ On peut demander à `jflex` de montrer ces trois automates (option `-dot`, visualiser les automates avec `xdot` par exemple)

Pourquoi utiliser plusieurs états ?

- ▶ Un premier exemple sont les commentaires : avec une expression régulière comme `"/*" . "*" */` on a un problème quand il y a plusieurs commentaires dans le texte (pourquoi ?)
- ▶ Dans ce cas on veut en fait trouver le mot *le plus court* décrit par l'expression régulière. Cela peut être simulé en utilisant deux états.

Les états de l'analyseur lexical

- ▶ Par défaut (comme sur le premier exemple), votre analyseur lexical a un seul état.
- ▶ Pour en avoir plusieurs :
 - └ les déclarer à l'aide de `%state` (sauf `YYINITIAL`)
 - └ mettre toutes les règles dans le contexte d'un état
 - └ dans les actions : changer d'état à l'aide de `yybegin`.

Reconnaître les commentaires (simplifié)

```
%%
%type Token
EspaceChar=[\n\r\f\t]
Letter=[a-zA-Z]
%state INCOMMENT
%%
<YYINITIAL>{
    {Letter}+ {          return token(Sym.IDENT,yytext());}
    {EspaceChar} {}
    "/" "*" {yybegin(INCOMMENT);}
}
<INCOMMENT>{
    "*/" {yybegin(YYINITIAL);}
    . {}
}
```



Reconnaître les commentaires

- ▶ YYINITIAL est l'état par défaut
- ▶ Il est crucial que la dernière règle s'applique à un mot de longueur 1 seulement.

Exemple : découper un mot en plusieurs parties

- ▶ Retour à notre premier exemple : on souhaite maintenant aussi reconnaître des entiers avec exposant (756e2, par exemple).
- ▶ On utilise deux états : quand on trouve un symbole "e" après une séquence de chiffres on stocke la valeur entière trouvée dans une variable, puis on va dans un deuxième état où on va lire l'exposant.

Nouvelle version de arith.flex I

```
%%

%public
%class Lexer
%type Token
%unicode

%{
    private Token token(Symtype){
        return new Token(type);
    }
    private StringToken token(Symtype, String value){
        return new StringToken(type, value);
    }
    private IntToken token(Symtype, int value){
        return new IntToken(type, value);
    }
    int intbuff=0;
```

Nouvelle version de arith.flex II

```
private String chop(Strings){
    return (s.substring(0,s.length() - 1));
}
private int expo( int base, int ex){
    int result=base;
    for ( int i=1; i<=ex; i++){
        result=result * 10;
    }
    return result;
}
%}

EspaceChar=[\n\r\f\t]
Ch=[0-9]
Le=[a-zA-Z]

%state EXPONENT

%%
```



Nouvelle version de arith.flex III

```

<YYINITIAL>{
  {Ch}+ {          return  token(Sym.INT,
                        Integer.parseInt(yytext()));}

  {Le}({Le}|{Ch}) * { return  token(Sym.IDENT,yytext());}
  "(" {          return  token(Sym.PARG);}
  ")" {          return  token(Sym.PARD);}
  "*" {          return  token(Sym.MULT);}
  "+" {          return  token(Sym.PLUS);}
  {EspaceChar}+ {}
  {Ch}+"e" { intbuff=Integer.parseInt(chop(yytext()));
              yybegin(EXPONENT);}
}
<EXPONENT>{
  {Ch}+{yybegin(YYINITIAL);
        return  (token(Sym.INT,
                        expo(intbuff,Integer.parseInt(yytext())));}
}

```

Attention à l'ordre des règles

Entrée: begbeginbeginner

- ▶ Premier appel à `yylex()` : seulement la quatrième règle s'applique \Rightarrow token IDENT.
- ▶ Deuxième appel à `yylex()` : les règles (2) et (4) s'appliquent au même lexème `begin`, c'est donc la première parmi ces deux qui gagne \Rightarrow token BEGIN.
- ▶ Troisième appel à `yylex()` : les règles (2) et (4) s'appliquent mais la dernière reconnaît un lexème plus long \Rightarrow token IDENT.

Mais regarder la taille de l'automate engendré!

Mots clefs d'un langage de programmation

Solution naïve: une règle par mot clef.

```

%%
%type Token
EspaceChar=[\n\r\f\t]
Letter=[a-zA-Z]
%%
"begin" {          return  token(Sym.BEGIN);}
"end" {          return  token(Sym.END);}
"class" {         return  token(Sym.CLASS);}
{Letter}+ {       return  token(Sym.IDENT,yytext());}
{EspaceChar} {}

```

Contrôler la taille de l'automate

- ▶ Techniques utilisées par les générateurs:
 - ↳ Utiliser des classes de caractères au lieu de la représentation de l'automate.
 - ↳ Minimiser l'automate engendré à partir d'expressions régulières.
- ▶ Optimisation dans la spécification: Éviter de créer une nouvelle classe lexicale pour chaque mot clef (Java: 46 mots clefs.)



Comment reconnaître les mots clefs sans catégories dédiées?

- ▶ En Java (et pareil dans les autres langages de programmation) : tous les mots clefs sont des séquences de lettres en minuscules.
- ▶ Mettre une seule catégorie pour les identificateurs.
- ▶ Dans l'action associée, on cherche (par ex. dans une table de hachage) si le lexème est un mot clef, et crée un jeton en fonction.

```
Lefichier keys.flex ||
```

```
%{  
private Keys keys = new Keys();  
private Token token(Sym sym, KT key, Tm t, Cc c, K k, I i) { return ee } private Token token(Sym sym, ee unnnnrT, Key sys et(unnnn);
```

```
Lefichier keys.flex |
```

```
import java.util.HashMap;
class Keys extends HashMap<String,Sym>{
    public Keys() {
        super ();
        this .put("end",Sym.END);
        this .put("begin",Sym.BEGIN);
        this .put("class",Sym.CLASS);
    }
}

%%
%public
%typeToken
%class Lexer
%unicode
EspaseChar=[\n\r\f\t]
Letter=[a -zA-Z]
```

Analyse de Données Structurées- Cours 4

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

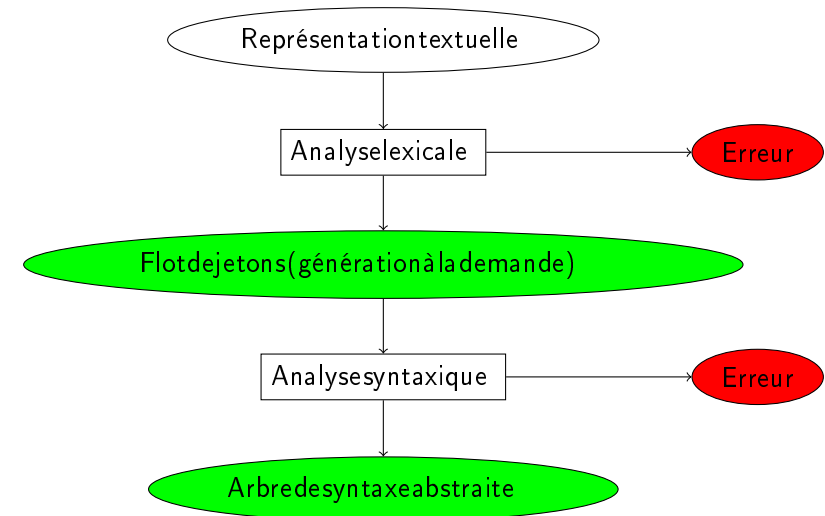
11 février 2015

© Ralf Treinen 2015

Objectif de l'analyse syntaxique

- ▶ Deux objectifs:
 - ┌ Détecter des textes d'entrée qui ne sont pas correctement formés (et donner des indications utiles sur la nature de l'erreur).
 - ┌ Si l'entrée est correcte, construire un arbre de syntaxe abstraite.
- ▶ Focalisons d'abord sur le premier objectif: distinguer des textes corrects de textes incorrects.

Le rôle de l'analyse syntaxique



Pourquoi deux étapes séparées?

- ▶ On essaye de faire tant d'analyse que possible dans l'analyse lexicale.
- ▶ Raison: l'exécution d'un automate est très efficace (temps linéaire dans la longueur du texte d'entrée).
- ▶ Problème: l'expressivité des automates est limitée (voir les transparents suivants).



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Reconnaître les expressions arithmétiques

- ▶ L'ensemble des expressions arithmétiques correctement parenthésées n'est pas régulier.
- ▶ Nous démontrons une simplification : l'ensemble des mots formés seulement de parenthèses '[' et ']' qui sont correctement imbriqués, n'est pas régulier.
- ▶ Exemple d'un mot correctement imbriqué :

$[[[]]]$

- ▶ Exemple d'un mot qui n'est pas correctement imbriqué :

$[[[]]$

Un lemme

Lemme :

Le langage $L = \{[n]^n \mid n \geq 1\}$ n'est pas régulier.

Intuition

- ▶ L consiste en tous les mots de la forme

$\underbrace{[\dots]}_n \underbrace{]\dots]}_n$

pour un $n \geq 1$ quelconque.

- ▶ Un automate qui accepte L devrait compter les '[' et puis comparer ce nombre avec le nombre des ']'.
- ▶ Or, un automate fini ne peut pas compter une quantité non bornée.

Reconnaître l'imbrication correcte des parenthèses

Définition du langage P

P est le plus petit langage tel que

- ▶ $\epsilon \in P$
- ▶ si $x, y \in P$ alors $xy \in P$
- ▶ si $x \in P$ alors $[x] \in P$

Théorème

P n'est pas régulier.

Du lemme au théorème

- ▶ Avant de démontrer le lemme, réfléchissons pourquoi le lemme implique le théorème :
- ▶ Supposons pour l'absurde que P soit régulier.
- ▶ Le langage $L^0 = \{[n]^m \mid n, m \geq 0\}$ est régulier.
- ▶ $L = P \cap L^0$
- ▶ L'intersection de deux langages réguliers est régulière : contradiction !



Edited with the demo version of
Infix Pro PDF Editor

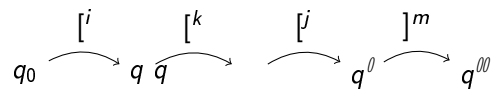
To remove this notice, visit:
www.iceni.com/unlock.htm

La preuve du lemme

- Supposons, pour l'absurde, que l'automate $A = (Q; q_0; F; \Delta)$ accepte L . Soit m le nombre d'éléments de Q .
 - A accepte donc le mot $[^m]^m$.
 - Regardons les actions de l'automate quand il lit m fois le symbole '[' :
- $$q_0 \xrightarrow{[} q_1 \xrightarrow{[} q_2 \dots q_{m-1} \xrightarrow{[} q_m$$
- Ça fait $m + 1$ états dans cette séquence, mais il n'y a que m états différents.

La preuve du lemme

- On a donc pour le mot $[^m]^m$:



avec $m = i + k + j$.

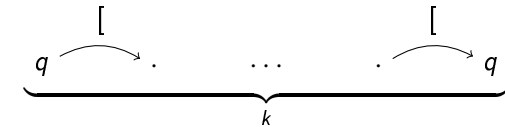
- Donc, on a aussi :



- Donc : $[^m \text{ } ^k]^m \in \mathcal{L}(A)$, ce qui est absurde !

La preuve du lemme

- Donc il y a un état, disant q qui paraît deux fois dans cette séquence :



- Soit k la longueur de cette séquence. Nous avons donc $\Delta(q; [^k) = q$.

Le résultat général

Lemme de l'étoile (angl. : *pumping lemma*)

Pour tout langage régulier L il existe un entier m tel que tout mot $w \in L$ de longueur $|w| \geq m$ a une décomposition $w = xyz$ telle que

- $0 < |y|$
- $|xy| \leq m$
- $xy^n z \in L$ pour tout $n \geq 0$

Démonstration

Comme sur l'exemple des transparents précédents, Edited with the demo version of

Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



Définition des grammaires algébriques

Une **grammaire algébrique** (où : **grammaire hors contexte**) est un tuple $G = (V_T; V_N; S; R)$ où

- ▶ V_T est un ensemble fini de symboles dits **terminaux**;
- ▶ V_N est un ensemble fini de symboles **non-terminaux**, où $V_N \cap V_T = \emptyset$;
- ▶ $S \in V_N$ est l'**axiome**;
- ▶ R est un ensemble fini de **règles** de la forme $N \rightarrow u$, où $N \in V_N$, et $u \in (V_N \cup V_T)^*$.

Dérivation

Réécriture

Étant donnée une grammaire $G = (V_T; V_N; S; R)$, on note $V = V_T \cup V_N$.

Le mot $u \in V^*$ se **réécrit** en mot $v \in V^*$ dans G , noté $u \rightarrow v$, si

1. $u = w_1 N w_2$, où $w_1 \in V^*$, $N \in V_N$ et $w_2 \in V^*$;
2. $N \rightarrow w$ est une règle de R ;
3. $v = w_1 w w_2$.

Dérivation

Le mot $v \in V^*$ **dérive** du mot $u \in V^*$ dans la grammaire G , noté $u \rightarrow^* v$, s'il existe une suite finie $w_0; w_1; \dots; w_n$ de mots de V^* telle que $w_0 = u$, $w_i \rightarrow w_{i+1}$ pour tout $i \in [0; n-1]$, et $w_n = v$.

Exemple d'une grammaire algébrique

$G_1 = (V_T; V_N; S; R)$ où

- ▶ $V_T = \{0; 1; +; *\}$
- ▶ $V_N = \{E; I\}$
- ▶ $S = E$
- ▶ R consiste en les règles suivantes:

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow I$

$I \rightarrow 0$

$I \rightarrow 1$

Exemple

- ▶ Pour la grammaire G donnée ci-dessus, on a les réécritures suivantes:
(Nous mettons en **rouge** le non-terminal qui est réécrit.)

| **E** $\rightarrow E + E$

| $E + \mathbf{E} \rightarrow E + E + E$

| $E + E + \mathbf{E} \rightarrow E + E + E * E$

| $E * E + \mathbf{E} * E \rightarrow E * E + I * E$



Exemple: Dérivation de $1+0*1$ dans G_1

$E \rightarrow E + E$
 $\rightarrow E + E * E$
 $\rightarrow I + E * E$
 $\rightarrow I + E * I$
 $\rightarrow I + E * 1$
 $\rightarrow 1 + E * 1$
 $\rightarrow 1 + I * 1$
 $\rightarrow 1 + 0 * 1$

Relation avec les langages réguliers

- ▶ Tout langage régulier est algébrique.
- ▶ Par contre, il y a des langages algébriques qui ne sont pas réguliers: $G_P = (\{[,]\}; \{E\}; R)$ où R consiste en

$E \rightarrow$
 $E \rightarrow EE$
 $E \rightarrow [E]$

On a que $\mathcal{L}(G_P) = P$, le langage des séquences de parenthèses correctement imbriquées.

Langage engendré

Définition

Soit $G = (V_T; V_N; S; R)$ une grammaire algébrique. Le *langage engendré* par G est

$$\mathcal{L}(G) = \{w \in V_T \mid S \rightarrow w\}$$

Un langage est *algébrique* s'il est engendré par une grammaire algébrique.

Exemple

Pour la grammaire G_1 on obtient que $\mathcal{L}(G_1)$ est l'ensemble des expressions arithmétiques formées à l'aide des opérateurs $+$ et $*$, et des constantes 0 et 1 , et sans parenthèses.

Notation: alternatives

- ▶ Une grammaire peut avoir plusieurs règles avec le même côté gauche:

$E \rightarrow E + E$
 $E \rightarrow E * E$

- ▶ Nous permettons dans la suite dans ce cas d'écrire une seule règle, avec plusieurs *alternatives* sur le côté droit:

$E \rightarrow E + E \mid E * E$



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Dérivation gauche

Soit une grammaire $G = (V_T; V_N; S; R)$; on note $V = V_T \cup V_N$.

Réécriture à gauche

Le mot $u \in V$ se *réécrit à gauche* en le mot $v \in V$ dans G , noté $u \rightarrow_g v$, si

1. $u = w_1 N w_2$, où $w_1 \in V_T^*$, $N \in V_N$ et $w_2 \in V^*$;
2. $N \rightarrow w$ est une règle de R ;
3. $v = w_1 w w_2$.

Dérivation gauche

Une *dérivation gauche* est une suite finie $w_0; w_1; \dots; w_n$ de mots de V^* telle que $w_i \rightarrow_g w_{i+1}$ pour tout $i \in [0; n-1]$.

Exemple: Dérivation gauche de $1+0*1$ dans G_1

$$\begin{aligned} E &\rightarrow E * E \\ &\rightarrow E + E * E \\ &\rightarrow I + E * E \\ &\rightarrow 1 + E * E \\ &\rightarrow 1 + I * E \\ &\rightarrow 1 + 0 * E \\ &\rightarrow 1 + 0 * I \\ &\rightarrow 1 + 0 * 1 \end{aligned}$$

Arbre de dérivation

Définition

Soit $G = (V_T; V_N; S; R)$ une grammaire. Un *arbre de dérivation* de G est un arbre tel que

- ▶ les nœuds internes sont étiquetés par des symboles de V_N ;
- ▶ les feuilles sont étiquetées par des symboles de $V_T \cup V_N$;
- ▶ si les fils pris de gauche à droite d'un nœud interne étiqueté par le non-terminal N sont étiquetés par les symboles respectifs $x_1; \dots; x_n$, alors $(N \rightarrow x_1 \dots x_n) \in R$.

Définition

Un arbre de dérivation dont la racine est étiquetée par l'axiome de la grammaire et dont le mot des feuilles u appartient à V_T^* est appelé *arbre de dérivation de u* .

Exemple: Un arbre de dérivation de $1+0*1$ dans G

D'une dérivation gauche à l'arbre de dérivation

- ▶ La racine de l'arbre est étiquetée par l'axiome de la grammaire (qui est aussi le premier élément de la dérivation gauche).
- ▶ On parcourt la dérivation de gauche à droite, en complétant l'arbre.
- ▶ À la réécriture du premier non-terminal dans un mot correspond l'extension de la feuille de l'arbre de dérivation en construction la plus gauche qui est étiquetée par un non-terminal.
- ▶ (Voir l'exemple donné au tableau)

Arbres de dérivation et Dérivation gauche

- ▶ Équivalence entre dérivation gauche et arbre de dérivation :
 - └ Pour chaque arbre de dérivation il y a une unique dérivation gauche.
 - └ Pour chaque dérivation gauche il y a un unique arbre de dérivation.
- ▶ Dans un premier temps, on peut imaginer l'arbre de dérivation comme résultat de l'analyse syntaxique.
- ▶ Il nous faut donc que chaque mot du langage possède un arbre de dérivation unique.

De l'arbre de dérivation à la dérivation gauche

- ▶ Le premier élément de la dérivation gauche est l'axiome de la grammaire (qui est aussi l'étiquette de la racine de l'arbre).
- ▶ On fait un parcours préfixe de l'arbre, en complétant la dérivation gauche.
- ▶ Aux fils d'un nœud interne correspond la réécriture du premier non-terminal dans un élément de la dérivation.
- ▶ (Voir l'exemple donné au tableau)

Grammaires ambiguës

Définition

Une grammaire G est non-ambiguë quand tout $w \in \mathcal{L}(G)$ a un seul arbre de dérivation.

Définition équivalente

Une grammaire G est non-ambiguë quand tout $w \in \mathcal{L}(G)$ a une seule dérivation gauche.

Sur l'exemple G_1

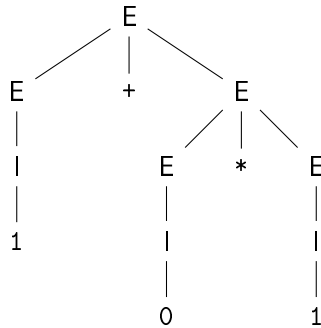
La grammaire G_1 est ambiguë : le mot $1+0*1$ a deux arbres de dérivation différents!



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Un autre arbre de dérivation de $1+0*1$ dans G_1



G_2 est non-ambiguë

- Pour le démontrer nous avons besoin de deux petits lemmes.
- Utilité de ces lemmes : Imaginez qu'un mot w a deux dérivations gauches qui sont de la forme

$$E \rightarrow (E + E) \rightarrow w$$

$$E \rightarrow (E * E) \rightarrow w$$

- Soit w_1 le mot dérivé du E gauche de la première ligne, w_2 le mot dérivé du E gauche de la deuxième ligne.
- Forcément, $w_1 \neq w_2$. Donc, $w_1 <_{\text{pref}} w_2$ (ou l'inverse).
- Les lemmes nous permettent de démontrer que c'est impossible d'avoir $w_1, w_2 \in \mathcal{L}(G_2)$ avec $w_1 <_{\text{pref}} w_2$.

Un autre exemple

- La grammaire $G_2 = (\{i; v; +; *; (;)\}; \{E\}; E; R)$, où R est

$$E \rightarrow i \mid v \mid (E + E) \mid (E * E)$$

- Cette grammaire décrit les expressions arithmétiques complètement parenthésées.
- Ici il y a un seul terminal i pour les entiers, et un seul terminal v pour les noms des variables. On imagine qu'ils s'agit des jetons issus d'une analyse lexicale.

Séparer $\mathcal{L}(G_2)$ en ses propres préfixes

Notation

- $|w|_{(}$ est le nombre de symboles (dans le mot w).
- $|w|_{)}$ est le nombre de symboles (dans le mot w).

Lemme 1

Pour tous $w \in \mathcal{L}(G_2)$: $|w|_{(} = |w|_{)}$.

Démonstration

Par induction sur la longueur de la dérivation la plus courte de w (au tableau).



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Séparer $\mathcal{L}(G_2)$ et ses propres préfixes

Lemme 2

Pour tous $w \in \mathcal{L}(G_2)$: pour tous $v <_{\text{pref}} w$ avec $v \neq \epsilon$:
 $|v|_c > |v|_l$.

Démonstration

Par induction sur la longueur de la dérivation la plus courte de w
(au tableau).

G_2 est non-ambiguë

- ▶ Soit $v_1 = u_1$, alors forcément $\circ = \diamond$, $v_2 = u_2$, et les deux dérivationes gauches sont les mêmes. Contradiction!
- ▶ Soit $v_1 <_{\text{pref}} u_1$. On a donc $v_1 \in \mathcal{L}(G_2)$, et $v_1 <_{\text{pref}} u_1 \in \mathcal{L}(G_2)$: Contradiction aux Lemmes 1 et 2!
- ▶ Soit $u_1 <_{\text{pref}} v_1$. On a donc $u_1 \in \mathcal{L}(G_2)$, et $u_1 <_{\text{pref}} v_1 \in \mathcal{L}(G_2)$: Contradiction aux Lemmes 1 et 2!

G_2 est non-ambiguë

- ▶ Supposons pour l'absurde que G_2 est ambiguë.
- ▶ Soit w un mot de longueur minimale qui a deux dérivationes gauches.
- ▶ Il n'est pas possible qu'une des deux commences sur $E \rightarrow i$ ou $E \rightarrow v$.
- ▶ Donc, une commences sur $E \rightarrow (E \circ E)$, et l'autre sur $E \rightarrow (E \diamond E)$, avec $\circ, \diamond \in \{+, *\}$.
- ▶ Donc :

$$w = (v_1 \circ v_2)$$

$$w = (u_1 \diamond u_2)$$

avec $E \rightarrow v_1; v_2; u_1; u_2$, de longueurs plus petites que w .
Chacun de ces mots a donc une seule dérivation gauche!

Backus-Naur Form

- ▶ Une notation très utilisée pour la documentation des langages de programmation est la notation **BNF** (Backus-Naur Form) qui est équivalente aux grammaires algébriques.
- ▶ Les non-terminals sont écrits entre chevrons : $\langle T \rangle$.
- ▶ La flèche est remplacée par $::=$.
- ▶ Les alternatives pour le même non-terminal sont regroupées comme indiquées au-dessus.
- ▶ Exemple :

$\langle \text{cond} \rangle ::= \text{if} (" \langle \text{condition} \rangle ") \{ " \langle \text{code} \rangle " \}$



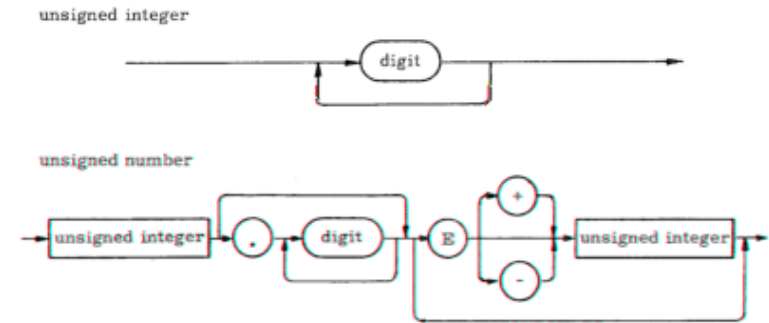
Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

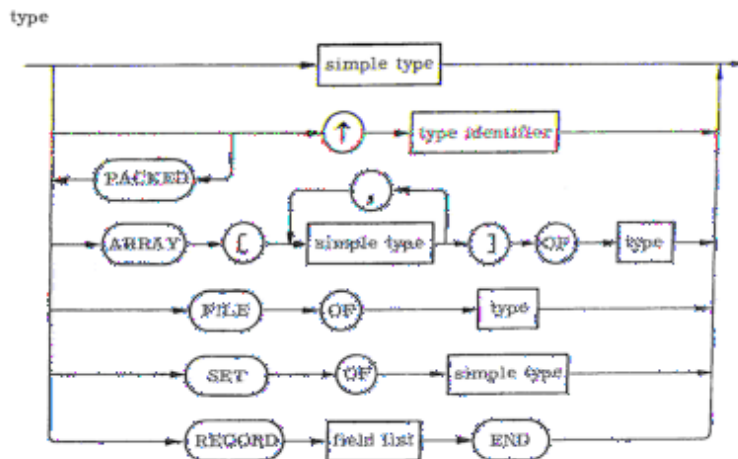
ExtendedBackus-NaurForm

- ▶ Laformeétenduepermetaussidesconstructionsdes expressionsrégulièrèsdansla grammaire:opérateurs * et +, etparenthèses.
- ▶ Unepartieentrecrochets[]estoptionnelle.
- ▶ Exemple: $\langle \text{exp list} \rangle ::= "(" \langle \text{exp} \rangle ("," \langle \text{exp} \rangle)^* ")"$
- ▶ C'estencoreéquivalentauxgrammairesalgébriques.

Diagrammedesyntaxe(non-récursif)



Diagrammedesyntaxe(récursif)



AnalysedeDonnéesStructurées- Cours5

RalfTreinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

18 février 2015

© RalfTreinen2015

L'objectif

- ▶ L'analysesyntaxiqueadeuxobjectifs:
 - └ détectersiletexteluestcorrect (est dans le langage engendré par la grammaire);
 - └ le cas échéant, construire un arbre de syntaxe abstraite.
- ▶ Pour l'instant, nous continuons à étudier le premier problème: reconnaissance de textes d'entrée correctes.

Déterminer un automate fini

- ▶ À chaque moment de l'exécution d'un automate fini, sa configuration consiste en :
 - └ l'information où on est dans la lecture du mot d'entrée, et
 - └ l'état de l'automate (nombre fini).
- ▶ L'ensemble des configurations possibles dans lesquelles l'automate peut se trouver après la lecture d'un mot donné est borné (par le nombre d'états de l'automate).
- ▶ Construction d'un automate déterministe de taille exponentielle.

Plusieurs techniques pour l'analyse syntaxique

- ▶ Analyse descendante : construction (virtuelle) de l'arbre de dérivation, à partir de l'axiome aux feuilles.
Ordre de construction : parcours préfixé de l'arbre.
C'est l'approche présentée dans ce cours.
- ▶ Analyse ascendante : construction (virtuelle) d'un arbre de dérivation à partir des feuilles jusqu'à l'axiome. Plus complexes à maîtriser, nécessitent des connaissances des automates à pile. Voir le cours de Compilation au M1.

Déterminer un automate à pile ?

- ▶ La configuration d'un automate à pile est plus riche :
 - └ l'information où on est dans la lecture du mot d'entrée,
 - └ l'état de l'automate,
 - └ et le *contenu de la pile*.
- ▶ La taille de la pile est *non-bornée*.
- ▶ Le nombre des configurations possibles dans lesquelles l'automate peut se trouver après la lecture d'un mot n'est plus borné !

Construction (virtuelle) d'un arbre de dérivation

- ▶ Dans la construction d'un arbre de dérivation (ou, d'une dérivation), il y a à chaque moment deux choix à faire :
 - └ d'un non-terminal qu'on va remplacer à l'aide d'une règle de la grammaire,
 - └ une fois le non-terminal choisi, de la règle parmi celles qui ont un non-terminal sur le côté gauche.
- ▶ Nous avons vu la semaine dernière que le premier choix n'est pas essentiel : on peut imposer une stratégie comment choisir le non-terminal à remplacer (par ex., celui qui est le plus à gauche).



Exploration complète de l'espace de recherche?

- ▶ Une façon de réaliser une analyse syntaxique est maintenant d'essayer simplement toutes les possibilités de choisir des règles.
- ▶ Cela donner lieu à un algorithme non-déterministe:
 - └ soit par *retour-en-arrière* (angl.: backtracking)
 - └ soit par *programmation dynamique*
- ▶ Approche complète: on est sûr de trouver un arbre de dérivation si le mot est dans le langage ☺
- ▶ Problème: efficacité ☹
- ▶ On cherche des solutions efficaces, éventuellement en imposant des restrictions aux grammaires qu'on peut traiter.

Exemple

- ▶ Grammaire $G = (V_T, V_N, S, R)$ où
- ▶ $V_T = \{i, +, [,]\}$
- ▶ $V_N = \{S\}$
- ▶ $S = S$
- ▶ R consiste en les règles suivantes:

$$S \rightarrow i \quad (1)$$

$$S \rightarrow [S+S] \quad (2)$$

- ▶ $\mathcal{L}(G)$: expressions complètement parenthésées, construites avec la constante i et l'opérateur binaire $+$.

Comment obtenir une solution efficace?

- ▶ Il faut maîtriser le choix de la règle de la grammaire par laquelle on va remplacer un non-terminal.
- ▶ On ne peut pas demander qu'il y ait une seule règle par non-terminal (car dans ce cas la grammaire est complètement triviale).
- ▶ Sur quoi baser le choix de la règle?
- ▶ Sur la suite du mot pour lequel on cherche à construire l'arbre de dérivation!

Construction d'un arbre de dérivation (1)

S

[[i	+	i]	+	[i	+	i]]
---	---	---	---	---	---	---	---	---	---	---	---	---

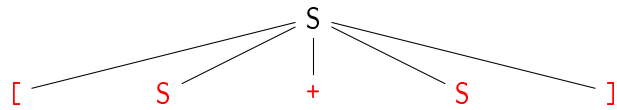
Choisir règle (2): c'est la seule qui peut produire à partir de S un mot qui commence sur $[$.



Edited with the demo version of
Infix Pro PDF Editor

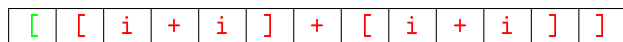
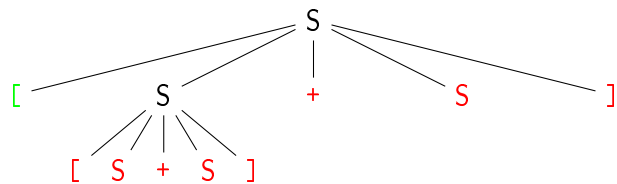
To remove this notice, visit:
www.iceni.com/unlock.htm

Construction d'un arbre de dérivation (2)



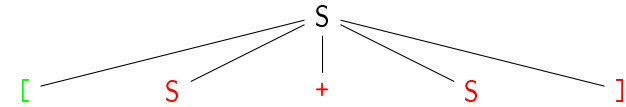
Le premier non-terminal du mot des feuilles est [.

Construction d'un arbre de dérivation (5)



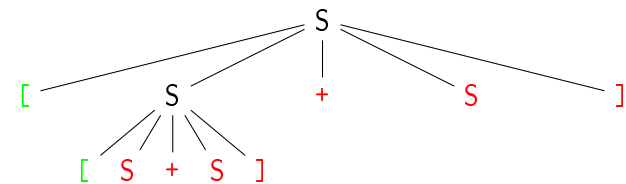
Le suivant non-terminal du mot des feuilles est [.

Construction d'un arbre de dérivation (4)



Choisir règle (2) : c'est la seule qui peut produire à partir de S un mot qui commence sur [.

Construction d'un arbre de dérivation (6)



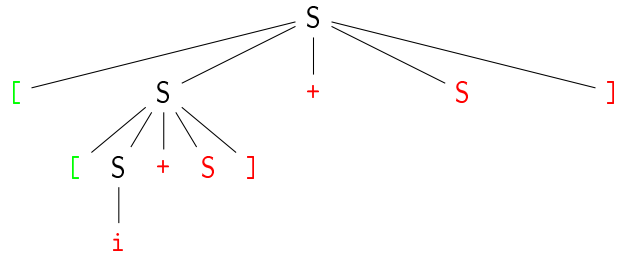
Choisir règle (1) : c'est la seule qui peut produire à partir de S un mot qui commence sur i.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

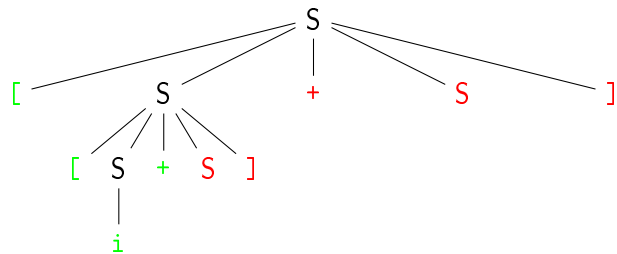
Constructiond'unarbrededérivation(7)



[[[i + i]] + [i + i]]

Lesuivantnon-terminaldumotdesfeuillesest i.

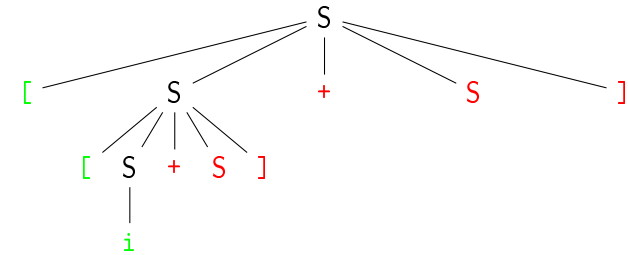
Constructiond'unarbrededérivation(9)



[[[i + i]] + [i + i]]

Choisirrègle(1):c'estlaseulequiputproduireàpartirdeSun motquicommmencesur i.

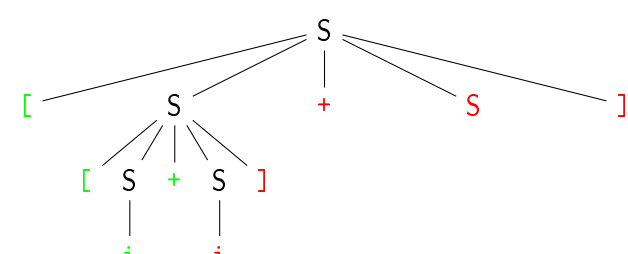
Constructiond'unarbrededérivation(8)



[[[i + i]] + [i + i]]

Lesuivantnon-terminaldumotdesfeuillesest +.

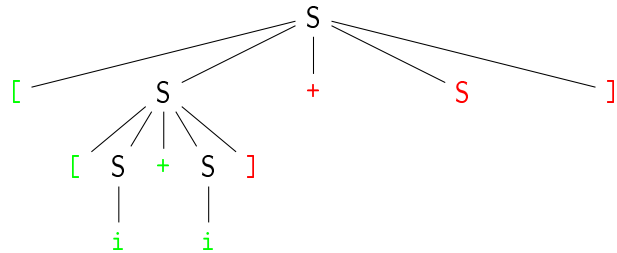
Constructiond'unarbrededérivation(10)



[[[i + i]] + [i + i]]

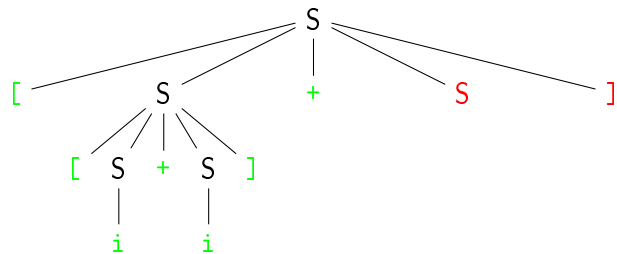
Lesuivantnon-terminaldumotdesfeuillesest

Constructiond'unarbrededérivation(11)



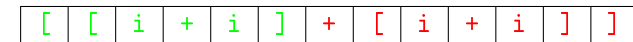
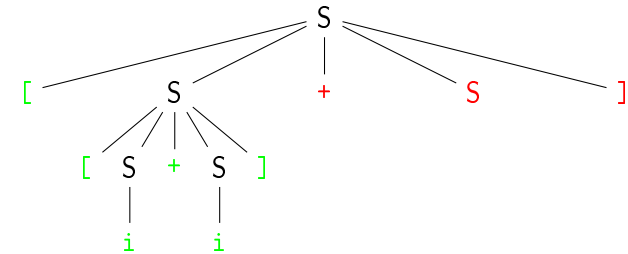
Lesuivantnon-terminaldumotdesfeuillesest `]`.

Constructiond'unarbrededérivation(13)



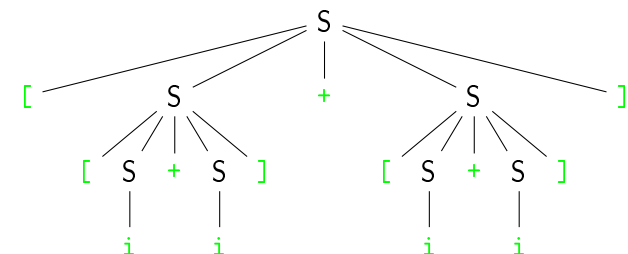
Choisirrègle(2):c'estlaseulequipeutproduireàpartirdeSun motquicommmencesur `[`.

Constructiond'unarbrededérivation(12)

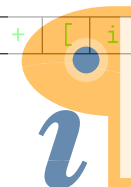


Seulemlarègle(2)peutproduireàpartirdeSun motqui commencesur `+`.

etc.etc.



Constructionterminée!



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Ce qu'on a vu sur l'exemple :

- ▶ Il y a deux types d'actions :
 - └ consommer en parallèle un non-terminal du préfixe du mot des feuilles déjà construit, et le même symbole de l'entrée;
 - └ ajouter des fils à une feuille de l'arbre de dérivation partiel.
- ▶ Pour choisir la règle de la grammaire, on regarde en avant quel est le symbole suivant de l'entrée que nous aurions à consommer (lookahead).

Notation: $w : k$

Définition

Soit $w \in \Sigma^*$ un mot, et $k \in \mathbb{N}$. On définit

- ▶ si $|w| \leq k$ alors $w : k = w$
- ▶ si $|w| > k$ alors $w : k = x$ tel que $w = xy$ et $|x| = k$

Explication

- ▶ $w : k$ est le préfixe de longueur k du mot w , ou le mot w entier si w est plus court que k .
- ▶ $abcdefg : 3 = abc$
- ▶ $abcd : 7 = abcd$

Grammaires LL(1)

En fait, l'algorithme que nous avons vu sur l'exemple appartient à la classe $LL(1)$:

- ▶ le premier L indique qu'on parcourt l'entrée de la gauche (angl. : left) à la droite;
- ▶ le deuxième L indique qu'on construit une dérivation gauche (angl. : left), c.-à-d. un arbre de dérivation dans un ordre préfixe;
- ▶ le nombre 1 indique que nous utilisons la connaissance de 1 caractère dans la partie de l'entrée qui reste à consommer, pour déterminer la règle à appliquer (lookahead=1).

Définition LL(k)

Définition

Soit $G = (V_T, V_N, S, R)$ une grammaire algébrique, $k \in \mathbb{N}$. G est dite $LL(k)$ ssi

- ▶ S'il existe deux dérivation gauche

$$\begin{aligned} S &\rightarrow^* uY\alpha \rightarrow u\beta\alpha \rightarrow^* ux \\ S &\rightarrow^* uY\alpha \rightarrow u\gamma\alpha \rightarrow^* uy \end{aligned}$$

où $Y \in V_N$, $u, x, y \in V_T^*$, $\alpha, \beta, \gamma \in (V_N \cup V_T)^*$, avec

$$x : k = y : k$$

- ▶ alors $\beta = \gamma$.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Explication de la définition de LL(k)

- ▶ On a déjà consommé le mot déterminaux u .
- ▶ Le non-terminal le plus à gauche à réécrire est maintenant Y .
- ▶ Dans les deux cas considérés, le mot d'entrée continu une fois par le mot x , l'autre fois par le mot y .
- ▶ En regardant les k premiers caractères de la suite du mot d'entrée, on peut maintenant décider comment réécrire le non-terminal y .

Un meilleur critère pour être LL(1)?

Problème

Le critère du transparent précédent est un peu trop restrictif car il ne permet pas des règles où le côté droit commence par un non-terminal:

$$\begin{aligned} A &\rightarrow BA \mid B \quad (3) \\ B &\rightarrow \dots \quad (4) \end{aligned}$$

Un premier critère simple pour être LL(1)

Lemme

Si pour tout non-terminal, les côtés droites de toutes les règles pour ce non-terminal commencent par des terminaux différents, alors la grammaire est LL(1).

Exemple

La grammaire de l'exemple précédent:

$$\begin{aligned} S &\rightarrow i \\ S &\rightarrow [S+S] \end{aligned}$$

satisfait le critère, et est donc LL(1).

La fonction FIRST_k

Définition

Soit $G = (V_T, V_N, S, R)$ une grammaire, et $k \in \mathbb{N}$. Nous définissons une fonction

$$\text{FIRST}_k : (V_T \cup V_N)^* \rightarrow 2^{V_T^*}$$

par

$$\text{FIRST}_k(\alpha) = \{w : k \mid w \in V_T^*, \alpha \rightarrow^* w\}$$

Explication

$\text{FIRST}_k(\alpha)$ est l'ensemble des préfixes de longueur k des mots terminaux qu'on peut obtenir à partir de α .

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



UnmeilleurbcritèrepourêtreLL(1)

Lemme

Soit $G = (V_T, V_N, S, R)$ unegrammaire *sansproductionsdela forme* $N \rightarrow \epsilon$. G estLL(1)siestseulementsipourtoutesrègles différentes:

$$N \rightarrow \alpha$$

$$N \rightarrow \beta$$

onaque $\text{FIRST}_1(\alpha) \cap \text{FIRST}_1(\beta) = \emptyset$.

Exemple

Toujourssurlemêmeexemple:

$$\begin{aligned}\text{FIRST}_1(i) &= \{i\} \\ \text{FIRST}_1([S+S]) &= \{[\]\}\end{aligned}$$

Exemple

- Grammaire $G = (\{ \ , (,), + \}, \{F, S\}, S, R)$ où R est

$$F \rightarrow$$

$$S \rightarrow (F+S)$$

$$S \rightarrow F$$

- Initialisation:

$$\begin{aligned}Fi() &= \{ \} \\ Fi((F+S)) &= \{ (\} \\ Fi(F) &= \emptyset\end{aligned}$$

- Complétion: Onaunerègle $F \rightarrow a$, mais $Fi() \not\subseteq Fi(F)$.
- Onaugmente: $Fi(F) = \{ \}$

Calculde FIRST_1

- Grammaire $G = (V_T, V_N, S, R)$. *Hypothèse: aucune production* $N \rightarrow \epsilon$.
- On calcule $Fi(\alpha)$ pour tout côté droit de R .
- $Fi(a\alpha) = \{a\}$ pour tout $a \in V_T$
 $Fi(N\alpha) = \emptyset$ pour tout $N \in V_N$
- Tantqu'il existe une règle $N \rightarrow \beta$ telle que $Fi(\beta) \not\subseteq Fi(N\alpha)$:

$$Fi(N\alpha) = Fi(N\alpha) \cup Fi(\beta)$$

- Pour tout α : $\text{FIRST}_1(\alpha)$ est la valeur finale de $Fi(\alpha)$.

Exemple(2)

- Grammaire $G = (\{ \ , (,), + \}, \{F, S\}, S, R)$ où R est

$$F \rightarrow$$

$$S \rightarrow (F+S)$$

$$S \rightarrow F$$

- On obtient donc:

$$\begin{aligned}\text{FIRST}_1() &= \{ \} \\ \text{FIRST}_1((F+S)) &= \{ (\} \\ \text{FIRST}_1(F) &= \{ \}\end{aligned}$$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



Structuredu code

- ▶ Nousavonsbesoin deregarderlesymbolesuivantdansleflot d'entréesansdeleconsommer:classe LookAheadReader.
- ▶ L'analyseursyntaxiquecontientuneméthode term(ch rc) quiconsommeunsymbole c duflotd'entrée.
- ▶ L'analyseursyntaxiquecontientuneméthode nonterm_N() pourchaquenon-terminal N.Cetteméthodeconsommeun motduflotd'entréequiestengendrépar N. L'implémentationdecesméthodesutilise FIRST₁ pour déterminerlarègledelagrammaireàappliquer.

Fichier LookAhead1Reader.java II

```
public void eat( char expected)
    throws ReadException, IOException{
    /* consumes c from the stream, exception */
    /* when the contents does not start on c. */
    char found=( char ) this .read();
    if (found!=expected){
        throw new ReadException(expected,found);
    }
}
```

Fichier LookAhead1Reader.java I

```
import java.io. *;

class LookAhead1Reader extends PushbackReader{
    /* Reader class with a lookahead of one character */

    public LookAhead1Reader(Reader r){
        super (r,1);
    }

    public boolean check( char c)
        throws IOException{
        /* check whether c is the first character */
        int lastread=0;
        lastread= this .read();
        this .unread(lastread);
        return (lastread==c);
    }
}
```

Fichier Parser.java I

```
import java.io. *;

/* simple LL(1) parser for the grammar: */
/* F -> a S -> F S -> (F+S) */
class Parser{

    protected LookAhead1Reader reader;

    public Parser(LookAhead1Reader r){
        reader=r;
    }
}
```



Fichier Parser.java II

```
public void term( char c) throws IOException, ReadException{
    /* consume the character c */
    reader.eat(c);
}

public void nonterm_S()
    throws IOException, IllegalArgumentException,
        ReadException, ParseException{
    /* parse a word generated from nonterminal S */
    if (reader.check('a')){
        this.nonterm_F();
    } elseif (reader.check('(')){
```

Fichier Parser.java IV

```
        this.term('a');
    }
}
```

Fichier Parser.java III

```
        this.term('(');
        this.nonterm_F();
        this.term('+');
        this.nonterm_S();
        this.term(')');
    } else {
        throw new ParseException("cannot reduce S");
    }
}

public void nonterm_F() throws IOException, ReadException{
    /* parse a word generated from nonterminal F */
}
```

Fichier Test.java

```
import java.io.*;

class Test{

    public static void main(String[] args) throws Exception{
        File input = new File(args[0]);
        Reader reader = new FileReader(input);
        LookAhead1Reader r = new LookAhead1Reader(reader);
        Parser p = new Parser(r);
        p.nonterm_S();
    }
}
```



Analyse de Données Structurées- Cours 6

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

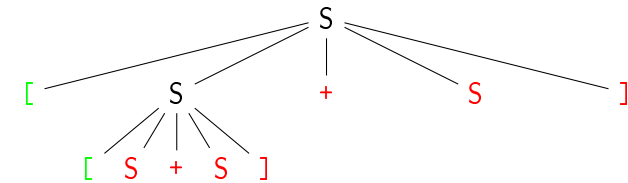
11 mars 2015

© Ralf Treinen 2015

Grammaires LL(1)

- ▶ Intuition derrière les grammaires LL(1): dans la construction d'une dérivation gauche, les symboles suivants de l'entrée nous indiquent quelle règle de la grammaire appliquera au non-terminal le plus à gauche de l'arbre de dérivation.
- ▶ Conséquence: toute grammaire LL(1) (ou même LL(k)) est non-ambiguë (oublié de préciser au dernier cours).
- ▶ Un critère très simple: Si tous les côtés droits de la grammaire pour le *même* non-terminal commencent sur des terminaux différents, alors la grammaire est LL(1).

Construction d'un arbre de dérivation



Choisir règle (1): c'est la seule qui peut produire à partir de S un mot qui commence sur i.

Une caractérisation des grammaires LL(k)

Théorème

Une grammaire $G = (V_T; V_N; S; R)$ est LL(k)ssi

- ▶ si $A \rightarrow \alpha$ et $A \rightarrow \beta$ sont deux règles différentes
- ▶ et $S \rightarrow^* wA$ une dérivation gauche, $w \in V_T^*$
- ▶ alors $\text{FIRST}_k(\alpha) \cap \text{FIRST}_k(\beta) = \emptyset$

- ▶ Preuve omise (conséquence immédiate de la définition)
- ▶ Attention, ce critère général prend en compte le "contexte" dans lequel on peut obtenir A.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Un meilleur critère pour être LL(1)

- Définition de $FIRST_1()$: l'ensemble des symboles avec lesquels un mot terminal dérivé à partir de S peut commencer (plus dans le cas où $S \rightarrow^*$ ϵ).
- Vu au dernier cours : calcul de $FIRST_1$ dans le cas où aucun côté droit n'est ϵ .
- Meilleur critère : Si tous les côtés droits de la grammaire pour le *même* non-terminal ont des ensembles $FIRST_1$ disjoints, alors la grammaire est LL(1).

Reconnaître la fin de l'entrée

- Le programme vu au dernier cours a un défaut : il accepte aussi des expressions correctes, avec un texte quelconque ajouté à la fin : $(+ (+)) \% \$ \# @$
- Solution :
 - └ l'analyse lexicale envoie un jeton qui signale la fin de l'entrée (par exemple, EOF)
 - └ remplacer l'axiome par $S' \rightarrow S \text{ EOF}$

$$S' \rightarrow S \text{ EOF}$$

où S est l'ancien axiome de la grammaire.

Exemple vu au dernier cours

- Grammaire $G = (\{ ; (;) ; + \}; \{ F ; S \}; S ; R)$ où R est

$$\begin{aligned} F &\rightarrow \\ S &\rightarrow (F+S) \\ S &\rightarrow F \end{aligned}$$

- Le premier critère simple ne s'applique pas.
- On obtient pour les côtés droits des règles :

$$\begin{aligned} FIRST_1() &= \{ \} \\ FIRST_1((F+S)) &= \{ (\} \\ FIRST_1(F) &= \{ \} \end{aligned}$$

Le même exemple avec reconnaissance de la fin

- Grammaire $G = (\{ ; (;) ; + ; \text{EOF} \}; \{ F ; S ; S' \}; S' ; R)$ où R est

$$\begin{aligned} F &\rightarrow \\ S &\rightarrow (F+S) \\ S &\rightarrow F \\ S' &\rightarrow S \text{ EOF} \end{aligned}$$

- On obtient pour les côtés droits des règles :

$$\begin{aligned} FIRST_1() &= \{ \} \\ FIRST_1((F+S)) &= \{ (\} \\ FIRST_1(F) &= \{ \} \\ FIRST_1(S \text{ EOF}) &= \{ \text{EOF} \} \end{aligned}$$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



Un exemple un peu plus avancé

- ▶ Une grammaire pour les expressions arithmétiques partiellement parenthésées.
- ▶ Terminaux: $\{i; +; *; (;) ; EOF\}$
- ▶ Règles:

$$\begin{aligned} S &\rightarrow E \text{ EOF} \\ E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid i \end{aligned}$$

- ▶ Axiome: S

Un critère pour ne pas être LL(k)

Définition

Une grammaire $G = (V_T; V_N; S; R)$ est *récursive à gauche* s'il y a un non-terminal $N \in V_N$ tel que $N \rightarrow^+ N$ pour un $\alpha \in (V_T \cup V_N)^*$.

- ▶ \rightarrow^+ : dérivation en au moins une étape.
- ▶ Exemple: notre grammaire pour les expressions partiellement parenthésées, car $E \rightarrow E + T$

Lemme

Si la grammaire G est récursive à gauche, alors G n'est pas LL(k), pour aucun $k \in \mathbb{N}$.

L'exemple des expressions partiellement parenthésées

- ▶ Cette grammaire est non-ambiguë ☺
- ▶ Intuition: Les " + " peuvent être reproduites seulement à partir du E. Tout mot engendré par E est protégé par des parenthèses (et).
- ▶ (Il y a aussi une preuve formelle mais je vous en fait grâce.)
- ▶ Cette grammaire, est-elle aussi LL(1)? Ou au moins LL(k) pour quelque $k \in \mathbb{N}$?
- ▶ Elle n'est pas LL(k), pour aucun k ! ☹

Preuve

- ▶ Supposons pour l'absurde que G est LL(k) et récursive à gauche.
- ▶ Il y a donc une dérivation gauche $S \rightarrow^* wX$ (sinon X n'est pas atteignable).
- ▶ Il y a aussi une règle $X \rightarrow X$ (hypothèse simplificatrice), et une règle différente $X \rightarrow \alpha$ (sinon X est non-productif).
- ▶ Il existe donc une dérivation gauche

$$S \rightarrow^* wX \rightarrow^* wX^k$$

- ▶ Par le théorème du début du cours:

$$\text{FIRST}_k(X^{k+1}) \cap \text{FIRST}_k(\alpha) = \emptyset$$

Preuve(2)

- ▶ On a :

$$\text{FIRST}_k(X^{k+1}) \cap \text{FIRST}_k(X^k) = \emptyset$$

- ▶ Donc, grâce à la règle $X \rightarrow \dots$:

$$\text{FIRST}_k(X^{k+1}) \cap \text{FIRST}_k(X^k) = \emptyset$$

- ▶ Contradiction (deux cas : \rightarrow^* ou pas).

Lagrammaire transformée

- ▶ Une grammaire pour les expressions arithmétiques partiellement parenthésées.
- ▶ Terminaux : $\{i, +, *, (,), \text{EOF}\}$
- ▶ Règles :

$$\begin{aligned} S &\rightarrow E \text{ EOF} \\ E &\rightarrow T E' \\ E' &\rightarrow \mid + E \\ T &\rightarrow F T' \\ T' &\rightarrow \mid * T \\ F &\rightarrow (E) \mid i \end{aligned}$$

- ▶ Axiome : S

Quoi faire?

- ▶ On peut transformer la grammaire en une grammaire équivalente (qui définit le même langage), et qui est LL(1).
- ▶ C'est toujours possible, mais il y a deux inconvénients :
 - └ la grammaire résultante peut être plus grande ;
 - └ la structure de l'arbre de dérivation peut changer.
- ▶ Il y a un troisième inconvénient : la transformation peut introduire des règles $N \rightarrow \dots$, il faut donc adapter la technique à ces cas.

Explication de la transformation

- ▶ Les deux règles originales pour le non-terminal E :

$$\begin{aligned} E &\rightarrow T \\ E &\rightarrow E + T \end{aligned}$$

- ▶ Dans la grammaire d'origine, le non-terminal E engendre une séquence non-vide de non-terminaux T, séparés par des +.
- ▶ Dans la nouvelle grammaire, le non-terminal E engendre la suite de cette séquence après un T :

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow \mid + E \end{aligned}$$

- ▶ Pareil pour le non-terminal T.

But d'une analyse d'une grammaire

- ▶ Détecter des propriétés des non-terminaux dans une grammaire, et des grammaires.
- ▶ Exemple: non-terminaux atteignables ou pas, productifs ou pas, pouvant engendrer le mot vide ou pas? Grammaires $LL(1)$, $LL(k)$ ou pas?
- ▶ Défis: récurrence entre les non-terminaux d'une grammaire.
- ▶ Dans certains cas on maîtrise une descente récursive *sur un mot donné* pour la construction d'une dérivation-cas des grammaires $LL(k)$.
- ▶ Dans une analyse de la grammaire, on a pas de mot d'entrée fixé on némaison s'intéresse à une propriété générale de la grammaire.

Calcul d'un point fixe

- ▶ Notre définition de productivité est correcte, mais il faut organiser le calcul différemment:
 - ┌ Au début, on pose que tout non-terminal N pour lequel existe une règle $N \rightarrow w$, où $w \in V_T^*$, est productif.
 - ┌ Puis, s'il y a une règle $N \rightarrow \dots$, où on a déjà reconnu tous les non-terminaux dans \dots comme étant productifs, alors on pose que N est aussi productif.
 - ┌ On s'arrête si on ne peut plus ajouter d'information de cette façon.
- ▶ Ils'agit d'un point fixe!

Illustration du problème d'une descente récursive

- ▶ Un non-terminal N est *productif* s'il y a une règle $N \rightarrow \dots$ telle que tous les non-terminaux dans \dots sont productifs. (C'est en particulier vrai quand \dots ne contient que des terminaux.)
- ▶ Imaginez les règles suivantes de la grammaire:
$$\begin{aligned} A &\rightarrow BC \mid \\ B &\rightarrow Cb \\ C &\rightarrow Ac \end{aligned}$$
- ▶ Il faut éviter une descente récursive qui mène dans une boucle:
$$\begin{aligned} A \text{ productif?} &\rightarrow B \text{ productif?} \rightarrow C \text{ productif?} \\ &\rightarrow A \text{ productif?} \rightarrow \dots \end{aligned}$$

Points fixes

- ▶ En maths: un *point fixe* d'une fonction $f: D \rightarrow D$ est une valeur $x \in D$ telle que $f(x) = x$.
- ▶ Application à l'algorithmique:
 - ┌ le domaine D est l'ensemble de toutes les affectations possibles à des variables d'intérêt. Dans l'exemple: toutes les affectations possibles de P du type $V_N \rightarrow \text{boolean}$.
 - ┌ la fonction f est une mise à jour des variables. Dans l'exemple: propagation de l'information de productivité d'un non-terminal à un autre.
 - ┌ on ne cherche pas un point fixe quelconque, mais on commence avec une valeur initiale, puis on applique f jusqu'à un point fixe.



Calcul d'un point fixe

Nous utiliserons dans nos algorithmes la construction "do...until $X_1, \dots, X_n \text{ fix}$ ". Exemple productivité:

```
forall  $N \in V_N$  :  
  if exists  $(N \rightarrow \alpha) \in R$  tel que  $\alpha \in V_T^*$   
  then  $P(N) = \text{true}$   
  else  $P(N) = \text{false}$   
do  
  forall  $(N \rightarrow \alpha) \in R$  :  
    if  $P(M)$  for all non  $\alpha$  — terminals Min  
    then  $P(N) = \text{true}$   
until P fix
```

Autres exemples symboles atteignables

- ▶ Un non-terminal N est **atteignable** (à partir de l'axiome S) s'il existe une dérivation $S \rightarrow^* N$.
- ▶ Calcul des non-terminals atteignables: transparents suivant.
- ▶ Une grammaire est **réduite** si
 - ┆ toutes les non-terminals sont productives et
 - ┆ toutes les non-terminals sont atteignables.
- ▶ Dans la suite nous supposons que toutes les grammaires sont réduites.

La construction do ... until ... fix

Le code

```
// code initialisation de X  
do  
  // code mise à jour de X  
until X fix
```

peut être traduit vers:

```
// code initialisation de X  
do  
  Xold = X  
  // code mise à jour de X  
while (X != Xold)
```

Calcul des non-terminals atteignables

Algorithme

Donnée une grammaire $G = (V_T; V_N; S; R)$.

```
forall  $N \in V_N$  :  
   $A(N) = \text{false}$   
 $A(S) = \text{true}$   
do  
  forall  $(N \rightarrow \alpha) \in R$  :  
    if  $A(N)$   
    then for all non  $\alpha$  — terminals Min :  
       $A(M) = \text{true}$   
until A fix
```

Lemma

Pour tout $N \in V_N$: $A(N) == \text{true}$ ssi N est atteignable.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Non-terminaux qui peuvent produire ϵ

Algorithme

Donnée une grammaire $G = (V_T; V_N; S; R)$.

```
EPS = { N ∈ V_N | (N → ε) ∈ R }
do
  EPS = EPS ∪ { N ∈ V_N | (N → N_1 :: N_n) ∈ R;
                        N_1, ..., N_n ∈ EPS }
until EPS fix
```

Lemme

Pour tout $N \in V_N$: $N \in \text{EPS}$ si et seulement si $N \rightarrow^* \epsilon$.

Calcul de EPS sur l'exemple

$$\begin{array}{lll} E \rightarrow & TE' & E' \rightarrow \mid +E \\ T \rightarrow & FT' & T' \rightarrow \mid *T \\ F \rightarrow & (E) \mid i & S \rightarrow E \text{ EOF} \end{array}$$

- Initialisation: $\text{EPS} = \{E'; T'\}$
- C'est déjà un point fixe!

Comment calculer FIRST_1 dans le cas général?

- Imaginez une règle $A \rightarrow BC \mid d \mid E$
- Tous les non-terminaux peuvent a priori engendrer ϵ .
- Si $B \notin \text{EPS}$: dans cette règle, seulement $\text{FIRST}_1(B)$ peut contribuer à $\text{FIRST}_1(A)$.
- Si $B \in \text{EPS}$: $\text{FIRST}_1(C)$ peut aussi contribuer à $\text{FIRST}_1(A)$.
- Si $B \in \text{EPS}$ et $C \in \text{EPS}$: d doit être dans $\text{FIRST}_1(A)$.
- Dans aucun des cas, $\text{FIRST}_1(E)$ ne peut contribuer car il se trouve derrière le terminal d .

Le calcul de FIRST_1 dans le cas général

- Pour des raisons techniques, il est plus facile de calculer d'abord une variante de FIRST_1 sans ϵ :

$$Fi(N) = \{a \in V_T \mid N \rightarrow^* aw; w \in V_T^*\}$$

- Puis on en obtient $\text{FIRST}_1(N)$ en utilisant EPS que nous avons déjà calculé.



Calculde $FIRST_1$ danslecasgénéral

Algorithme

Donnéeunegrammaire $G = (V_T; V_N; S; R)$.

pour tout $N \in V_N$:
 $Fi(N) = \{a \in V_T \mid (N \rightarrow N_1 \dots N_n a) \in R; N_1, \dots, N_n \in EPS\}$

do

pour tout $(N \rightarrow N_1 \dots N_n M) \in R$
 tel que $M \in V_N; N_1, \dots, N_n \in EPS$:

$Fi(N) = Fi(N) \cup Fi(M)$

until Ffix

Lemme

Pour tout $N \in V_N : Fi(N) = FIRST_1(N) - \{ \}$

Calculde $FIRST_1$ danslecasgénéral

Pour tout $N \in V_N$:

$$FIRST_i(N) = \begin{cases} Fi(N) \cup \{ \} & \text{si } N \in EPS \\ Fi(N) & \text{sinon} \end{cases}$$

Calculde Fisurl'exemple

$$\begin{aligned} E &\rightarrow TE' & E' &\rightarrow \mid +E \\ T &\rightarrow FT' & T' &\rightarrow \mid *T \\ F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\ EPS &= \{E'; T'\} \end{aligned}$$

Calcul

	Initial	Iter1	Iter2	Iter3
S	\emptyset	\emptyset	\emptyset	$\{i; (\}$
E	\emptyset	\emptyset	$\{i; (\}$	$\{i; (\}$
E'	$\{+\}$	$\{+\}$	$\{+\}$	$\{+\}$
T	\emptyset	$\{i; (\}$	$\{i; (\}$	$\{i; (\}$
T'	$\{*\}$	$\{*\}$	$\{*\}$	$\{*\}$
F	$\{i; (\}$	$\{i; (\}$	$\{i; (\}$	$\{i; (\}$

Calculde $FIRST_1$ surl'exemple

$$\begin{aligned} E &\rightarrow TE' & E' &\rightarrow \mid +E \\ T &\rightarrow FT' & T' &\rightarrow \mid *T \\ F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\ EPS &= \{E'; T'\} \end{aligned}$$

Calcul

	Fi	$FIRST_1$
S	$\{i; (\}$	$\{i; (\}$
E	$\{i; (\}$	$\{i; (\}$
E'	$\{+\}$	$\{+\}$
T	$\{i; (\}$	$\{i; (\}$
T'	$\{*\}$	$\{*\}$
F	$\{i; (\}$	$\{i; (\}$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Calcul de $FIRST_1$ dans le cas général

On étend maintenant $FIRST_1$ à des mots déterminaux et non-terminaux :

- ▶ $FIRST_1(\epsilon) = \{\epsilon\}$
- ▶ pour tout $a \in V_T$: $FIRST_1(a) = \{a\}$
- ▶ pour tout $N \in V_N$:

$$FIRST_1(N) = \begin{cases} FIRST_1(N) & \text{si } N \notin EPS \\ (FIRST_1(N) \setminus \{\epsilon\}) \cup FIRST_1(\epsilon) & \text{si } N \in EPS \end{cases}$$

La fonction $FOLLOW_k$

Définition

Soit $G = (V_T; V_N; S; R)$ une grammaire, $k \in \mathbb{N}$. La fonction $FOLLOW_k: V_N \rightarrow 2^{V_T^*}$ est définie par

$$FOLLOW_k(N) = \{w \mid S \xrightarrow{*} N ; ; \in (V_T \cup V_N)^*; w \in FIRST_k(\epsilon)\}$$

Remarques

- ▶ $FOLLOW_k(N)$ est l'ensemble de tous les mots déterminaux de longueur k qui peuvent, dans des mots de $\mathcal{L}(G)$, suivre à un mot dérivé de N .

Nous avons besoin de plus d'information !

- ▶ Le calcul de $FIRST_1$ n'est plus suffisant pour savoir quelle production appliquer !
- ▶ Exemple : $E' \rightarrow \mid +E$
Si nous voyons $+$ alors il faut utiliser la deuxième alternative pour réécrire E' . Mais quand faut-il appliquer la première ?
- ▶ Il nous manque une information : quels sont les symboles terminaux qui peuvent *suivre* à un mot produit par un non-terminal ?

Calcul de $FOLLOW_1$

Algorithme

Donnée une grammaire $G = (V_T; V_N; S; R)$.

```

for all  $N \in V_N$  :
   $Fo(N) = \{a \in V_T \mid (M \rightarrow \dots NN_1 \dots N_k a \dots) \in R; N_1, \dots, N_n \in EPS\}$ 
  for all  $(M \rightarrow \dots NN_1 \dots N_k N' \dots) \in R$ 
  with  $N_1, \dots, N_k \in EPS$  :
     $Fo(N) = Fo(N) \cup Fi(N')$ 
  do
    for all  $(M \rightarrow \dots NN_1 \dots N_k) \in R$ 
    with  $N_1, \dots, N_n \in EPS$  :
       $Fo(N) = Fo(N) \cup Fo(M)$ 
until  $Fo$  fix
  
```

Lemme

Si $(S \rightarrow N EOF) \in R$, alors $\forall N \in V_N : Fo(N) = FOLLOW_1(N)$



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Calcul de $FOLLOW_1$ sur l'exemple

$$\begin{aligned} E &\rightarrow TE' & E' &\rightarrow \mid +E \\ T &\rightarrow FT' & T' &\rightarrow \mid *T \\ F &\rightarrow (E) \mid i & S &\rightarrow E EOF \end{aligned}$$

$$EPS = \{E'; T'\} \quad Fi(E') = \{+\} \quad Fi(T') = \{*\}$$

Calcul

	Initial	Iter1	Iter2
S	\emptyset	\emptyset	\emptyset
E	$\{EOF; \}$	$\{EOF; \}$	$\{EOF; \}$
E'	\emptyset	$\{EOF; \}$	$\{EOF; \}$
T	$\{+\}$	$\{+; EOF; \}$	$\{+; EOF; \}$
T'	\emptyset	$\{+\}$	$\{EOF; \}; +\}$
F	$\{*\}$	$\{*; EOF; \}; +\}$	$\{*; EOF; \}; +\}$

Comment choisir la règle dans l'analyse lexicale

Soit $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ une alternative. Il y a deux cas :

1. Soit, aucun des $FIRST_1(\alpha_i)$ ne contient ϵ :
comme avant :
 - └ On choisit la règle $A \rightarrow \alpha_i$ quand les symboles suivants sont dans $FIRST_1(\alpha_i)$ (ils sont tous disjoints).
 - └ Erreurs si aucun tel α_i existe
2. Soit $\epsilon \in FIRST_1(\alpha_i)$:
 - └ si les symboles suivants sont dans $FIRST_1(\alpha_j)$: choisir $A \rightarrow \alpha_j$, pour $1 \leq j \leq n$.
 - └ si les symboles suivants sont dans $FOLLOW_1(A)$: choisir $A \rightarrow \alpha_i$.
 - └ sinon Erreur.

Enfin (!) : le critère pour être LL(1)

Théorème

La grammaire $G = (V_T; V_N; S; R)$ est LL(1) si pour toutes les alternatives $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$:

- ▶ $FIRST_1(\alpha_1); \dots; FIRST_1(\alpha_n)$ sont disjoints entre eux;
- ▶ Si $\epsilon \in FIRST_1(\alpha_i)$, alors pour tous $j \neq i$:

$$FIRST_1(\alpha_j) \cap FOLLOW_1(A) = \emptyset$$

Remarque

Condition (1) implique qu'au plus un des ensembles $FIRST_1(\alpha_i)$ contient ϵ .



Analyse de Données Structurées - Cours 7

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

18 mars 2015

© Ralf Treinen 2015

Rappel: Calcul de Fi

Algorithme

Donnée une grammaire $G = (V_T, V_N, S, R)$.

pour tout $N \in V_N$:
 $Fi(N) = \{a \in V_T \mid N \rightarrow N_1 \dots N_n a \alpha \in R, N_1, \dots, N_n \in EPS\}$

do
pour tout $N \rightarrow N_1 \dots N_n M \alpha \in R$
tel que $M \in V_N, N_1, \dots, N_n \in EPS$:
 $Fi(N) = Fi(N) \cup Fi(M)$
until Fi fixe

Vu au dernier cours

- ▶ Calcul d'un point fixe par $do \dots until$ X fixe
- ▶ Définition: $Fi(N) = \{a \in V_T \mid N \rightarrow^* a \alpha\}$
- ▶ On a que $Fi(N) = FIRST_1(N) \setminus \{\epsilon\}$
- ▶ Algorithme pour le calcul de $Fi(N)$ en présence de règles $M \rightarrow \epsilon$
- ▶ Reste à faire: calcul de $FIRST_1(\alpha)$ pour $\alpha \in V_T \cup V_N^*$

Exemple

Grammaire

$$\begin{array}{ll} A \rightarrow \epsilon \mid a & B \rightarrow \epsilon \mid Ab \\ C \rightarrow ABc & D \rightarrow BCd \end{array}$$

Non-terminaux qui peuvent produire ϵ

$EPS = \{A, B\}$

Calcul

	Initial	Iter1	Iter2
A	{ }	{ }	{ }
B	{b}	{b}	{b, c}
C	{c}	{b, c}	{b, c}
D	{ }	{b, c}	{b, c}

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Calculde $FIRST_1$ danslecasgénéral

Pourtout $N \in V_N$:

$$FIRST_i N = \begin{cases} Fi N \cup \{\epsilon\} & \text{si } N \in EPS \\ Fi N & \text{sinon} \end{cases}$$

Calculde $FIRST_1$ danslecasgénéral

Onétendmaintenant $FIRST_1$ àdesmotsdeterminauxet non-terminaux:

- $FIRST_1 \epsilon = \{\epsilon\}$
- pourtout $a \in V_T$: $FIRST_1 a\alpha = \{a\}$
- pourtout $N \in V_N$:

$$FIRST_1 N\alpha = \begin{cases} FIRST_1 N & \text{si } N \notin EPS \\ FIRST_1 N \setminus \{\epsilon\} \cup FIRST_1 \alpha & \text{si } N \in EPS \end{cases}$$

Calculde $FIRST_1$ surl'exemple

$$\begin{aligned} E &\rightarrow TE' & E' &\rightarrow \epsilon \mid +E \\ T &\rightarrow FT' & T' &\rightarrow \epsilon \mid *T \\ F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\ & & EPS & \{E', T'\} \end{aligned}$$

Calcul

	Fi	$FIRST_1$
S	{i, (}	{i, (}
E	{i, (}	{i, (}
E'	{+}	{+, }
T	{i, (}	{i, (}
T'	{*}	{*, }
F	{i, (}	{i, (}

Calculde $FIRST_1$ descôtésdroitesdansl'exemple

$$\begin{aligned} FIRST_1 TE' &= \{i, (\} \\ FIRST_1 \epsilon &= \{\epsilon\} \\ FIRST_1 + E &= \{+\} \\ FIRST_1 FT' &= \{i, (\} \\ FIRST_1 * T &= \{*\} \\ FIRST_1 (E) &= \{(\} \\ FIRST_1 i &= \{i\} \\ FIRST_1 S EOF &= \{i, (\} \end{aligned}$$

Nousavonsbesoindeplused'information!

- Le calcul de $FIRST_1$ n'est plus suffisant pour savoir quelle production appliquer!
- Exemple: $E' \rightarrow \epsilon \mid +E$
Sinous voyons $+$ alors il faut utiliser la deuxième alternative pour écrire E' . Mais quand faut-il appliquer la première?
- Il nous manque une information: quels sont les symboles terminaux qui peuvent *suivre* à un mot produit par un non-terminal?

Calcul de $FOLLOW_1$

Algorithme

Donnée une grammaire $G = (V_T, V_N, S, R)$.

```

for all  $N \in V_N$ :
   $Fo(N) = \{a \in V_T \mid M \rightarrow \dots NN_1 \dots N_k a \dots \in R, N_1, \dots, N_n \in EPS\}$ 
for all  $M \rightarrow \dots NN_1 \dots N_k N^0 \dots \in R$ 
with  $N_1, \dots, N_k \in EPS$ :
   $Fo(N) = Fo(N) \cup Fi(N^0)$ 
do
  for all  $M \rightarrow \dots NN_1 \dots N_k \in R$ 
  with  $N_1, \dots, N_n \in EPS$ :
     $Fo(N) = Fo(N) \cup Fo(M)$ 
until  $Fo$  fix
  
```

Lemme

Si $S \rightarrow N EOF \in R$, alors $\forall N \in V_N : Fo(N) = FOLLOW_1 N$

La fonction $FOLLOW_k$

Définition

Soit $G = (V_T, V_N, S, R)$ une grammaire, $k \in \mathbb{N}$. La fonction $FOLLOW_k : V_N \rightarrow 2^{V_T^*}$ est définie par

$$FOLLOW_k N = \{w \mid S \xrightarrow{*} \beta N \gamma, \beta, \gamma \in V_T \cup V_N^*, w \in FIRST_k \gamma\}$$

Remarques

- $FOLLOW_k N$ est l'ensemble de tous les mots déterminaux de longueur k qui peuvent, dans des mots de $\mathcal{L}(G)$, suivre à un mot dérivé de N .

Calcul de $FOLLOW_1$ sur l'exemple

$$\begin{aligned}
 E &\rightarrow TE' & E' &\rightarrow \epsilon \mid +E \\
 T &\rightarrow FT' & T' &\rightarrow \epsilon \mid *T \\
 F &\rightarrow (E) \mid i & S &\rightarrow E EOF \\
 EPS &= \{E', T'\} & Fi(E') &= \{+\} & Fi(T') &= \{*\}
 \end{aligned}$$

Calcul

	Initial	Iter1	Iter2
S	\emptyset	\emptyset	\emptyset
E	$\{EOF,)\}$	$\{EOF,)\}$	$\{EOF,)\}$
E'	\emptyset	$\{EOF,)\}$	$\{EOF,)\}$
T	$\{+\}$	$\{+, EOF,)\}$	$\{+, EOF,)\}$
T'	\emptyset	$\{+\}$	$\{EOF,), +\}$
F	$\{*\}$	$\{*, +\}$	$\{*, EOF,)\}$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Finalement(!): le critère pour être LL(1)

Théorème

La grammaire $G = (V_T, V_N, S, R)$ est LL(1) si pour toutes les alternatives $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$:

- ▶ $FIRST_1 \alpha_1, \dots, FIRST_1 \alpha_n$ sont disjoints entre eux;
- ▶ Si $\epsilon \in FIRST_1 \alpha_i$, alors pour tous $j \neq i$:

$$FIRST_1 \alpha_j \cap FOLLOW_1 A = \emptyset$$

Remarque

Condition (1) implique qu'au plus un des ensembles $FIRST_1 \alpha_i$ contient ϵ .

Le critère sur l'exemple

Non-terminal	Cas1	$FIRST_1$	Cas2	$FIRST_1$
E	TE'	{i, (}		
E'	ϵ	{ ϵ }	+ E	{+}
T	FT'	{i, (}		
T'	ϵ	{ ϵ }	* T	{*}
F	(E)	{(}	i	{i}
S	E EOF	{i, (}		

- ▶ $FOLLOW_1 E' = \{), EOF\}$ disjoint avec $\{+\}$ ☺
- ▶ $FOLLOW_1 T' = \{EOF,), +\}$ disjoint avec $\{*\}$ ☺
- ▶ $\{(}$ disjoint avec $\{i\}$ ☺
- ▶ Conclusion : la grammaire est LL(1)!

Comment choisir la règle dans l'analyse syntaxique

Soit $A \rightarrow \alpha_1 \mid \dots \mid \alpha_n$ une alternative. Il y a deux cas :

1. Soit, aucun des $FIRST_1 \alpha_i$ ne contient ϵ :
comme avant :
 - ▶ On choisit la règle $A \rightarrow \alpha_i$ quand les symboles suivants sont dans $FIRST_1 \alpha_i$ (ils sont tous disjoints).
 - ▶ Erreur si aucun tel i existe
2. Soit α_i avec $\epsilon \in FIRST_1 \alpha_i$:
 - ▶ si les symboles suivants sont dans $FOLLOW_1 A$: choisir $A \rightarrow \alpha_j$, pour $1 \leq j \leq n$.
 - ▶ si les symboles suivants sont dans $FOLLOW_1 A$: choisir $A \rightarrow \alpha_j$.
 - ▶ sinon Erreur.

Le code de l'analyseur de syntaxe sur l'exemple

```
import java.io.*;

/* LL(1) parser for expressions with priorities, */
/* terminated by the symbol '$' */
class Parser {

    protected LookAhead1Reader reader;

    public Parser(LookAhead1Reader r) {
        reader = r;
    }

    public void term(char c) throws IOException {
        /* consume the character c */
    }
}
```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Lecodedul'analyseurdesyntaxesurl'exempleII

```
        reader.eat(c);
    }

    public void      nonterm_S()      throws      IOException,ReadException{
        /* S ->E$ */
        this .nonterm_E();
        this .term('$');
    }

    public void      nonterm_E()      throws      IOException,ReadException{
        /* E ->TE1 */
        this .nonterm_T();
        this .nonterm_E1();
    }
}
```

Lecodedul'analyseurdesyntaxesurl'exempleIV

```
    public void      nonterm_T1()      throws      IOException,ReadException{
        if (reader.check('$')||reader.check('')||reader.check('+')){
            /* T1 ->epsilon */
        } elseif (reader.check(' * ')){
            /* T1 -> * T */
            this .term(' * ');
            this .nonterm_T();
        }
    }

    public void      nonterm_F()      throws      IOException,ReadException{
        if (reader.check('(')){
            /* F ->(E) */
            this .term('(');
        }
    }
}
```

Lecodedul'analyseurdesyntaxesurl'exempleIII

```
    public void      nonterm_E1()      throws      IOException,ReadException{
        if (reader.check('$')||reader.check('')){
            /* E1 ->epsilon */
        } elseif (reader.check('+')){
            /* E1 ->+E */
            this .term('+');
            this .nonterm_E();
        }
    }

    public void      nonterm_T()      throws      IOException,ReadException{
        /* T ->FT1 */
        this .nonterm_F();
        this .nonterm_T1();
    }
}
```

Lecodedul'analyseurdesyntaxesurl'exempleV

```
        this .nonterm_E();
        this .term('');
    } elseif (reader.check('i')){
        /* F ->i */
        this .term('i');
    }
}
}
```



Syntaxe concrète

- ▶ **Syntaxe concrète** : c'est la définition de la forme correcte de l'écriture (représentation textuelle). Elle est normalement définie par des expressions régulières, et une grammaire hors-contexte.
- ▶ La vérification qu'un texte d'entrée suit la syntaxe concrète d'un langage est réalisée en coopération par l'analyse lexicale et l'analyse syntaxique.
- ▶ Or, normalement on ne veut pas seulement savoir si l'entrée est correcte ou pas, mais aussi faire quelque chose avec, par exemple :
 - ▶ Document XML représentant des données géographiques : afficher une carte.
 - ▶ Programme écrit dans un langage de programmation : l'exécuter (cas d'un interpréteur), ou engendrer du code exécutable (cas d'un compilateur).

Abstraction

- ▶ Quels sont les détails de la syntaxe concrète dont on peut faire abstraction : ça dépend de ce qu'on compte faire avec !
- ▶ Nous avons déjà vu certaines abstractions faites par l'analyse lexicale, dans le contexte des langages de programmation :
 - ▶ les espaces (mais attention aux cas où l'incrémention des lignes indique la structure du programme)
 - ▶ les commentaires (mais attention aux cas où les commentaires sont importantes pour la suite, par exemple on écrit un **prettyprinter** de code source)
- ▶ Dans un premier temps on peut utiliser l'arbre de dérivation construit par l'analyse syntaxique comme syntaxe abstraite.

- ▶ On va donc construire une représentation de la structure que les analyses lexicales et syntaxiques ont découvertes : c'est la **syntaxe abstraite**.
- ▶ Le passage de la syntaxe concrète à la syntaxe abstraite a deux fonctions :
 - ▶ Vérifier que l'entrée correspond aux règles de la syntaxe concrète,
 - ▶ si c'est le cas, construire une représentation en forme de la syntaxe abstraite.
- ▶ Il y a ici souvent une **abstraction** : certains détails de la représentation textuelle (en syntaxe concrète) sont omis dans la syntaxe abstraite car ils ne sont pas pertinents pour la suite.
- ▶ On essaye d'abstraire de tout ce qui n'est pas nécessaire pour la suite : simplifier tant que possible !

Arbre de dérivation vs. Syntaxe abstraite

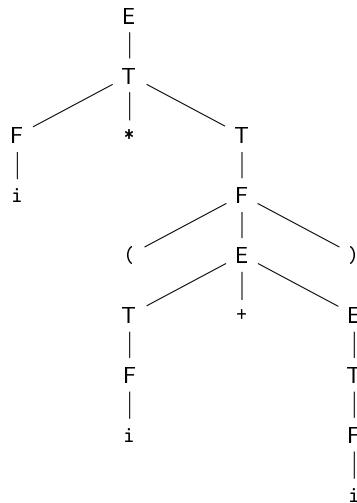
- ▶ En général, l'arbre de dérivation contient trop d'information.
- ▶ Raison : l'arbre de dérivation dit **comment** la structure a été découverte dans le texte d'entrée. Or, c'est seulement la structure elle-même qui est pertinente pour la suite.
- ▶ Exemple : grammaire

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow (E) \mid i \end{aligned}$$

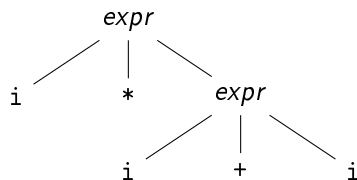
- ▶ Arbre de dérivation pour $i * i$



Arbre de dérivation pour $i * (i + i)$



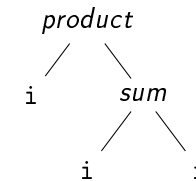
Autre possibilité pour la syntaxe abstraite



- Il y aura une seule classe, avec l'opérateur comme argument supplémentaire du constructeur.
- Solution à préférer si dans la suite tous les opérateurs sont traités de la même façon.

Quelle information retenir dans la syntaxe abstraite?

- La distinction entre E, T et F n'est plus pertinente.
- Les parenthèses ne sont plus pertinentes non plus : elles servent à indiquer la structure dans la syntaxe concrète, mais une fois la structure découverte elles ne servent plus à rien.
- Structure à retenir :



Définition de la Syntaxe Abstraite sur l'exemple

```
abstract class Expression {
    public boolean isIdentifier() {
        return false;
    }
    public boolean isSum() {
        return false;
    }
    public boolean isProduct() {
        return false;
    }
    abstract void print();
}
```



DéfinitiondelaSyntaxeAbstraitesurl'exempleII

```
class Identifier extends Expression{
    public boolean isIdentifier(){
        return true ;
    }
    public Identifier(){
    }
    public void print(){
        System.out.print("I");
    }
}

class Sum extends Expression{
    private Expression left;
    private Expression right;
    public boolean isSum(){
```

DéfinitiondelaSyntaxeAbstraitesurl'exempleIV

```
class Product extends Expression{
    private Expression left;
    private Expression right;
    public boolean isProduct(){
        return true ;
    }
    public Product(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }

    public void print(){
        System.out.print("Product");
        left.print();
        System.out.print(',');
    }
}
```

DéfinitiondelaSyntaxeAbstraitesurl'exempleIII

```
        return true ;
    }
    public Sum(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }
    public void print(){
        System.out.print("Sum");
        left.print();
        System.out.print(',');
        right.print();
        System.out.print(")");
    }
}
```

DéfinitiondelaSyntaxeAbstraitesurl'exempleV

```
        right.print();
        System.out.print(")");
    }
}
```



Constructionde la Syntaxe Abstraite sur l'exemple I

```
import java.io. *;

/* LL(1) parser for expressions with priorities, */
/* terminated by the symbol '$' */
class Parser{

    protected LookAhead1Reader reader;

    public Parser(LookAhead1Reader r){
        reader=r;
    }

    public void term( char c) throws IOException, ReadException{
        /* consume the character c */
    }
}
```

Constructionde la Syntaxe Abstraite sur l'exemple II

```
Expression e2= this.nonterm_E1();
if (e2== null){
    return e1;
} else {
    return new Sum(e1,e2);
}

public Expression nonterm_E1()
throws IOException, ReadException, ParseException{
    if (reader.check('$') || reader.check('')){
        /* E1 -> epsilon */
        return null;
    } elseif (reader.check('+')){
        /* E1 -> +E */
    }
}
```

Constructionde la Syntaxe Abstraite sur l'exemple II

```
        reader.eat(c);
    }

    public Expression nonterm_S()
    throws IOException, ReadException, ParseException{
        /* S -> E$ */
        Expression e= this.nonterm_E();
        this.term('$');
        return e;
    }

    public Expression nonterm_E()
    throws IOException, ReadException, ParseException{
        /* E -> TE1 */
        Expression e1= this.nonterm_T();
    }
}
```

Constructionde la Syntaxe Abstraite sur l'exemple IV

```
        this.term('+');
        return this.nonterm_E();
    } else {
        throw new ParseException("cannot reduce E");
    }
}

public Expression nonterm_T()
throws IOException, ReadException, ParseException{
    /* T -> FT1 */
    Expression e1= this.nonterm_F();
    Expression e2= this.nonterm_T1();
    if (e2== null){
        return e1;
    } else {
    }
}
```



Constructionde la Syntaxe Abstraite sur l'exemple V

```
        return new Product(e1,e2);
    }
}

public Expression nonterm_T1()
    throws IOException, ReadException, ParseException {
    if (reader.check('$') || reader.check('(') || reader.check('+')) {
        /* T1 -> epsilon */
        return null;
    } else if (reader.check('*')) {
        /* T1 -> * T */
        this.term('*');
        return this.nonterm_T();
    } else {
        throw new ParseException("cannot reduce T");
    }
}
```

Constructionde la Syntaxe Abstraite sur l'exemple VI

```
    } else {
        throw new ParseException("cannot reduce F");
    }
}
```

Constructionde la Syntaxe Abstraite sur l'exemple VI

```
    }
}

public Expression nonterm_F()
    throws IOException, ReadException, ParseException {
    if (reader.check('(')) {
        /* F -> (E) */
        this.term('(');
        Expression e = this.nonterm_E();
        this.term('');
        return e;
    } else if (reader.check('i')) {
        /* F -> i */
        this.term('i');
        return new Identifier();
    }
}
```



Analyse de Données Structurées- Cours 8

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

25 mars 2015

© Ralf Treinen 2015

Le format JSON

- ▶ **JavaScript Object Notation**
- ▶ Le document de standardisation commence avec le résumé suivant:

JavaScript Object Notation (JSON) is a lightweight, text-based, language-independent data interchange format.

[...]

JSON defines a small set of formatting rules for the portable representation of structured data.

Formats génériques pour des données structurées

- ▶ Permettent la représentation textuelle de données structurées
- ▶ Utiles pour:
 - | Stocker des données dans un fichier: persistance
 - | Transmission de données entre des applications hétérogènes (éventuellement écrites dans des langages différents)
 - | Transmission de données sur le web (par exemple, AJAX)
 - | Des fichiers de configuration (car facile à modifier pour des humains, et facile à interpréter pour des programmes)
- ▶ Exemples dans ce cours: JSON et XML

Le format JSON

- ▶ Valeurs de base: `true`, `false`, `null`, valeurs numériques, et chaînes de caractères.
- ▶ Deux mécanismes pour combiner des valeurs:
 - | *array* : c'est simplement une liste de valeurs entre crochets `[]`, séparées par des virgules
 - | *object* : c'est une séquence de paires, chacune consistant en une chaîne de caractères, et une valeur. La séquence est notée entre accolades `{ et }`, les paires sont séparées par des virgules.
- ▶ Attention à la différence entre *terminé* et *séparé* par des virgules.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Exemple

```
{
  "menu":{
    "id":"file",
    "value":"File",
    "popup":{
      "menuitem":[
        {"value":"New","onclick":"CreateNewDoc()"},
        {"value":"Open","onclick":"OpenDoc()"},
        {"value":"Close","onclick":"CloseDoc()"}
      ]
    }
  }
}
```

Source: http://fr.wikipedia.org/wiki/JavaScript_Object_Notation

Ladéfinitionofficielle

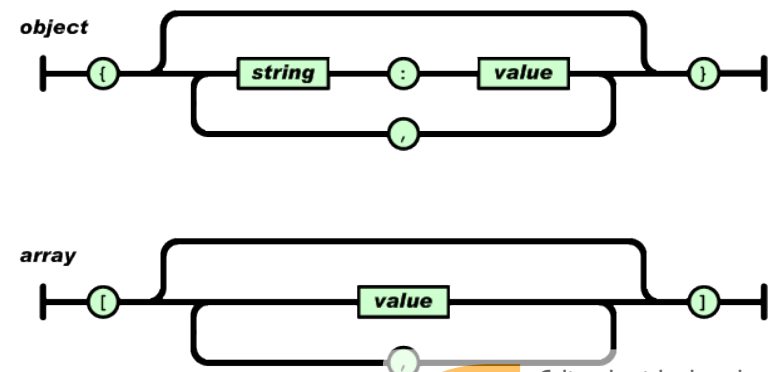
- ▶ Ontrouvesurle“siteofficiel”deJSON, <http://json.org/>, unesériededidiagrammesdesyntaxe.
- ▶ Ilyaaussiunepropositiondestandardisationdelapartdela IETF(InternetEngineeringTaskForce)“RFC7159”,disponible àl'adresse <http://tools.ietf.org/html/rfc7159>.

Unautreexemple

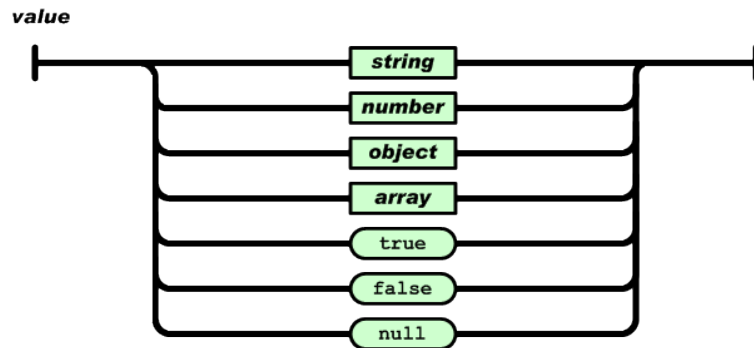
```
{
  "Image":{
    "Width":800,
    "Height":600,
    "Title": "View from 15th Floor",
    "Thumbnail":{
      "Url": "http://www.example.com/image/481989",
      "Height":125,
      "Width":100
    },
    "Animated":false,
    "IDs":[116,943,234,38793]
  }
}
```

Source: <http://tools.ietf.org/html/rfc7159>

Diagrammesdesyntaxe

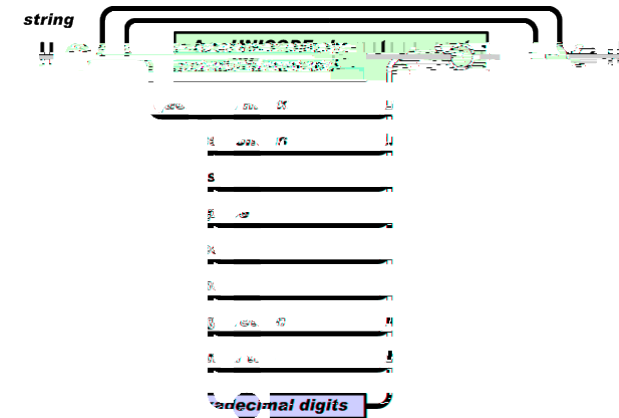


Diagrammesdesyntaxe



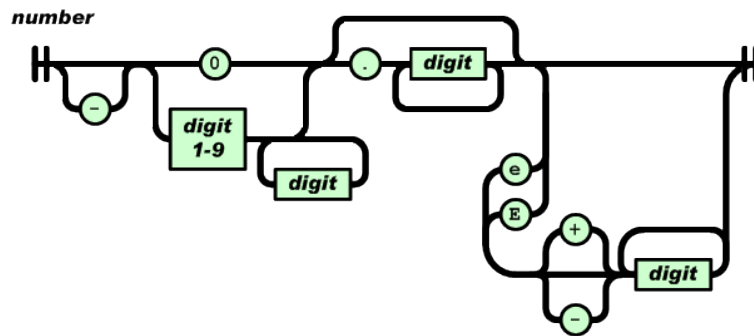
Source: <http://json.org/>

Diagrammesdesyntaxe



Source: <http://json.org/>

Diagrammesdesyntaxe



Source: <http://json.org/>

Commentfaireanalyselexicale/syntaxique?

- ▶ Expressionsrégulièrespourleschaînesetlesvaleurs numériques.
- ▶ Jetons(délivrésparl'analyselexicale):
 - | string, number,(munisdevaleurs)
 - | true, false, null,
 - | , , : , { , } , [,]



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Lagammaire (première tentative)

- ▶ $V_N = \{\text{Value}, \text{Object}, \text{Array}, \text{Aseq}, \text{Oseq}, \text{AseqNV}, \text{OseqNV}\}$
- ▶ Axiome: value
- ▶ Règles:
 - Value \rightarrow string | number | Object | Array | true | false | null
 - Array \rightarrow [Aseq]
 - Aseq \rightarrow ϵ | AseqNV
 - AseqNV \rightarrow Value | Value , AseqNV
 - Object \rightarrow { Oseq }
 - Oseq \rightarrow ϵ | OseqNV
 - OseqNV \rightarrow string: Value | string: Value , Oseq
- ▶ Est-ce LL(1)?

Analyser un document écrit en JSON

- ▶ Il est donc assez facile d'écrire un analyseur syntaxique pour le format JSON.
- ▶ Il existe aussi de nombreuses bibliothèques, pour presque toutes les langages de programmation (voir la liste sur json.org).

Lagammaire (deuxième tentative)

- ▶ $V_N = \{\text{Value}, \text{Object}, \text{Array}, \text{Aseq}, \text{Oseq}, \text{Aseq}', \text{Oseq}'\}$
- ▶ Axiome: value
- ▶ Règles:
 - Value \rightarrow string | number | Object | Array | true | false | null
 - Array \rightarrow [Aseq]
 - Aseq \rightarrow ϵ | Value Aseq'
 - Aseq' \rightarrow ϵ | , Value Aseq'
 - Object \rightarrow { Oseq }
 - Oseq \rightarrow ϵ | string: Value Oseq'
 - Oseq' \rightarrow ϵ | , string: Value Oseq'
- ▶ Est-ce maintenant LL(1)?

Leformat XML

- ▶ Extensible Markup Language (fr.: *langage de balisage extensible*)
- ▶ Objectifs: les mêmes que pour JSON
- ▶ Ils agissent d'une instance d'un format plus général, d'un nom SGML (Standard Generalized Markup Language).
- ▶ Langage plus riche que JSON, mais aussi plus verbeux et moins élégant.
- ▶ Le langage XML est standardisé par le W3C (World Wide Web Consortium).



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

La structure d'un document XML

- ▶ Un document XML décrit un arbre.
- ▶ Un nœud peut avoir un nombre arbitraire d'enfants.
- ▶ Un nœud peut avoir des *attributs*.
- ▶ Les enfants d'un nœud sont appelés *éléments*; un élément peut être un morceau de texte, ou un autre nœud.

Remarques

- ▶ Ça ressemble à du HTML, et pour cause: XML et HTML sont tous les deux des dérivés du langage plus général SGML.
- ▶ L'utilisation des chevrons <et> et assez caractéristique pour ces langages.
- ▶ Le XML est plus strict que le HTML, par exemple:
 - └ XML exige que les valeurs des attributs soient écrites entre guillemets, contrairement aux premières versions de HTML.
 - └ XML exige que toutes les balises soient proprement fermées, contrairement à HTML qui est très libéral en ce regard.

Exemple

Le même exemple que dans la section sur JSON:

```
menuid="file" value="File">
  popup>
    menuitem value="New" onclick="CreateNewDoc()">
    /menuitem>
    menuitem value="Open" onclick="OpenDoc()">
    /menuitem>
    menuitem value="Close" onclick="CloseDoc()">
    /menuitem>
  /popup>
/menu>
```

(sans utiliser l' raccourci pour les nœuds sans enfants)

Version de la définition de la syntaxe

- ▶ Un nœud commence avec une balise (angl.: *tag*) de la forme

`<m>`

et se termine avec

`</m>`

nom est mot quelconque, appelée *nom* de cet élément.

- ▶ Le nom est obligatoirement le même dans la balise de début et de fin du même nœud. Il est permis d'utiliser le même nom pour plusieurs nœuds du même arbre.
- ▶ Une balise de début peut en plus contenir des attributs avec leurs valeurs.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Vers une définition de la syntaxe

- Une paire attribut/valeur est donnée dans la forme

attr= valeur

où *attr* est un mot (appelé l'attribut), et *valeur* est une chaîne de caractères entre guillemets (appelée la valeur).

- Il n'est pas permis de définir dans la même balise deux fois le même attribut.

Exemple

Le même exemple avec raccourci :

```
menuid="file" value="File">
  popup>
    menuitem value="New" onclick="CreateNewDoc()" />
    menuitem value="Open" onclick="OpenDoc()" />
    menuitem value="Close" onclick="CloseDoc()" />
  /popup>
/menu>
```

Source: http://fr.wikipedia.org/wiki/JavaScript_Object_Notation

Raccourci pour des nœuds sans enfants

La syntaxe

```
<nom attr1="value1" ... attrn="valuen"/>
```

est un raccourci pour

```
<nom attr1="value1" ... attrn="valuen">
</nom>
```

Exemple: Données OpenStreetMap

```
<?xml version="1.0" encoding="UTF-8"?>
<osm version="0.6"
  copyright="OpenStreetMap and contributors"
  license="http://opendatacommons.org/licenses/odbl/1-0/">
  <way id="62378611"
    visible="true" version="8" changeset="20691652"
    timestamp="2014-02-21T11:23:38Z" user="thibdrev" uid="1279506">
    <ndref="779143878"/>
    <ndref="2198721646"/>
    <ndref="2198721727"/>
    .....
    <tag k="amenity" v="university"/>
    <tag k="building" v="yes"/>
    <tag k="name" v="Halle aux Farines (Université Paris Diderot)"/>
    <tag k="source" v="cadastre-dgi-fr" source:Direction Générale des Impôts-Cadastre.
    <tag k="wheelchair" v="yes"/>
    <tag k="wikipedia" v="fr:Université Paris VII - Diderot"/>
  </way>
</osm>
```


Une grammaire pour XML?

- ▶ Terminaux (jetons): `mot`, `string`, `texte`, `/`, `=`, `<`, `>`
- ▶ Non-terminaux: `Arbre`, `Start`, `End`, `Attrbs`, `Elements`
- ▶ Axiome: `Arbre`
- ▶ Règles:

`Arbre` \rightarrow `StartElementsEnd`

`Start` \rightarrow `<mot Attrbs >`

`End` \rightarrow `</mot>`

`Attrbs` \rightarrow ϵ | `mot=string Attrbs`

`Elements` \rightarrow ϵ | `ArbreElements` | `texte Elements`

Vers une preuve

- ▶ Nous allons montrer que c'est impossible de définir le langage XML entièrement par une grammaire hors-contexte.
- ▶ Cela implique que c'est également impossible en utilisant une combinaison d'expressions régulières-grammaire hors-contexte, car toute expression régulière peut être transformée en une grammaire.

Une grammaire pour XML?

- ▶ Attention: Cette grammaire ne définit pas complètement le langage XML.
- ▶ Il y a deux éléments de la définition qui ne sont pas exprimés par la grammaire:

- └ Toute balise de début doit être fermée par une balise de fin *du même nom*.

- └ Tous les attributs dans une balise de début doivent être *différents*.

- ▶ Peut-on améliorer la grammaire pour exprimer aussi ces deux restrictions?
- ▶ Réponse: non! voir le transparent suivants.

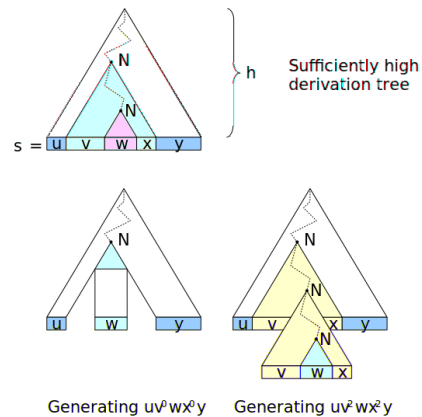
Rappel: lemme de pumping pour les langages algébriques

Soit L un langage algébrique. Alors il existe un $p \geq 0$ tel que tout mot $s \in L$ avec $|s| \geq p$ peut être décomposé en $z = uvwx$ avec:

- ▶ $|vwx| \leq p$
- ▶ $|vx| \geq 1$
- ▶ $uv^iwx^iy \in L$ pour tout $i \geq 0$



Idée de la preuve du lemme de pumping



Source: Jochen Burghardt

Application du lemme de pumping à XML (2)

$$s = \langle \underbrace{a \cdots a}_p \underbrace{b \cdots b}_p \rangle \langle \underbrace{a \cdots a}_p \underbrace{b \cdots b}_p \rangle$$

- $s = uvwxy$, $|vwx| \leq p$, $|vx| \geq 1$, $uwy \in L$.
- Cas 2: vwx

Application du lemme de pumping à XML (1)

Considérons le mot suivant qui est en L_{XML} :

$$s = \langle \underbrace{a \cdots a}_p \underbrace{b \cdots b}_p \rangle \langle \underbrace{a \cdots a}_p \underbrace{b \cdots b}_p \rangle$$

- On peut écrire $s = uvwxy$, avec $|vwx| \leq p$ et $|vx| \geq 1$, tel que $uwy \in L_{XML}$.
- Cas 1: vwx est entièrement dans la partie verte ou la partie rouge: contradiction car, après passage à uwy , les parties verte et rouge sont de longueur différente!

Pour aller plus loin

- ▶ Il existe un langage de *schéma* pour des documents XML. Un schéma restreint la forme d'un document XML. DTD (Document Type Definition).
- ▶ Un schéma peut par exemple définir quels attributs sont obligatoires (autorisés) pour quel type, ou quels types d'éléments sont obligatoires (autorisés) pour un nœud d'un certain type.
- ▶ Il existe aussi des langages riches pour des requêtes (extraction de données d'un document XML), par ex. XPath, XQuery.
- ▶ Pour en savoir plus : voir le cours XML du M1



Analyse de Données Structurées- Cours 9

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

2 avril 2015

© Ralf Treinen 2015

Théorie et Pratique

Tâche	Modèle théorique	Réalisation
Analyse lexicale	Expressions régulières	Fichier JFlex
Analyse syntaxique	Grammaires	Implémentation d'un parseur LL(1)
Syntaxe abstraite	Définition inductive	Classes Java
Sémantique	Règles sémantiques	Parcours d'arbre

Analyse Sémantique

- ▶ Nous savons maintenant construire un arbre de syntaxe abstraite à partir d'une représentation textuelle (syntaxe concrète).
- ▶ Par exemple : syntaxe abstraite pour un langage de programmation, pour un langage de données structurées comme XML ou JSON.
- ▶ Nous savons aussi vu la semaine dernière que certaines restrictions syntaxiques ne peuvent pas être réalisées par une grammaire.
- ▶ C'est maintenant le rôle de l'analyse sémantique.
- ▶ Les techniques sont les mêmes que pour l'interprétation.

Définition Inductive

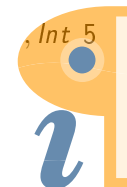
L'ensemble E des arbres de syntaxe abstraite représentant des expressions régulières est défini comme suit :

- ▶ tout $Int\ n$, où $n \in \mathbb{N}$, est un élément de E
- ▶ si $e_1, e_2 \in E$, alors $Sum\ e_1, e_2 \in E$
- ▶ si $e_1, e_2 \in E$, alors $Product\ e_1, e_2 \in E$

Ces sont les seules façons de construire un élément de E .

Exemples

- ▶ $Int\ 42$
- ▶ $Product\ Sum\ Int\ 2, Int\ 49$



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Remarques

- ▶ Attention : ils'agit d'une notation mathématique pour des *arbres*. Ce n'est pas une définition d'une syntaxe concrète.
- ▶ On utilise parfois aussi pour la définition des arbres de syntaxe abstraite un formalisme appelé *grammaire d'arbres*, ce qui risque créer une confusion. Pour nous, une grammaire défini toujours une syntaxe concrète.
- ▶ L'intérêt de cette notation est qu'elle nous permet de rédiger les règles de calcul sur la syntaxe abstraite. Cette notation est plus compacte que de dessiner un arbre.
- ▶ Souvent on définit plusieurs types d'arbres de syntaxe abstraite, par exemple un type pour les expressions arithmétiques, un autre pour les instructions d'un langage de programmation.

Définition de la Syntaxe Abstraite sur l'exemple I

```
abstract class Expression{
    public boolean isInt(){
        return false;
    }
    public boolean isSum(){
        return false;
    }
    public boolean isProduct(){
        return false;
    }
}

class Int extends Expression{
    private int value;
```

Traduction en JAVA

- ▶ Chaque type d'arbre correspond à une classe abstraite.
- ▶ Pour chacun des cas dans la définition il y a une classe dérivée.
- ▶ Le constructeur de cette classe dérivée a comme arguments les valeurs différentes utilisées dans la construction de l'arbre.
- ▶ Le constructeur stocke ces valeurs dans des champs privés de la classe.
- ▶ Éventuellement : des méthodes pour savoir dans quel cas on est, et pour récupérer les valeurs utilisées dans la construction.

Définition de la Syntaxe Abstraite sur l'exemple II

```
public boolean isInt(){
    return true;
}

public Int(int i){
    value = i;
}

public int getValue(){
    return value;
}

class Sum extends Expression{
    private Expression left;
    private Expression right;
    public boolean isSum(){
```



DéfinitiondelaSyntaxeAbstraitesurl'exempleIII

```
        returntrue    ;
    }
    public    Sum(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }
    public    ExpressiongetLeft(){
        return    left;
    }
    public    ExpressiongetRight(){
        return    right;
    }
}

class    Product    extends    Expression{
```

DéfinitiondelaSyntaxeAbstraitesurl'exempleV

```
}
```

DéfinitiondelaSyntaxeAbstraitesurl'exempleIV

```
private    Expressionleft;
private    Expressionright;
publicboolean    isProduct(){
    returntrue    ;
}
public    Product(Expressione1,Expressione2){
    left=e1;
    right=e2;
}
public    ExpressiongetLeft(){
    return    left;
}
public    ExpressiongetRight(){
    return    right;
}
```

LamêmechosedanslelangageOCaml

```
type    expression=
| Integer    of    int
| Sum    of    expression    expression
| Product    of    expression    expression;;
```

(voirlecours *ProgrammationFonctionnelle* duL3)



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Objectif

- ▶ Jusqu'à maintenant: syntaxe concrète et syntaxe abstraite
- ▶ Reste à faire:
 - ┌ Les restrictions supplémentaires qu'on ne peut pas exprimer par une grammaire: *analyse sémantique*
 - ┌ Exécution ou compilation d'un programme: *sémantique*
- ▶ Il nous faut d'abord une notation qui nous permet de spécifier précisément les restrictions supplémentaires, et la sémantique.
- ▶ Puis: codage (en Java, par exemple)

L'enchaînement des trois étapes

1. Normalement: Construction de la syntaxe abstraite pour le programme entier avant toute analyse sémantique ou exécution (attention: il y a des cas où analyse syntaxique et exécution sont intercalées)
2. La vérification des conditions sémantiques est faite *avant* l'exécution. Pour cette raison on parle aussi parfois d'une sémantique *statique*.
3. L'exécution (ou la compilation) peut utiliser des résultats de l'étape (2). Dans ce cas, on l'étape (3) se base sur une syntaxe abstraite *annotée* (exemple: inférence de types).

Langages de programmation: Trois étapes

1. De la syntaxe concrète à la syntaxe abstraite.
2. Analyse sémantique: par exemple assurer que toutes les variables sont déclarées, toutes les méthodes définies, que le type est correct.
3. Exécution de la syntaxe abstraite (interprétation), ou traduction vers un autre langage (compilation, ou transformation de code)

Notion de jugement

- ▶ Un jugement est une affirmation portant sur une expression abstraite.
- ▶ Quelques exemples, en langage naturel:
 - ┌ «l'expression *Sum Int 1, Product Int 2, Int 3* a pour valeur 7»;
 - ┌ «l'expression *Sum Int 1, Product Int 2, Int 3* a pour type int»;
 - ┌ «l'instruction *Write Sum Int 3, Int 4* affiche 7»;
 - ┌ «l'expression *Sum Int 3, Int 5* a pour code machine [PUSHI3; PUSHI5; ADD]».



Jugements avec contexte

- ▶ « dans le contexte où x est une variable valant 4, l'expression $\text{Sum Int } 1, \text{Product Var } x, \text{Int } 3$ a pour valeur 13 »;
- ▶ « dans le contexte où x est une variable de type réel, l'expression $\text{Sum Int } 1, \text{Product Var } x, \text{Int } 3$ a pour type réel »;
- ▶ pour un compilateur devant engendrer du code machine à une adresse donnée, on pourra avoir des jugements du style « dans le contexte où l'adresse courante du code est 250 et x est une variable dont l'adresse est 100, l'expression $\text{Product Int } 3, \text{Var } x$ a pour code machine [250 PUSH 13; 251 PUSH 100; 252 MUL] ».

Notation et exemples

- ▶ On note le jugement $C, a, v \in D^c \times A \times D^v$ par « $C \vdash a \Rightarrow v$ ».
- ▶ Dans le cas d'un domaine de contexte trivial: $a \Rightarrow v$.
- ▶ Exemples:
 - $\text{Sum Int } 1, \text{Product Int } 2, \text{Int } 3 \Rightarrow 7$
est le jugement « l'expression $1 + 2 \times 3$ a pour valeur 7 ».
 - $\{x \Rightarrow 4\} \vdash \text{Sum Int } 1, \text{Product Var } x, \text{Int } 3 \Rightarrow 13$
est le jugement « dans le contexte où x est une variable valant 4, l'expression $1 + x \times 3$ a pour valeur 13 », ici $D^v = \mathbb{N}$ et D^c est l'ensemble des *environnements d'évaluation*.

Définition

Un *jugement* est un triplet $C, a, v \in D^c \times A \times D^v$ tel que:

- ▶ A est un ensemble d'arbres de syntaxe abstraite, contenant l'arbre a ;
- ▶ D^v est un ensemble appelé *domaine de valeurs*; la valeur v appartient au sous-ensemble D^v des valeurs;
- ▶ D^c est un ensemble appelé *domaine de contextes*; la valeur C appartient au sous-ensemble D^c des valeurs.

Remarques

- ▶ Quand le domaine de contextes est trivial, le jugement se réduit à une paire a, v .
- ▶ Les contextes et les valeurs peuvent être des paires, triplets, ...

Variantes de notation

Lorsque le domaine de valeurs D^v est un ensemble de types, par exemple $D^v = \{\text{int}, \text{bool}, \text{real}\}$, nous utiliserons une variante classique de cette notation, $C \vdash a : t$, où $t \in D^v$. On écrira donc:

- ▶ $\text{Sum Int } 1, \text{Product Int } 2, \text{Int } 3 : \text{int}$
pour le jugement « l'expression $1 + 2 \times 3$ a pour type int »;
- ▶ $\{x : \text{real}\} \vdash \text{Sum Int } 1, \text{Product Var } x, \text{Int } 3 : \text{real}$
pour le jugement « dans le contexte où x est une variable de type réel, l'expression $1 + x \times 3$ a pour type réel ».

Jugements vrais et faux

La définition précédente autorise l'écriture de jugements vrais aussi bien que de jugements faux:

- ▶ « l'expression $\text{Sum Int } 1, \text{Mult Int } 2, \text{Int } 3$ a pour valeur 12 »;
- ▶ « l'expression $\text{Sum Int } 1, \text{Mult Int } 2, \text{Int } 3$ a pour type bool »;

Les règles générales

- ▶ Les règles proprement dites correspondent au cas général de la récurrence.
- ▶ Elles affirment la vérité d'un jugement qui constitue la **conclusion** de la règle sous l'hypothèse (de récurrence) de la vérité d'un ou plusieurs jugements appartenant aux **prémisses** de la règle.
- ▶ Par exemple:
Si le jugement $e_1 \text{ int}$ est vrai et si le jugement $e_2 \text{ int}$ est vrai alors le jugement $\text{Sum } e_1, e_2 \text{ int}$ est également vrai.

Axiomes

- ▶ Les règles sémantiques permettent d'établir qu'un certain jugement est vrai.
- ▶ Cas de base: Les axiomes: ce sont des jugements pour lesquels on affirme qu'ils sont vrais.
- ▶ Exemples (voir plus tard pour la notation):
 - ┆ si n est une constante entière, alors le jugement $\text{Int } n \Rightarrow n$ est vrai;
 - ┆ si n est une constante entière, alors le jugement $\text{Int } n \text{ int}$ est vrai;
 - ┆ dans un contexte C où x est une variable de type bool, le jugement $C \vdash \text{Var } x \text{ bool}$ est vrai.

Règles sémantiques

Définition

Une **règle sémantique** est un k -uplet de jugements, avec $k \geq 1$. Si $k = 1$, la règle est appelée un **axiome**. Dans le cas d'une règle J_1, \dots, J_k avec $k \geq 2$, J_1, \dots, J_{k-1} constituent les **prémisses** de la règle et J_k en est la **conclusion**.

Notation habituelle

Pour la règle J_1, \dots, J_k

$$\frac{J_1, \dots, J_{k-1}}{J_k}$$

Dans le cas d'un axiome, rien ne figure au-dessus du trait de fraction.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Les schémas de règles

- Dans le cas de l'évaluation d'une expression arithmétique il y a déjà une infinité d'axiomes:

$$\frac{}{Int\ 0 \Rightarrow 0} \quad \frac{}{Int\ 1 \Rightarrow 1} \quad \text{etc.}$$

- Pour cette raison, on utilise des *schéma d'axiomes*

$$\text{pour tout } n \in \mathbb{N} : \frac{}{Int\ n \Rightarrow n}$$

- De la même façon, on a les schémas de règles suivantes pour l'addition et la multiplication d'entiers:

$$\frac{e_1 \Rightarrow v_1 \quad e_2 \Rightarrow v_2}{Sum\ e_1, e_2 \Rightarrow v_1 +_{int} v_2}$$

et

$$\frac{e_1 \Rightarrow v_1 \quad e_2 \Rightarrow v_2}{Product\ e_1, e_2 \Rightarrow v_1 \times_{int} v_2}$$

Des règles sémantiques à des méthodes récursives

- Dans des cas simples: les jugements spécifient une fonction, par exemple la fonction d'évaluation qui associe à chaque arbre de syntaxe abstraite une valeur *unique*.
- Cette fonction peut être partielle.
- Dans notre cas de l'évaluation: trois schémas de règles, chacune correspondant à un cas dans la définition des arbres de syntaxe abstraite.
- Il y a une méthode abstraite `ev 1`.
- Chaque sous-classe définit cette méthode selon le schéma donnée:
 - └ L'argument de la méthode est le contexte.
 - └ Les prémisses correspondent à des appels récursifs.
 - └ La conclusion dit alors comment calculer le résultat.

Les schémas de règles

- Dans un schéma de règles comme

$$\frac{e_1 \Rightarrow v_1 \quad e_2 \Rightarrow v_2}{Sum\ e_1, e_2 \Rightarrow v_1 +_{int} v_2}$$

e_1 et e_2 sont des *méta-variables* qu'il faut remplacer par des arbres de syntaxe abstraite pour obtenir une instance.

- Souvent on confond les notions de règle et de schéma de règle.
- Les opérateurs `int` et `×int` sont les opérations arithmétiques.

Évaluation: code

```
abstract class Expression {
    abstract int eval();
}

class Int extends Expression {
    private int value;
    public Int(int i) {
        value = i;
    }
    public int eval() {
        return value;
    }
}
```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Évaluation: code I

```

class Sum extends Expression{
    private Expression left;
    private Expression right;
    public Sum(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }
    public int eval(){
        return left.eval()+right.eval();
    }
}

```

```

class Product extends Expression{
    private Expression left;
    private Expression right;

```

Évaluation: code II

```

    public Product(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }
    public int eval(){
        return left.eval() + right.eval();
    }
}

```



Analyse de Données Structurées- Cours 10

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

9 avril 2015

© Ralf Treinen 2015

Expressions avec variables

- ▶ L'ensemble E' de arbres de syntaxe abstraite représentant des expressions régulières est défini comme suit:
 - └ tout $Int(n)$, où $n \in \mathbb{N}$, est un élément de E'
 - └ tout $Ident(s)$, où $s \in \Sigma^*$, est un élément de E'
 - └ si $e_1, e_2 \in E'$, alors $Sum(e_1, e_2) \in E'$
 - └ si $e_1, e_2 \in E'$, alors $Product(e_1, e_2) \in E'$
- ▶ Σ est l'alphabet utilisé pour former des noms de variables (par exemple, tous les caractères UNICODE).
- ▶ Éventuellement, l'ensemble des noms de variables possibles dans la syntaxe abstraite est plus riche que les noms de variables permis par la syntaxe concrète.

Rappel: évaluation d'expressions sans variables

- ▶ L'ensemble E de arbres de syntaxe abstraite représentant des expressions régulières est défini comme suit:
 - └ tout $Int(n)$, où $n \in \mathbb{N}$, est un élément de E
 - └ si $e_1, e_2 \in E$, alors $Sum(e_1, e_2) \in E$
 - └ si $e_1, e_2 \in E$, alors $Product(e_1, e_2) \in E$
- ▶ Règles sémantiques:

$$\frac{}{Int(n) \Rightarrow n} \quad n \in \mathbb{N}$$

$$\frac{e_1 \Rightarrow v_1 \quad e_2 \Rightarrow v_2}{Sum(e_1, e_2) \Rightarrow v_1 +_{int} v_2}$$

$$\frac{e_1 \Rightarrow v_1 \quad e_2 \Rightarrow v_2}{Product(e_1, e_2) \Rightarrow v_1 \times_{int} v_2}$$

Évaluation d'expressions avec variables

- ▶ Maintenant, l'évaluation d'une expression dépend aussi d'un *contexte*.
- ▶ La définition du contexte dépend de l'application. Pour l'évaluation: ils'agit d'une fonction partielle

$$\Sigma^* \rightarrow \mathbb{N}$$

aussi appelée *environnement de valeurs*.

- ▶ Le jugement d'évaluation:

$$\Gamma \vdash e \Rightarrow v$$

exprime: Dans l'environnement de valeur Γ , la valeur de l'expression e est v .

- ▶ Exemple d'un jugement vrai:

$$\{x = 3, y = 5\} \vdash Sum(Product(Ident(x), Int(2)), Ident(y)) \Rightarrow 11$$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm



Règles pour l'évaluation des expressions avec variables

► Axiomes:

$$\frac{}{\Gamma \vdash \text{Int}(n) \Rightarrow n} \quad n \in \mathbb{N}$$
$$\frac{}{\Gamma \vdash \text{Ident}(s) \Rightarrow n} \quad \Gamma(s) = n$$

► Règles récursives:

$$\frac{\Gamma \vdash e_1 \Rightarrow v_1 \quad \Gamma \vdash e_2 \Rightarrow v_2}{\Gamma \vdash \text{Sum}(e_1, e_2) \Rightarrow v_1 +_{\text{int}} v_2}$$
$$\frac{\Gamma \vdash e_1 \Rightarrow v_1 \quad \Gamma \vdash e_2 \Rightarrow v_2}{\Gamma \vdash \text{Product}(e_1, e_2) \Rightarrow v_1 \times_{\text{int}} v_2}$$

Remarques

- Quand $s \notin \text{domain}(\Gamma)$, alors le schéma d'axiome ne s'applique pas à $\text{Ident}(s)$. Dans ce cas on ne peut donc pas calculer la valeur d'une expression qui contient $\text{Ident}(s)$.
- Dans les règles récursives, comme

$$\frac{\Gamma \vdash e_1 \Rightarrow v_1 \quad \Gamma \vdash e_2 \Rightarrow v_2}{\Gamma \vdash \text{Sum}(e_1, e_2) \Rightarrow v_1 +_{\text{int}} v_2}$$

le contexte Γ utilisé dans les hypothèses est le même que dans la conclusion. Le calcul correspondant aux deux hypothèses peut donc se faire en parallèle.

Implémentation: ValueEnvironment.java

- Implémentation des environnements de valeur: table de hachage.

```
import java.util.HashMap;

class ValueEnvironment extends HashMap<String,Integer>{
    public final static long serialVersionUID=44;
}
```

Implémentation: Expression.java

```
abstract class Expression{
    abstract int eval(ValueEnvironment env);
}

class Int extends Expression{
    private int value;
    public Int(int i){
        value=i;
    }
    public int eval(ValueEnvironment env){
        return value;
    }
}
```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Implémentation: Expression.java II

```

class Ident extends Expression{
    private String name;
    public Ident(String s){
        name=s;
    }
    public int eval(ValueEnvironment env){
        return env.get(name);
    }
}

class Sum extends Expression{
    private Expression left;
    private Expression right;
    public Sum(Expression e1, Expression e2){
        left=e1;

```

Implémentation: Expression.java IV

```

        return left.eval(env) + right.eval(env);
    }
}

```

Implémentation: Expression.java III

```

        right=e2;
    }
    public int eval(ValueEnvironment env){
        return left.eval(env) + right.eval(env);
    }
}

class Product extends Expression{
    private Expression left;
    private Expression right;
    public Product(Expression e1, Expression e2){
        left=e1;
        right=e2;
    }
    public int eval(ValueEnvironment env){

```

Affectation et Affichage

- ▶ Objectif: exécution de programmes.
- ▶ Première étape: séquences d'affectation à des variables, et d'instructions d'affichage.
- ▶ Instruction "print" pour afficher, symbole ":" pour l'affectation.
- ▶ Instruction terminée par le symbole ";"
- ▶ Fragment de grammaire (le reste est comme avant):

$$S \rightarrow IL \text{ EOF}$$

$$IL \rightarrow \epsilon \mid I ; IL$$

$$I \rightarrow \text{ident} := E \mid \text{print } E$$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Règles sémantiques

- ▶ On s'intéresse surtout à la modification des environnements.
- ▶ Jugement $\Gamma_1 \vdash i \Rightarrow \Gamma_2$: l'exécution de l'instruction (ou de la liste d'instructions) i dans l'environnement Γ_1 résulte dans l'environnement Γ_2 .
- ▶ Modélisation des affichages : omise.

Remarques

- ▶ Maintenant, le contexte dans les hypothèses n'est plus toujours le même que dans la conclusion.
- ▶ Les hypothèses des règles peuvent être des jugements qui portent sur des expressions de syntaxe abstraite de sorte différente que la conclusion (exemple : règle pour $Affect$).
- ▶ $\Gamma[s \mapsto v]$ est l'environnement qui associe à s la valeur v , et qui associe $\Gamma(x)$ à tout $x \neq s$.
- ▶ La deuxième règle exprime que l'exécution de $Print$ ne modifie pas l'environnement.

Instructions

- ▶ Syntaxe abstraite : L'ensemble I des arbres de syntaxe abstraite de sorte $Instruction$ est défini comme suit :
 - ┆ tout $Assign(s, e)$, où $s \in \Sigma^*$ et $e \in E'$, est un élément de I ;
 - ┆ tout $Print(e)$, où $e \in E'$, est un élément de I .
- ▶ Règles sémantiques :

$$\frac{\Gamma_1 \vdash e \Rightarrow v_1}{\Gamma_1 \vdash Affect(s, e) \Rightarrow \Gamma_1[s \mapsto v]}$$

$$\frac{}{\Gamma_1 \vdash Print(s, e) \Rightarrow \Gamma_1}$$

Listes d'instructions

- ▶ Syntaxe abstraite : L'ensemble IL des arbres de syntaxe abstraite de sorte $InstructionList$ est défini comme suit :
 - ┆ Nil est un élément de IL ;
 - ┆ Si $il \in IL$ et $i \in I$, alors $Seq(i, il)$ est un élément de IL .
- ▶ Règles sémantiques :

$$\frac{\Gamma_1 \vdash Nil \Rightarrow \Gamma_1 \quad \Gamma_1 \vdash i \Rightarrow \Gamma_2 \quad \Gamma_2 \vdash il \Rightarrow \Gamma_3}{\Gamma_1 \vdash Seq(i, il) \Rightarrow \Gamma_3}$$



Remarques

- ▶ La première règle (qui est un axiome) exprime que l'exécution de *Nil* ne modifie pas l'environnement.
- ▶ La deuxième règle impose maintenant un ordre sur l'exécution des deux composantes *i* et *il* :
 - ┆ Donnée initialement : Γ_1 et $Seq(i, il)$.
 - ┆ Il faut d'abord exécuter la première hypothèse qui donne, à partir de Γ_1 et de *i*, l'environnement Γ_2 .
 - ┆ Puis on peut exécuter la deuxième hypothèse qui donne, à partir de Γ_2 et de *il*, l'environnement Γ_3 .
 - ┆ L'implémentation utilise des tables de hachage modifiables, il faut bien respecter l'ordre !

Implémentation: Instruction.java II

```
}  
  
class Print extends Instruction{  
    private Expression expression;  
    public Print(Expression e){  
        expression=e;  
    }  
    public void exec(ValueEnvironment env){  
        System.out.println(expression.eval(env));  
    }  
}
```

Implémentation: Instruction.java I

```
abstract class Instruction{  
    abstract void exec(ValueEnvironment env);  
}  
  
class Assignment extends Instruction{  
    private Expression expression;  
    private String variable;  
    public Assignment(String v, Expression e){  
        expression=e;  
        variable=v;  
    }  
    public void exec(ValueEnvironment env){  
        env.put(variable, expression.eval(env));  
    }  
}
```

Implémentation: InstructionList.java I

```
abstract class InstructionList{  
    abstract void exec(ValueEnvironment env);  
}  
  
class Nil extends InstructionList{  
    public Nil(){  
    }  
    public void exec(ValueEnvironment env){  
    }  
}  
  
class Seq extends InstructionList{  
    private Instruction head;  
    private InstructionList rest;
```



Implémentation: InstructionList.java II

```
public Seq(Instruction i, InstructionList il) {
    head=i;
    rest=il;
}
public void exec(ValueEnvironment env){
    head.exec(env);
    rest.exec(env);
}
}
```

Typage statique et typage dynamique

- ▶ Typage statique: vérification des types après construction de la syntaxe abstraite, et *avant* l'exécution.
- ▶ Avantage: Plus de sécurité, une fois le typage vérifié par le compilateur on sait que des erreurs de *type* ne peuvent plus se produire (sauf certains langages de programmation qui permettent de contourner le typage).
- ▶ Avantage: le typage peut annoter la syntaxe abstraite par des types des expressions, ce qui simplifie l'exécution. Désambiguïson des opérateurs surchargés.
- ▶ Exemples: Java, Ocaml, ...

Expressionset Types

- ▶ Pour l'exécution des programmes avec, par exemple des instructions conditionnelles et des boucles, nous avons besoin des expressions booléennes.
- ▶ On ne peut pas faire la distinction entre expressions entières et expressions booléennes par la grammaire à cause des variables.

Typage statique et typage dynamique

- ▶ Typage dynamique: vérification des types pendant *exécution* du programme, quand on applique un opérateur à des valeurs.
- ▶ Avantage: plus de flexibilité (mais: flexibilité aussi possible dans le cas statique, par ex. polymorphie).
- ▶ Avantage: programmes moins verbeux, car il n'est pas nécessaire de déclarer les variables avec leur type (mais: le même avantage peut être obtenu par une *inférence de type*).
- ▶ Exemples: Python, bash, Perl.



Unpetittourdequelqueslangagesdeprogrammation

- ▶ bash : analysesyntaxiqueàlavolée. Possibilitédefaire interpréterunechaînedecaractèrecommeunecommandede lshell (commande `ev 1`).
- ▶ Python : analysesyntaxiquecomplèteavantexécution, typage dynamique.
- ▶ Java : typagestatique, déclarationsexplicites
- ▶ Ocaml : typagestatique, inférencedetype.

Bash:eval

```
#!/bin/bash  
#evaluationd'unechaîneconstruitependant  
#l'exécutionduprogramme.
```

```
mot1="ecrire"  
mot2="hello"  
mot3="bonjour"
```

```
mot4=' sed  -e's/crire/cho/'<<<$mot1'
```

```
eval  "$mot4_$mot2_$mot3"
```

Bash: analysesyntaxiqueàlavolée

```
#!/bin/bash  
#lesdeuxpremieresslignessontexecuteesavant  
#quel'erreurdesyntaxeestdetectee.
```

```
echo  "hello"  
echo  "bonjour"  
if fi
```

Python: analysesyntaxiquecomplète

```
#!/usr/bin/python  
#erreurdesyntaxedetecteeavantexécution
```

```
print  "hello"  
print  "bonjour"  
if  fi
```



Python: typage dynamique

```
#!/usr/bin/python
#une variable peut prendre des valeurs de
#type différent.
```

```
x=42
print x
x=True
print x
```

OCaml: inférence de types

```
(* le type de la fonction fest infere *)
let somme x y = x + y;;
```

```
(* erreur de typage *)
somme "hello" "bonjour";;
```

Python: typage dynamique

```
#!/usr/bin/python
```

```
#erreur de typage seulement detectee quand
#la branche else est executee
```

```
import sys
read=sys.stdin.readline()
if read!='coocoo\n':
    print 17+21
else :
    print ('abc' & 'def')
```

Typage statique et typage dynamique

- ▶ Dans le cas d'un typage statique, des types sont initialement par des déclarations associés à des *identificateurs*. Puis, une inférence de type peut déduire le type pour chaque expression, et vérifier que tous les opérateurs sont utilisés en cohérence avec le type des expressions.
- ▶ Dans le cas d'un typage dynamique, les types sont associés aux *valeurs* (éventuellement résultat d'un calcul). C'est au moment de l'application de l'opérateur qu'une incompatibilité de types peut être détectée.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Analyse de Données Structurées- Cours 11

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

15 avril 2015

© Ralf Treinen 2015

Règles de typage

- Environnement de typage: fonction partielle $\Sigma^* \rightarrow \{int, bool\}$
- Jugement $\Gamma \vdash e : t$: Dans l'environnement Γ , l'expression e a le type t .

$$\frac{}{\Gamma \vdash Int(n) : int} \quad n \in \mathbb{N}$$

$$\frac{}{\Gamma \vdash True : bool}$$

$$\frac{}{\Gamma \vdash False : bool}$$

$$\frac{}{\Gamma \vdash Ident(s) : \Gamma(s)} \quad s \in \Sigma^*$$

Syntaxe Abstraite

L'ensemble E'' des arbres de syntaxe abstraite représentant des expressions arithmétiques et booléennes est défini comme suit:

- tout $Int(n)$, où $n \in \mathbb{N}$, est un élément de E''
- $True$ et $False$ sont des éléments de E''
- tout $Ident(s)$, où $s \in \Sigma^*$, est un élément de E''
- si $e_1, e_2 \in E''$, alors $Plus(e_1, e_2) \in E''$
- si $e_1, e_2 \in E''$, alors $Mult(e_1, e_2) \in E''$
- si $e_1, e_2 \in E''$, alors $Greater(e_1, e_2) \in E''$
- si $e_1, e_2 \in E''$, alors $Equal(e_1, e_2) \in E''$
- si $e_1, e_2 \in E''$, alors $And(e_1, e_2) \in E''$
- si $e_1, e_2 \in E''$, alors $Or(e_1, e_2) \in E''$
- si $e \in E''$, alors $Not(e) \in E''$

Règles de typage

$$\frac{\Gamma \vdash e_1 : int \quad \Gamma \vdash e_2 : int}{\Gamma \vdash Plus(e_1, e_2) : int}$$

$$\frac{\Gamma \vdash e_1 : int \quad \Gamma \vdash e_2 : int}{\Gamma \vdash Mult(e_1, e_2) : int}$$

$$\frac{\Gamma \vdash e_1 : int \quad \Gamma \vdash e_2 : int}{\Gamma \vdash Equal(e_1, e_2) : bool}$$

$$\frac{\Gamma \vdash e_1 : int \quad \Gamma \vdash e_2 : int}{\Gamma \vdash Greater(e_1, e_2) : bool}$$

Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Règlesdetypage

$$\frac{\Gamma \vdash e_1 : bool \quad \Gamma \vdash e_2 : bool}{\Gamma \vdash And(e_1, e_2) : bool}$$

$$\frac{\Gamma \vdash e_1 : bool \quad \Gamma \vdash e_2 : bool}{\Gamma \vdash Or(e_1, e_2) : bool}$$

$$\frac{\Gamma \vdash e : bool}{\Gamma \vdash Not(e) : bool}$$

Implémentation: Expression.java (début)II

```
        this .type=Type.Int;
    }
}

class True extends Expression{
    public True(){
    }
    public void setType(TypeEnvironment env){
        this .type=Type.Bool;
    }
}

class False extends Expression{
    public False(){
    }
}
```

Implémentation: Expression.java (début)I

```
abstract class Expression{
    Type type;
    final public Type getType(){
        return this .type;
    }
    abstract void setType(TypeEnvironment env) throws TypeException;
}

class Int extends Expression{
    private int value;
    public Int( int i){
        value=i;
    }
    public void setType(TypeEnvironment env){
    }
}
```

Implémentation: Expression.java (début)III

```
    public void setType(TypeEnvironment env){
        this .type=Type.Bool;
    }
}

class Ident extends Expression{
    private String name;
    public Ident(Strings){
        name=s;
    }
    public void setType(TypeEnvironment env){
        this .type=env.get(name);
    }
}
```



Implémentation: Expression.java (début)IV

```
class Sum extends Expression{
    private Expression left;
    private Expression right;
    public Sum(Expressione1,Expressione2){
        left=e1;
        right=e2;
    }
    public void setType(TypeEnvironment env) throws TypeException{
        left.setType(env);
        right.setType(env);
        if (left.getType()==Type.Int
            && right.getType()==Type.Int){
            this.type=Type.Int;
        } else {
            thrownew TypeException("Sum");
        }
    }
}
```

Unreprésentationplusgénérique?

- ▶ Celadevientassezfastidieuxetrépétitif.
- ▶ Imaginezqu'est-cequeçadonnerasionajouteencoreplus d'opérateurs,ou desnouvellessortes(flottants,chaînesde caractère, ...)
- ▶ Solutionpluscompacte,etplusfacileàgénéraliser:
 - └ Définirunesorted'opérateurs(uneparnombred'arguments)
 - └ Définirunopérateur d'application d'unopérateur aubon nombred'expressions.

Implémentation: Expression.java (début)V

```
    }
    }
}
```

Syntaxeabstraitegénérique:opérateuretconstantes

- ▶ L'ensemble O_2 desopérateursensyntaxeabstraiteest:

$$O_2 = \{Plus, Mult, Greater, Equal, And, Or\}$$

- ▶ L'ensemble O_1 desopérateursunairesensyntaxeabstraiteest

$$O_1 = \{Not\}$$

- ▶ L'ensemble O_0 desconstantesensyntaxeabstraiteest

$$O_0 = \{True, False\} \cup \{Int(n) \mid n \in \mathbb{N}\}$$



Syntaxe abstraite: expressions

L'ensemble E''' des expressions en syntaxe abstraite est:

- ▶ $Ident(s) \in E'''$ pour tout $s \in \Sigma^*$
- ▶ si $o \in O_0$ alors $Const(o) \in E'''$ pour tout $n \in \mathbb{N}$
- ▶ si $o \in O_1$, $e \in E'''$ alors $Apply_1(o, e) \in E'''$
- ▶ si $o \in O_2$, $e_1, e_2 \in E'''$ alors $Apply_2(o, e_1, e_2) \in E'''$

Typage des constantes et opérateurs

- ▶ Constantes:

$$\begin{aligned} ot(True) = ot(False) &= bool \\ ot(Int(n)) &= int \end{aligned}$$

- ▶ Opérateurs unaires:

$$ot(Not) = (bool, bool)$$

- ▶ Opérateurs binaires:

$$\begin{aligned} ot(Plus) = ot(Mult) &= (int, int, int) \\ ot(And) = ot(Or) &= (bool, bool, bool) \\ ot(Greater) = ot(Equal) &= (int, int, bool) \end{aligned}$$

Typage pour la forme générique

- ▶ L'ensemble des types de base est $BT = \{int, bool\}$
- ▶ On va définir une fonction ot qui associe:
 - └ à toute constante $o \in O_0$ un élément de BT ,
 - └ à tout opérateur unaire $o \in O_1$ une paire dans $BT \times BT$
 - └ à tout opérateur binaire $o \in O_2$ un triplet dans $BT \times BT \times BT$.
- ▶ Dans tous les cas, le dernier type est celui du résultat de l'application de l'opérateur, et les types précédents ceux des arguments.
- ▶ Parfois notation: $t_1 \times t_2 \rightarrow t_3$ au lieu de (t_1, t_2, t_3) .

Typage des expressions

$$\frac{}{\Gamma \vdash Ident(s) : \Gamma(s)} \quad s \in \Sigma^*$$

$$\frac{}{\Gamma \vdash Const(o) : t} \quad o \in O_0, ot(o) = t$$

$$\frac{\Gamma \vdash e_1 : t_1}{\Gamma \vdash Apply_1(o, e) : t} \quad o \in O_1, ot(o) = (t_1, t)$$

$$\frac{\Gamma \vdash e_1 : t_1 \quad \Gamma \vdash e_2 : t_2}{\Gamma \vdash Apply_2(o, e_1, e_2) : t} \quad o \in O_2, ot(o) = (t_1, t_2, t)$$



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Syntaxe Abstraite

- Définitions mutuellement récursives des Instructions et des Listes d'Instructions.
- L'ensemble I des Instructions syntaxe abstraite est:
 - └ si $e \in E'''$ alors $Print(e) \in I$,
 - └ si $s \in \Sigma^*$, $e \in E'''$ alors $Assign(s, e) \in I$,
 - └ si $e \in E'''$, $il_1, il_2 \in IL$ alors $Cond(e, il_1, il_2) \in I$,
 - └ si $e \in E'''$, $il \in IL$ alors $While(e, il) \in I$.
- L'ensemble IL des Listes d'Instructions syntaxe abstraite est:
 - └ $Nil \in IL$,
 - └ si $i \in I$ et $il \in IL$ alors $Seq(i, il) \in IL$.

Déclaration et typage

- Jugement $dl : \Gamma$: la liste de déclaration dl donne l'environnement de typage Γ .
- Règles:

$$\overline{DeclNil : \emptyset}$$

$$\frac{dl : \Gamma}{DeclSeq(Decl(s, t), dl) : \Gamma[s \mapsto t]} \quad \text{si } s \notin \text{domain}(\Gamma)$$

Syntaxe Abstraite

- L'ensemble D des Déclarations syntaxe abstraite est:
 - └ si $s \in \Sigma^*$ et $t \in BT$ alors $Decl(s, t) \in D$.
- L'ensemble DL des Listes de Déclarations syntaxe abstraite est:
 - └ $DeclNil \in DL$,
 - └ si $d \in D$ et $dl \in DL$ alors $DeclSeq(d, dl) \in DL$.
- L'ensemble P des Programmes syntaxe abstraite est:
 - └ si $dl \in DL$ et $il \in IL$ alors $Program(dl, il) \in P$.

Implémentation: DeclarationList.java |

```

abstract class DeclarationList {
    abstract TypeEnvironment extract() throws DeclException;
}

class DeclNil extends DeclarationList {
    public DeclNil() {}
    public TypeEnvironment extract() {
        return new TypeEnvironment();
    }
}

class DeclSeq extends DeclarationList {
    private Declaration head;

```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Implémentation: DeclarationList.java II

```

private DeclarationList rest;
public DeclSeq(Declaration d, DeclarationList dl){
    head=d;
    rest=dl;
}
public TypeEnvironment extract() throws DeclException{
    TypeEnvironment env=rest.extract();
    String head_var=head.getVariable();
    if (env.containsKey(head_var)){
        throw new DeclException(head_var);
    } else {
        env.put(head_var, head.getType());
        return env;
    }
}

```

Typage des listes d'instructions

- ▶ Jugement $\Gamma \vdash il : il$: il est correctement typée dans Γ .
- ▶ Règles:

$$\frac{}{\Gamma \vdash Nil}$$

$$\frac{\Gamma \vdash i \quad \Gamma \vdash il}{\Gamma \vdash Seq(i, il)}$$

Implémentation: DeclarationList.java III

```

}

```

Implémentation: InstructionList.java I

```

abstract class InstructionList{
    abstract void checkTypes(TypeEnvironment env) throws TypeException;
}

class Nil extends InstructionList{
    public Nil(){
    }
    public void checkTypes(TypeEnvironment env){
    }
}

class Seq extends InstructionList{
    private Instruction head;
    private InstructionList rest;
}

```



Implémentation: InstructionList.java II

```

public Seq(Instruction i, InstructionList il) {
    head=i;
    rest=il;
}
public void checkTypes(TypeEnvironment env) throws TypeException{
    head.checkTypes(env);
    rest.checkTypes(env);
}
}

```

Typed instructions

- ▶ Règle pour la conditionnelle:

$$\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash il_1 \quad \Gamma \vdash il_2}{\Gamma \vdash \text{Cond}(e, il_1, il_2)}$$

- ▶ Règle pour la boucle:

$$\frac{\Gamma \vdash e : \text{bool} \quad \Gamma \vdash il}{\Gamma \vdash \text{While}(e, il)}$$

Typed instructions

- ▶ Jugement $\Gamma \vdash i : i$ est correctement typé dans Γ .
- ▶ Règle pour l'affichage:

$$\frac{\Gamma \vdash e : \text{int}}{\Gamma \vdash \text{Print}(e)}$$

- ▶ Règle pour l'affectation:

$$\frac{\Gamma \vdash e : \Gamma(s)}{\Gamma \vdash \text{Affect}(s, e)}$$

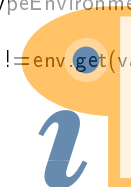
Implémentation: Instruction.java I

```

abstract class Instruction {
    abstract void checkTypes(TypeEnvironment env) throws TypeException;
}

class Assignment extends Instruction {
    private Expression expression;
    private String variable;
    public Assignment(String v, Expression e) {
        expression=e;
        variable=v;
    }
    public void checkTypes(TypeEnvironment env) throws TypeException {
        expression.setType(env);
        if (expression.getType() != env.get(variable)) {

```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Implémentation: Instruction.java II

```

        }
        }
    }

    class Print extends Instruction{
        private Expression expression;
        public Print(Expression e){
            expression=e;
        }
        public void checkTypes(TypeEnvironment env) throws TypeException{
            expression.setType(env);
            if (expression.getType() != Type.Int){
                throw new TypeException("Print");
            }
        }
    }

```

Implémentation: Instruction.java IV

```

        }
        }
    }

    class Conditional extends Instruction{
        private Expression condition;
        private InstructionList positiv;
        private InstructionList negativ;
        public Conditional(Expression e, InstructionList il1, InstructionList il2){
            condition=e;
            positiv=il1;
            negativ=il2;
        }
        public void checkTypes(TypeEnvironment env) throws TypeException{

```

Implémentation: Instruction.java III

```

        }
    }

    class While extends Instruction{
        private Expression condition;
        private InstructionList body;
        public While(Expression e, InstructionList il){
            condition=e;
            body=il;
        }
        public void checkTypes(TypeEnvironment env) throws TypeException{
            condition.setType(env);
            if (condition.getType() != Type.Bool){
                throw new TypeException("While");
            } else {

```

Implémentation: Instruction.java V

```

            condition.setType(env);
            if (condition.getType() != Type.Bool){
                throw new TypeException("Conditional");
            } else {
                positiv.checkTypes(env);
                negativ.checkTypes(env);
            }
        }
    }

```



Deux approches

- ▶ Les valeurs obtenues de l'évaluation d'une expression peuvent être de type différent (ici: *bool* ou *int*).
- ▶ Première solution: Construire un type *somme* qui réunit toutes les valeurs de tous les types.
- ▶ Deuxième solution: Des environnements de valeurs et des méthodes d'évaluation différents pour les types différents.
- ▶ La deuxième solution est plus simple tant qu'on n'a que deux types, la première solution est plus générique.

Jugements pour l'évaluation et l'exécution

- ▶ Formes des jugements pour l'évaluation:

$$\begin{aligned}\Gamma_i, \Gamma_b \vdash e \Rightarrow_i v_i \\ \Gamma_i, \Gamma_b \vdash e \Rightarrow_b v_b\end{aligned}$$

où $v_i \in \mathbb{N}$, $v_b \in \{true, false\}$.

- ▶ Jugement et règles pour le traitement des déclarations: omis.
- ▶ Formes des jugements pour les instructions et les listes d'instructions:

$$\begin{aligned}\Gamma_i, \Gamma_b \vdash i \Rightarrow \Gamma'_i, \Gamma'_b \\ \Gamma_i, \Gamma_b \vdash il \Rightarrow \Gamma'_i, \Gamma'_b\end{aligned}$$

Règles: transparents suivants.

Approche 2: utiliser deux environnements

- ▶ Deux environnements:

$$\begin{aligned}\Gamma_i &: \Sigma^* \rightarrow \mathbb{N} \\ \Gamma_b &: \Sigma^* \rightarrow \{true, false\}\end{aligned}$$

- ▶ Les deux fonctions sont partielles.
- ▶ Le type nous garantit que

$$\text{domain}(\Gamma_i) \cap \text{domain}(\Gamma_b) = \emptyset$$

Règles pour l'évaluation des expressions

$$\frac{}{\Gamma_i, \Gamma_b \vdash \text{Int}(n) \rightarrow_i n} \quad n \in \mathbb{N}$$

$$\frac{}{\Gamma_i, \Gamma_b \vdash \text{True} \rightarrow_b true}$$

$$\frac{}{\Gamma_i, \Gamma_b \vdash \text{False} \rightarrow_b false}$$

$$\frac{}{\Gamma_i, \Gamma_b \vdash \text{Ident}(s) \rightarrow_i \Gamma_i(s)}$$

$$\frac{}{\Gamma_i, \Gamma_b \vdash \text{Ident}(s) \rightarrow_b \Gamma_b(s)}$$

$s \in \Sigma^*$, si $s \in \text{domain}(\Gamma_i)$



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Règles pour l'évaluation des expressions

$$\frac{\Gamma_i, \Gamma_b \vdash e_1 \rightarrow_i v_1 \quad \Gamma_i, \Gamma_b \vdash e_2 \rightarrow_i v_2}{\Gamma_i, \Gamma_b \vdash Plus(e_1, e_2) \rightarrow_i v_1 +_{int} v_2}$$

$$\frac{\Gamma_i, \Gamma_b \vdash e_1 \rightarrow_i v_1 \quad \Gamma_i, \Gamma_b \vdash e_2 \rightarrow_i v_2}{\Gamma_i, \Gamma_b \vdash Equal(e_1, e_2) \rightarrow_b true} \quad \text{si } v_1 = v_2$$

$$\frac{\Gamma_i, \Gamma_b \vdash e_1 \rightarrow_i v_1 \quad \Gamma_i, \Gamma_b \vdash e_2 \rightarrow_i v_2}{\Gamma_i, \Gamma_b \vdash Equal(e_1, e_2) \rightarrow_b false} \quad \text{si } v_1 \neq v_2$$

et pareil pour les autres opérateurs.

Implémentation: Expression.java (début)I

```
abstract class Expression{
    Type type;
    final public Type getType(){
        return this.type;
    }
    abstract void setType(TypeEnvironment env) throws TypeException;
    //default, will be redefined for integer valued expressions
    public int evalInt(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{
        throw new EvalException();
    }
    //default, will be redefined for bool valued expressions
    public boolean evalBool(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{

```

Implémentation: Expression.java (début)II

```
        throw new EvalException();
    }
}

class Int extends Expression{
    private int value;
    public Int( int i){
        value=i;
    }
    public void setType(TypeEnvironment env){
        this.type=Type.Int;
    }
    public int evalInt(IntEnvironment ienv, BoolEnvironment benv){
        return value;
    }
}
```

Implémentation: Expression.java (début)III

```

    }

class True extends Expression{
    public True(){
    }
    public void setType(TypeEnvironment env){
        this.type=Type.Bool;
    }
    public boolean evalBool(IntEnvironment ienv, BoolEnvironment benv){
        return true;
    }
}


```

```
class Ident extends Expression{
    private String name;
```



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Implémentation: Expression.java (début)IV

```

    public Ident(Strings){
        name=s;
    }
    public void setType(TypeEnvironment env){
        this.type=env.get(name);
    }
    public int evalInt(IntEnvironment ienv, BoolEnvironment benv){
        return ienv.get(name);
    }
    public boolean evalBool(IntEnvironment ienv, BoolEnvironment benv){
        return benv.get(name);
    }
}

class Sum extends Expression{

```

Implémentation: Expression.java (début)V

```

    private Expression left;
    private Expression right;
    public Sum(Expressione1, Expressione2){
        left=e1;
        right=e2;
    }
    public void setType(TypeEnvironment env) throws TypeException{
        left.setType(env);
        right.setType(env);
        if (left.getType()==Type.Int
            && right.getType()==Type.Int){
            this.type=Type.Int;
        } else {
            throw new TypeException("Sum");
        }
    }

```

Implémentation: Expression.java (début)VI

```

    }
    public int evalInt(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{
        int value_left=left.evalInt(ienv,benv);
        int value_right=right.evalInt(ienv,benv);
        return value_left+value_right;
    }
}

class Equal extends Expression{
    private Expression left;
    private Expression right;
    public Equal(Expressione1, Expressione2){
        left=e1;
        right=e2;
    }

```

Implémentation: Expression.java (début)VII

```

    }
    public void setType(TypeEnvironment env) throws TypeException{
        left.setType(env);
        right.setType(env);
        if (left.getType()==Type.Int
            && right.getType()==Type.Int){
            this.type=Type.Bool;
        } else {
            throw new TypeException("Equal");
        }
    }
}

    public boolean evalBool(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{
        int value_left=left.evalInt(ienv,benv);
        int value_right=right.evalInt(ienv,benv);
    }

```



Implémentation: Expression.java (début) VIII

```

    return value_left == value_right;
}

```

Implémentation: InstructionList.java I

```

abstract class InstructionList {
    abstract void checkTypes(TypeEnvironment env) throws TypeException;
    abstract void exec(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException;
}

class Nil extends InstructionList {
    public Nil() {}
    public void checkTypes(TypeEnvironment env) {}
    public void exec(IntEnvironment ienv, BoolEnvironment benv) {}
}

```

Règles d'exécution de listes d'instructions

$$\overline{\Gamma_i, \Gamma_b \vdash Nil} \Rightarrow \Gamma_i, \Gamma_b$$

$$\frac{\overline{\Gamma_i, \Gamma_b \vdash i} \Rightarrow \Gamma'_i, \Gamma'_b \quad \overline{\Gamma'_i, \Gamma'_b \vdash il} \Rightarrow \Gamma''_i, \Gamma''_b}{\overline{\Gamma_i, \Gamma_b \vdash Seq(i, il)} \Rightarrow \Gamma''_i, \Gamma''_b}$$

Implémentation: InstructionList.java II

```

class Seq extends InstructionList {
    private Instruction head;
    private InstructionList rest;
    public Seq(Instruction i, InstructionList il) {
        head = i;
        rest = il;
    }
    public void checkTypes(TypeEnvironment env) throws TypeException {
        head.checkTypes(env);
        rest.checkTypes(env);
    }
    public void exec(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException {
        head.exec(ienv, benv);
    }
}

```



Implémentation: `InstructionList.java` III

```

    rest.exec(ienv, benv);
  }
}

```

Règles pour l'exécution des instructions

- Première règle pour la conditionnelle:

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_b \text{true} \quad \Gamma_i, \Gamma_b \vdash il_1 \Rightarrow \Gamma'_i, \Gamma'_b}{\Gamma_i, \Gamma_b \vdash \text{Cond}(e, il_1, il_2) \Rightarrow \Gamma'_i, \Gamma'_b}$$

- Deuxième règle pour la conditionnelle:

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_b \text{false} \quad \Gamma_i, \Gamma_b \vdash il_2 \Rightarrow \Gamma'_i, \Gamma'_b}{\Gamma_i, \Gamma_b \vdash \text{Cond}(e, il_1, il_2) \Rightarrow \Gamma'_i, \Gamma'_b}$$

Règles pour l'exécution des instructions

- Règle pour l'affichage:

$$\overline{\Gamma_i, \Gamma_b \vdash \text{Print}(e) \Rightarrow \Gamma_i, \Gamma_b}$$

- Règle pour l'affectation:

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_i v}{\Gamma_i, \Gamma_b \vdash \text{Affect}(s, e) \Rightarrow \Gamma_i[s \mapsto v], \Gamma_b}$$

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_b v}{\Gamma_i, \Gamma_b \vdash \text{Affect}(s, e) \Rightarrow \Gamma_i, \Gamma_b[s \mapsto v]}$$

Règles pour l'exécution des instructions

- Règles pour la boucle:

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_b \text{false}}{\Gamma_i, \Gamma_b \vdash \text{While}(e, il) \Rightarrow \Gamma_i, \Gamma_b}$$

$$\frac{\Gamma_i, \Gamma_b \vdash e \Rightarrow_b \text{true} \quad \Gamma_i, \Gamma_b \vdash il \Rightarrow \Gamma'_i, \Gamma'_b \quad \Gamma'_i, \Gamma'_b \vdash \text{While}(e, il) \Rightarrow \Gamma''_i, \Gamma''_b}{\Gamma_i, \Gamma_b \vdash \text{While}(e, il) \Rightarrow \Gamma''_i, \Gamma''_b}$$



Implémentation: Instruction.java I

```

abstract class Instruction{
    abstract void checkTypes(TypeEnvironment env) throws TypeException;
    abstract void exec(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException;
}

class Assignment extends Instruction{
    private Expression expression;
    private String variable;
    public Assignment(String v, Expression e){
        expression=e;
        variable=v;
    }
    public void checkTypes(TypeEnvironment env) throws TypeException{

```

Implémentation: Instruction.java II

```

        expression.setType(env);
        if (expression.getType()!=env.get(variable)){
            throw new TypeException("Assignment└to└"+variable);
        }
    }
    public void exec(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{
        if (expression.getType()==Type.Int){
            ienv.put(variable,expression.evalInt(ienv,benv));
        } else {
            benv.put(variable,expression.evalBool(ienv,benv));
        }
    }
}

```

Implémentation: Instruction.java III

```

class Print extends Instruction{
    private Expression expression;
    public Print(Expression e){
        expression=e;
    }
    public void checkTypes(TypeEnvironment env) throws TypeException{
        expression.setType(env);
        if (expression.getType()!=Type.Int){
            throw new TypeException("Print");
        }
    }
    public void exec(IntEnvironment ienv, BoolEnvironment benv)
        throws EvalException{
        System.out.println(expression.evalInt(ienv,benv));
    }
}

```

Implémentation: Instruction.java IV

```

}

class Conditional extends Instruction{
    private Expression condition;
    private InstructionList positiv;
    private InstructionList negativ;
    public Conditional(Expression e, InstructionList il1, InstructionList il2){
        condition=e;
        positiv=il1;
        negativ=il2;
    }
    public void checkTypes(TypeEnvironment env) throws TypeException{
        condition.setType(env);
        if (condition.getType()!=Type.Bool){
            throw new TypeException("Conditional");
        }
    }
}

```



Implémentation: Instruction.java V

```

    } else {
        positiv.checkTypes(env);
        negativ.checkTypes(env);
    }
}
public void      exec(IntEnvironment ienv, BoolEnvironment benv)
    throws      EvalException{
    if  (condition.evalBool(ienv,benv)){
        positiv.exec(ienv,benv);
    } else {
        negativ.exec(ienv,benv);
    }
}
}

```

Implémentation: Instruction.java VI

```

class While extends Instruction{
    private Expression condition;
    private InstructionList body;
    public While(Expression e, InstructionList il){
        condition=e;
        body=il;
    }
    public void      checkTypes(TypeEnvironment env)      throws      TypeException{
        condition.setType(env);
        if  (condition.getType()!=Type.Bool){
            throw new      TypeException("While");
        } else {
            body.checkTypes(env);
        }
    }
}

```

Implémentation: Instruction.java VII

```

public void      exec(IntEnvironment ienv, BoolEnvironment benv)
    throws      EvalException{
    if  (condition.evalBool(ienv,benv)){
        body.exec(ienv,benv);
        this .exec(ienv,benv);
    }
}
}

```



Analyse de Données Structurées- Cours 12

Ralf Treinen



Université Paris Diderot
UFR Informatique
Laboratoire Preuves, Programmes et Systèmes

treinen@pps.univ-paris-diderot.fr

6 mai 2015

© Ralf Treinen 2015

Exemple d'un programme

```
var
  x: integer;
  y: integer;
  b: boolean;
  c: boolean;
begin
  x := 73;
  if b & c || !b then
    print 42 + 1 * x + y;
  y := 42;
  else
    y := x + x;
  fi;
  while 42 > 41 do print 73; od;
end
```

Rappel

- ▶ Jusqu'à maintenant: un petit langage de programmation impératif
 - | Deux types: `integer` et `boolean`
 - | Expressions avec des opérateurs arithmétiques et logiques
 - | Déclaration des variables avec leur type au début du programme
 - | Instructions d'affichage et d'affectation
 - | Instructions composées: boucle `while...do...od` , conditionnelle `if...then...else...fi`

Un langage avec des déclarations locales

- ▶ Extension du langage: déclarations locales
- ▶ On veut maintenant aussi permettre des déclarations de variables avec portée limitée.
- ▶ Pour cela il y a une décision fondamentale à prendre:

est-ce qu'on veut permettre une déclaration locale d'une variable?



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Exemple: variable locale

```
1 var
2 x: integer;
3 y: integer;
4 begin
5   x := 73;
6   y := x + 42;
7   var
8     z: integer;
9   begin
10    z := x + y;
11    print(z);
12  end;
13  print(x + y);
14 end
```

Java: redéclaration pas autorisée

```
/* Erreur de compilation */
class Scope {
    public static void main(String[] args) {
        int x = 42;
        {
            int x = 73;
        }
    }
}
```

Exemple: redéclaration d'une variable

```
1 var
2 x: integer;
3 y: integer;
4 begin
5   x := 73;
6   y := x + 42;
7   var
8     x: boolean;
9   begin
10    x := y > 17;
11    if x then print(1); else print(2); fi;
12  end;
13  print(x + y);
14 end
```

Redéclaration de variables

- ▶ En Java: la redéclaration d'une variable dans un bloc n'est pas permise (erreur de compilation).
- ▶ Il y a des autres langages de programmation qui admettent des redéclarations locales de variable.
- ▶ Les redéclarations sont aussi pertinentes pour l'implémentation de procédures qui peuvent accéder à des variables globales.
- ▶ Pour cette raison nous allons étudier les deux solutions.



Extension de la syntaxe abstraite

- L'ensemble DL des Listes de Déclaration est ...
- L'ensemble IL des listes d'Instruction est ...
- L'ensemble I des Instructions en syntaxe abstraite est:
 - └ si $e \in E'''$ alors $Print(e) \in I$,
 - └ si $s \in *$, $e \in E'''$ alors $Assign(s, e) \in I$,
 - └ si $e \in E'''$, $il_1, il_2 \in IL$ alors $Cond(e, il_1, il_2) \in I$,
 - └ si $e \in E'''$, $il \in IL$ alors $While(e, il) \in I$,
 - └ si $dl \in DL$ et $il \in IL$ alors $Bloc(dl, il) \in I$.

Modification: exécution de déclarations

- Jugement $dl \Rightarrow i, b$: l'exécution de la liste de déclarations dl crée un environnement entier i , et un environnement booléen b .
- Règles:

$$\overline{DeclNil \Rightarrow \emptyset, \emptyset}$$

$$\frac{dl \Rightarrow i, b}{DeclSeq(Decl(s, int), dl) \Rightarrow i[s \mapsto 0], b}$$

$$\frac{dl \Rightarrow i, b}{DeclSeq(Decl(s, bool), dl) \Rightarrow i, b[s \mapsto false]}$$

Petite modification de la sémantique

- Jusqu'à maintenant, une variable n'a pas reçu une valeur initiale au moment de la déclaration.
- Modification: au moment de la déclaration, la variable reçoit une valeur initiale (0 pour les entiers, *false* pour les booléens).
- Raison: simplification technique.
- On aurait aussi pu étendre la syntaxe abstraite par une expression d'initialisation au moment de la déclaration d'une variable.

Modification: exécution d'un programme

- Jugement $i, b \vdash il \Rightarrow i', b'$: l'exécution de la liste d'instructions il dans un environnement initial i, b donne un nouvel environnement i', b'
- Règle modifiée pour un programme complet:

$$\frac{dl \Rightarrow i, b \quad i, b \vdash il \Rightarrow i', b'}{Prog(dl, il) \Rightarrow i', b'}$$



Deux opérations sur les environnements

- Union disjointe: si Γ_1 et Γ_2 sont deux environnements (de typage ou de valeur) avec des *domaines disjoints*, alors $\Gamma_1 \oplus \Gamma_2$ est défini comme

$$(\Gamma_1 \oplus \Gamma_2)(x) = \begin{cases} \Gamma_1(x) & \text{si } x \in \text{domaine}(\Gamma_1) \\ \Gamma_2(x) & \text{sinon} \end{cases}$$

- Restriction: si Γ est un environnement, et $D \subseteq \text{domaine}(\Gamma)$, alors la restriction de Γ à D est défini comme

$$\Gamma|_D(x) = \begin{cases} \Gamma(x) & \text{si } x \in D \\ \text{non-défini} & \text{sinon} \end{cases}$$

Exécution des blocs

$$\frac{dl \Rightarrow \Gamma'_i, \Gamma'_b \quad \Gamma'_i \oplus \Gamma'_i, \Gamma'_b \oplus \Gamma'_b \vdash il \Rightarrow \Gamma''_i, \Gamma''_b}{\Gamma'_i, \Gamma'_b \vdash \text{Bloc}(dl, il) \Rightarrow \Gamma''_i|_{\text{domaine } \Gamma'_i}, \Gamma''_b|_{\text{domaine } \Gamma'_b}}$$

Typage pour les blocs

- Rappel des jugements de typage:
 - $dl : \Gamma$: la liste de déclarations Γ donne lieu à l'environnement de typage Γ .
 - $\vdash il : \tau$: la liste d'instructions il est bien typée par rapport à l'environnement Γ .

- Règle

$$\frac{dl \Rightarrow \Gamma' \quad \vdash il : \tau}{\vdash \text{Bloc}(dl, il) : \tau}$$

Une modélisation plus complexe

- La technique des extensions et restrictions des environnements ne fonctionne plus dans le cas des redéclarations de variable.
- Raison: Une redéclaration d'une variable x masque, pour la durée d'exécution du bloc, une déclaration globale de x . Après avoir terminé le bloc, la variable x globale devient de nouveau visible.
- Nous avons donc maintenant besoin d'une liste d'environnements: quand on entre dans un bloc, un nouvel environnement est ajouté à la fin de la liste; quand on quitte un bloc, l'environnement le plus récent est supprimé.



Edited with the demo version of
Infix Pro PDF Editor

To remove this notice, visit:
www.iceni.com/unlock.htm

Listes d'environnements

- ▶ Dans une liste d'environnements, l'environnement le plus global est trouvé à gauche, et l'environnement le plus local est trouvé à droite.
- ▶ Pour obtenir la valeur d'une variable, ou pour la modifier, on parcourt la liste de droite à gauche : priorité à la déclaration la plus locale.
- ▶ Pour ne pas trop alourdir la notation nous ignorons maintenant la distinction entre environnement entier et environnement booléen.

Règles à modifier pour les listes d'environnements

- ▶ Évaluation d'une variable

$$\frac{}{1, \dots, n \vdash \text{Ident}(s) \rightarrow \text{get}(s, (1, \dots, n))}$$

- ▶ Affectation d'une variable

$$\frac{1, \dots, n \vdash e \Rightarrow v}{1, \dots, n \vdash \text{Assign}(s, e) \Rightarrow \text{set}(s, v, (1, \dots, n))}$$

Opérations sur les listes d'environnements

- ▶ L'opération *get* permet de récupérer la valeur d'une variable:

$$\text{get}(s, (1, \dots, n)) = \begin{cases} n(s) & \text{si } s \in \text{domaine}(n) \\ \text{get}(s, (1, \dots, n-1)) & \text{sinon} \end{cases}$$

- ▶ L'opération *set* permet de modifier la valeur d'une variable:

$$\text{set}(s, v, (1, \dots, n)) = \begin{cases} (1, \dots, n-1, n[s \mapsto v]) & \text{si } s \in \text{domaine}(n) \\ \text{set}(s, v, (1, \dots, n-1)), n & \text{sinon} \end{cases}$$

Nouvelle règle pour l'exécution d'un bloc

$$\frac{dl \Rightarrow n-1 \quad 1, \dots, n, n-1 \vdash il \Rightarrow l'_1, \dots, l'_n, l'_{n-1}}{1, \dots, n \vdash \text{Bloc}(dl, il) \Rightarrow l'_1, \dots, l'_n}$$



Extensiondelasyntaxe

- ▶ Syntaxeconcrète:définitiond'uneprocédure(àunargument)
par
proc *nom(variable : type)*
bloc
- ▶ Syntaxeabstraite:
 - | définitiond'uneprocédure:nomdelaprocédure,nomdu paramètre,typeduparamètre,bloc.
 - | listededéfinitionsdeprocédures: ...
 - | bloc:listededéclarationsdevariables,listededéfinitionsde procédures,bloc.

Exemple:procédureaccédantàunevariableglobale

```
1 var x: integer;  
2 proc p(y: integer)  
3   var z: integer;  
4   begin  
5     z := 2 * y;  
6     x := z;  
7   end  
8   begin  
9     x := 5;  
10  p(x);  
11  print x;  
12 end
```

Exemple:procédure

```
1 var x: integer;  
2 proc p(y: integer)  
3   var z: integer;  
4   begin  
5     z := 2 * y;  
6     print z;  
7   end  
8   begin  
9     x := 5;  
10  p(x);  
11 end
```

Exemple:liaisonstatiqueoudynamique

```
1 var x: integer;  
2 proc p(y: integer)  
3   begin  
4     x := x + y;  
5   end  
6   begin  
7     x := 1;  
8     var x: integer;  
9     begin  
10      x := 2;  
11      p(1);  
12      print x;  
13    end  
14 end
```



Liaisonstatiqueoudynamique

- ▶ Liaisonenquestion :laconnexionentrel'utilisationd'unnom (d'unevariable,parex.),etladéclarationdumêmenom.
- ▶ Problème: onpeutavoirplusieursdéclarationsdumêmenom, laquellechoisir?
- ▶ Liaisondynamique: laliaisonestfaitesuivantlaprioritéàla déclarationlapluslocale *aumomentdel'exécution* .
- ▶ Liaisonstatique: laliaisonestfaitesuivantlaprioritéàla déclarationlapluslocale *aumomentdeladéfinitiondela procédure*.
Tousleslangagesdeprogrammationmodernessuiventcette politique.

