



POLYTECHNIQUE
MONTREAL

LE GÉNIE
EN PREMIÈRE CLASSE

LOG3210

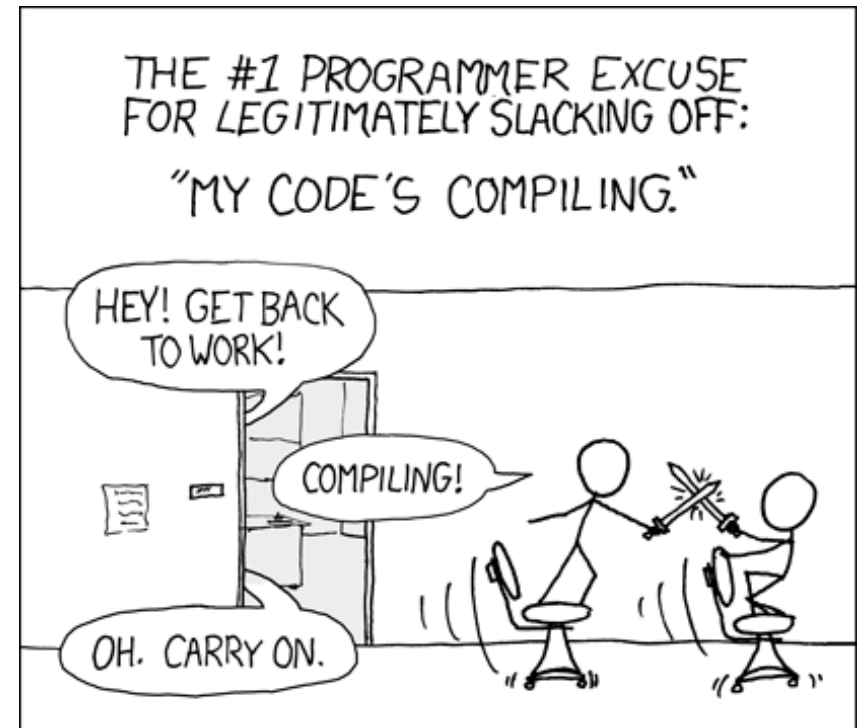
Éléments de langage et compilateurs

Introduction
Analyse lexicale
Grammaires

Semaine 1

PLAN

1. Présentation du cours
2. Chapitre 3 – Analyse lexicale
3. Chapitre 2 – Grammaires



QUI PROGRAMME EN C++ ?



- Génération de code intermédiaire (e.g. LLVM-IR)
- Génération de code machine (ou d'assembleur)
- Allocation de registres
- Optimisations (si possible)



QUI PROGRAMME EN JAVA ?



- Génération de code intermédiaire (Java *bytecode*)
- Gestion de la mémoire (*garbage collection*)



QUI PROGRAMME EN ML ?

Programmation fonctionnelle

- ML, Haskell, Lisp, R
- Environnements run-time
 - *Nested procedures*
 - Chapitre 7 : *Access links*



QUI PROGRAMME EN PHP ?



Langages interprétés

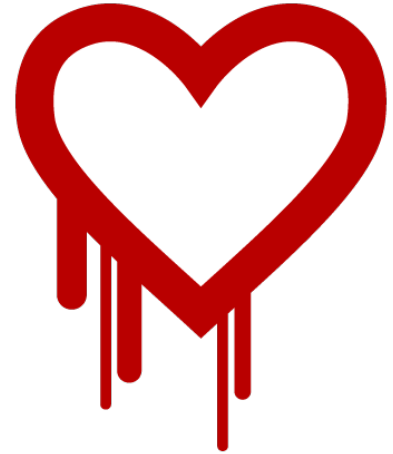
- PHP, JavaScript, Python, Ruby
- Analyse lexicale et syntaxique
 - Exécution du code « *on-the-fly* »
- Analyse de flux (LOG6302 Réingénierie)



ANALYSE STATIQUE DE PROGRAMMES

Heartbleed

```
hbtype = *p++;  
n2s(p, payload);  
p1 = p;
```



ANALYSE STATIQUE DE PROGRAMMES

Heartbleed : Boundary check!

```
/* Read type and payload length first */  
if (1 + 2 + 16 > s->s3->rrec.length)  
    return 0; /* silently discard */  
hbtype = *p++;  
n2s(p, payload);  
if (1 + 2 + payload + 16 > s->s3->rrec.length)  
    return 0; /* silently discard per RFC 6520  
    sec. 4 */  
p1 = p;
```



POURQUOI UN COURS DE COMPILATEURS ?

- Pour comprendre un outil indispensable au programmeur :
 - Pourquoi est-ce que le `;` manquant n'est détecté que plusieurs lignes plus loin ?
 - Quelles sont les optimisations d'un compilateur ?
 - Pourquoi utiliser un compilateur plutôt qu'un autre ?
 - Quel est l'impact de la *garbage collection* sur un programme ?



POURQUOI UN COURS DE COMPILATEURS ?

- Pour écrire un nouveau langage :
 - *Domain Specific Language* (DSL) ;
 - Systèmes embarqués.
(*Cross compilation*, port de Clang sur ARM)
- Pour faire de l'analyse de programmes :
 - Intégration dans un IDE ;
 - Analyses de sécurité.



POURQUOI UN COURS DE COMPILATEURS ?

Parce que c'est recherché dans l'industrie ! (\$\$)



LOG3210

Cours 1

Introduction et grammaires

Compilateurs vs Interpréteurs

- ▶ Quelle est la différence?

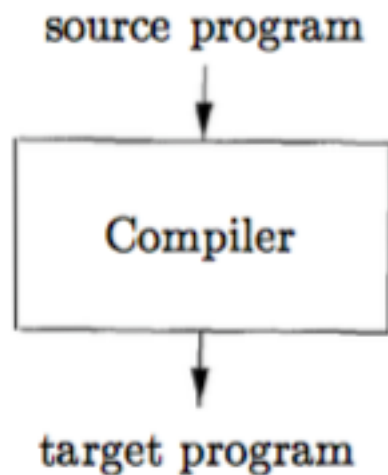


Figure 1.2: Running the target program

Figure 1.1: A compiler

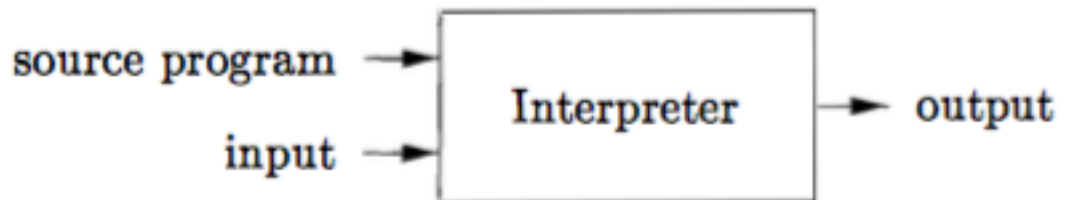


Figure 1.3: An interpreter

Compilateurs vs Interpréteurs

► Avantages – Inconvénients

- Rapidité ? (Compilateur)
- Diagnostiques / traitements des erreurs ?
- Portabilité ? (Interpréteur)

► Exemples

- C++ vs Python
- Langages web

► Approches mixtes

- Java (Compilation et interprétation)
- Compilateurs *just-in-time* (JIT)

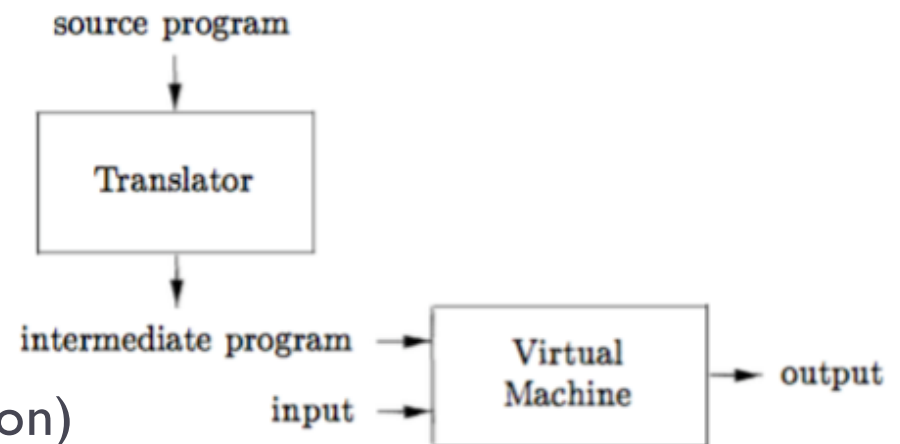
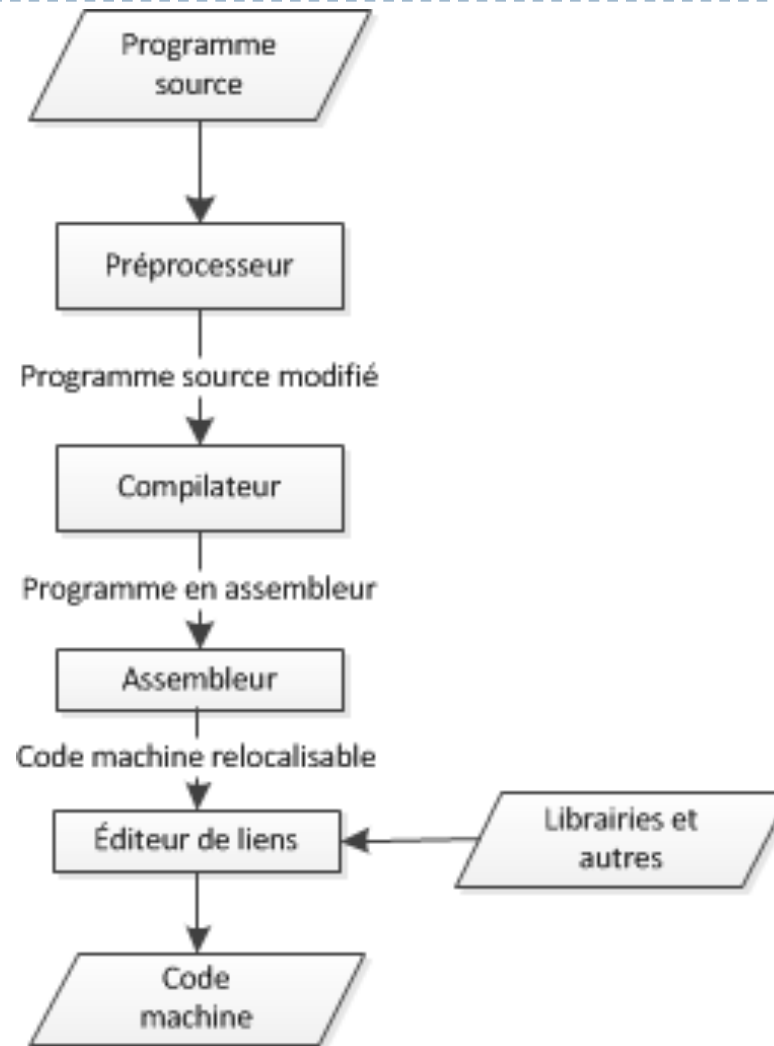
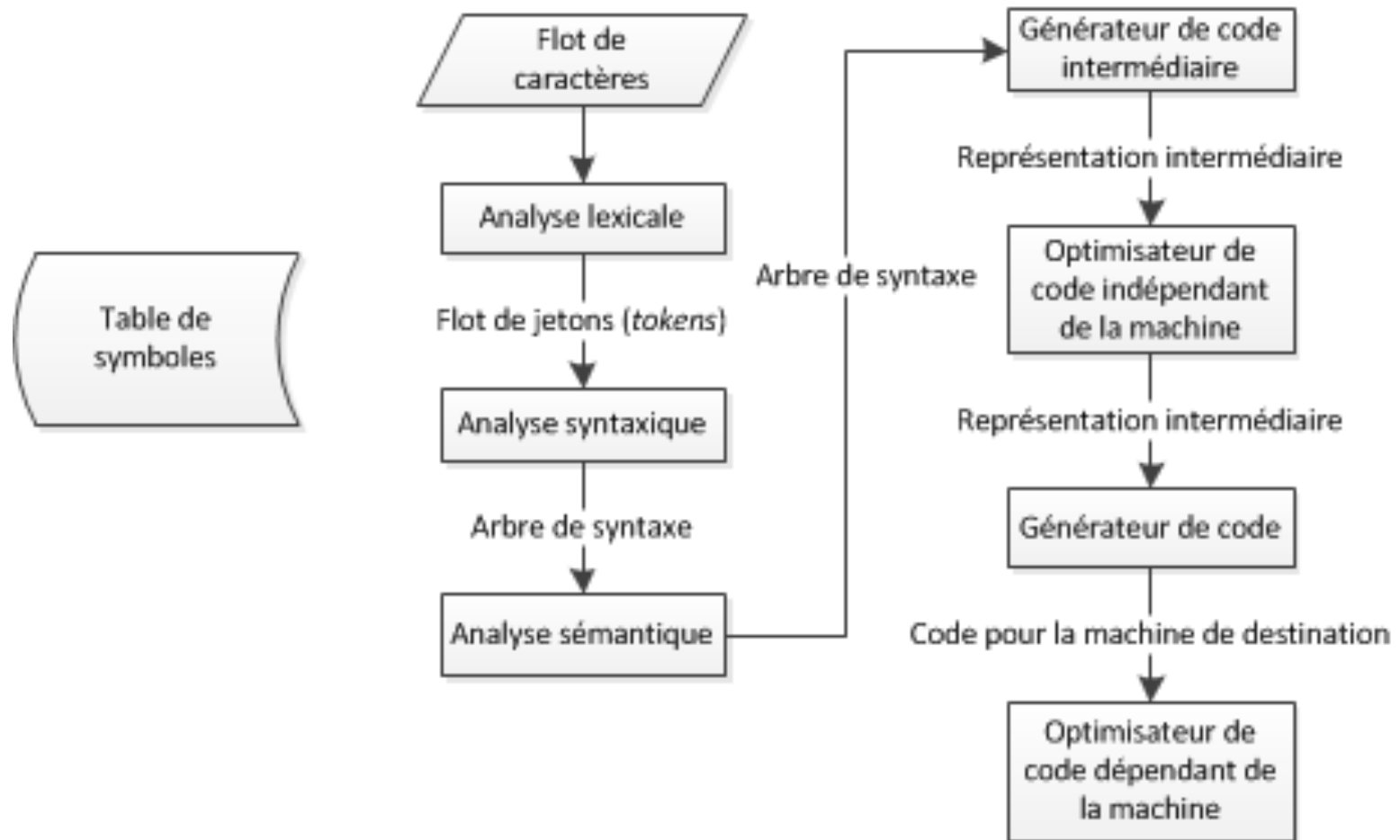


Figure 1.4: A hybrid compiler

Compilation d'un langage typique



Phases d'un compilateur



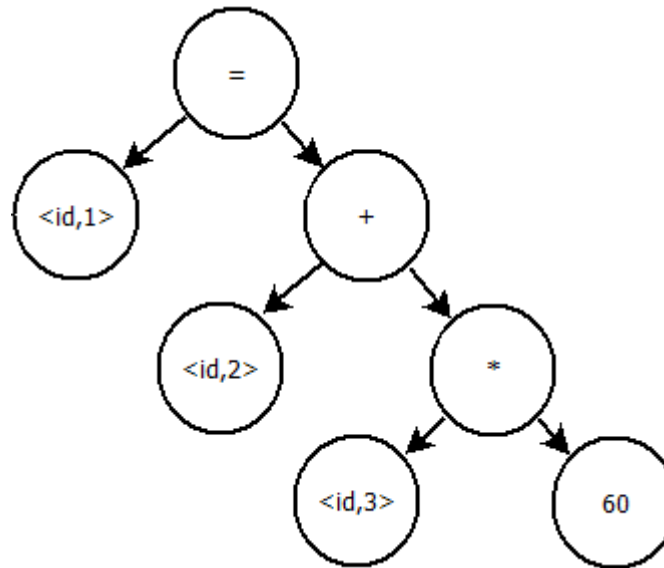
Analyse lexicale

- ▶ Permet de lire le programme source et produire des jetons (*tokens*) qui seront insérés dans la table de symboles.
- ▶ Par exemple:
 - ▶ *Position = initial + rate * 60*
- ▶ Devient:
 - ▶ **<id, 1>** **<=>** **<id, 2>** **<+>** **<id, 3>** **<*>** **<60>**

Analyse syntaxique

- ▶ Communément appelé *parsage* (*parsing*).
- ▶ Crée un arbre de syntaxe qui permet de connaître l'ordre dans lequel les expressions sont exécutées.

▶ Exemple:



Analyse sémantique

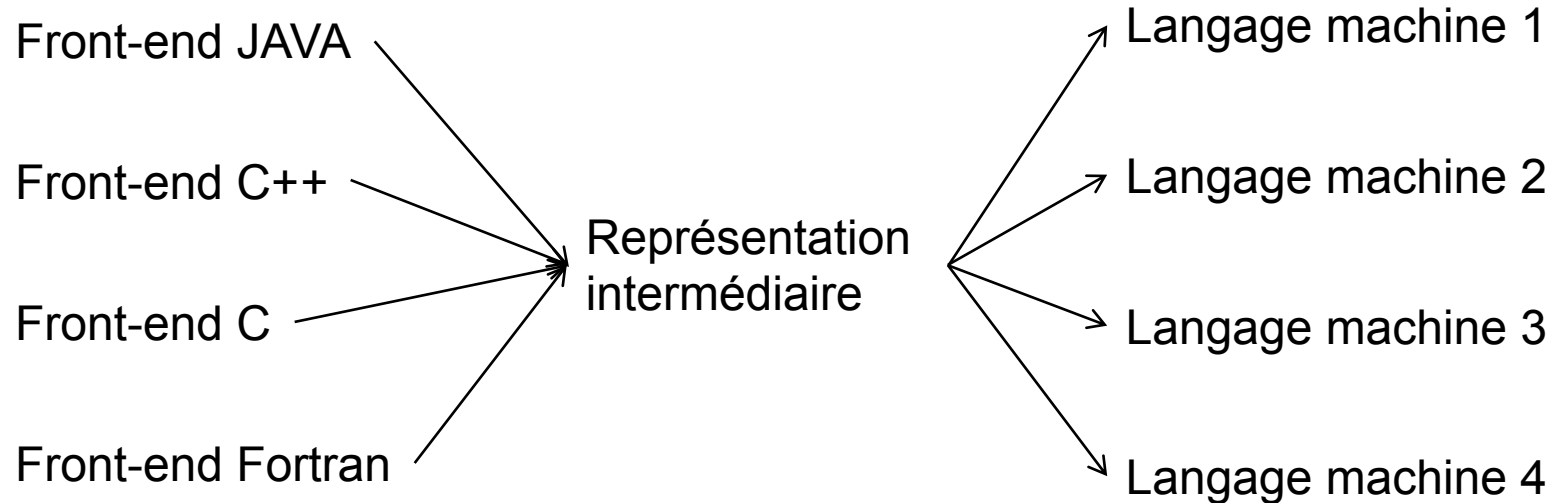
- ▶ Vérifie que le programme en entrée est consistant avec le langage défini.
- ▶ Effectue les diverses vérifications sémantiques telles que:
 - ▶ Vérifications des types
 - ▶ Vérification de l'initialisation des variables
- ▶ Renvoie les erreurs et avertissements de compilation, si nécessaire.

Génération de code

- ▶ Utilise l'arbre de syntaxe pour générer
 - ▶ du code intermédiaire et,
 - ▶ éventuellement, le code machine.

Génération de code intermédiaire

- But de la représentation intermédiaire: interface entre le front-end (spécifique au langage) et le back-end (spécifique à la machine).



Écrire un compilateur

- ▶ L'écriture d'un compilateur est une tâche très complexe.
- ▶ L'optimisation du compilateur doit:
 - ▶ Être correcte (conserver l'intention du code original)
 - ▶ Augmenter la performance de plusieurs programmes
 - ▶ Maintenir un temps de compilation raisonnable
 - ▶ Garder le compilateur relativement simple pour permettre sa maintenance
- ▶ Les compilateurs avec optimisations contiennent parfois des erreurs. (e.g. Option `-O4` sur gcc)

Applications des compilateurs

- ▶ Permettre le développement de langages de programmation de haut niveau.
- ▶ Optimiser le code pour diverses architectures
 - ▶ Parallélisme
 - ▶ Gestion des niveaux de mémoire
- ▶ Traduction de programmes (SQL, traduction binaires...)
- ▶ Outils de productivité

Analyse lexicale

Analyse lexicale

- ▶ L'analyseur lexicale lit la chaîne en entrée et la convertit en unités lexicales: les jetons (*tokens*).
- ▶ Deux tâches principales:
 - ▶ *Scanning* : tâches qui ne nécessitent pas la création de jetons. Par exemple, supprimer les commentaires et les espaces inutiles.
 - ▶ Analyse lexicale: produire les jetons

Analyse lexicale - Jetons

- ▶ Pour la plupart des langages de programmation, les classes de jetons suivantes sont suffisantes:
 - ▶ Un jeton pour chaque mot clé (**if**, **else**, etc.)
 - ▶ Un jeton pour chaque opérateur
 - ▶ Un jeton pour les identificateurs
 - ▶ Un ou plusieurs jetons pour les constantes (nombres, chaînes de caractères...)
 - ▶ Un jeton pour chaque symbole de ponctuation (point, parenthèses, point-virgule...)

Analyse lexicale - Jetons

TOKEN	INFORMAL DESCRIPTION	SAMPLE LEXEMES
if	characters <code>i</code> , <code>f</code>	<code>if</code>
else	characters <code>e</code> , <code>l</code> , <code>s</code> , <code>e</code>	<code>else</code>
comparison	<code><</code> or <code>></code> or <code><=</code> or <code>>=</code> or <code>==</code> or <code>!=</code>	<code><=</code> , <code>!=</code>
id	letter followed by letters and digits	<code>pi</code> , <code>score</code> , <code>D2</code>
number	any numeric constant	<code>3.14159</code> , <code>0</code> , <code>6.02e23</code>
literal	anything but <code>"</code> , surrounded by <code>"</code> 's	<code>"core dumped"</code>

Figure 3.2: Examples of tokens

Attributs des jetons

- ▶ Les jetons peuvent avoir un attribut permettant de les différencier.
 - ▶ Pointeur vers la table de symboles
 - ▶ Valeur d'une constante numérique
- ▶ Par exemple, $E = M * C ** 2$ devient:

```
<id, pointer to symbol-table entry for E>  
<assign_op>  
<id, pointer to symbol-table entry for M>  
<mult_op>  
<id, pointer to symbol-table entry for C>  
<exp_op>  
<number, integer value 2>
```

Spécification des jetons

- ▶ On peut spécifier les jetons via des expressions régulières:

$$\begin{aligned} \text{letter_} &\rightarrow A \mid B \mid \dots \mid Z \mid a \mid b \mid \dots \mid z \mid _ \\ \text{digit} &\rightarrow 0 \mid 1 \mid \dots \mid 9 \\ \text{id} &\rightarrow \text{letter_} (\text{letter_} \mid \text{digit})^* \end{aligned}$$

- ▶ **Exercice** : Donnez le langage spécifié par les expressions suivantes, en considérant les terminaux $T = \{a, b\}$:

1. $a \mid b$
2. $(a \mid b) (a \mid b)$
3. a^*
4. $(a \mid b)^*$
5. $a \mid a^* b$

Reconnaissance de jetons

- Le but de l'analyseur lexicale est de transformer l'entrée en une chaîne de jetons, tel quel spécifié dans à la figure 3.12 :

LEXEMES	TOKEN NAME	ATTRIBUTE VALUE
Any <i>ws</i>	–	–
if	if	–
then	then	–
else	else	–
Any <i>id</i>	id	Pointer to table entry
Any <i>number</i>	number	Pointer to table entry
<	relop	LT
<=	relop	LE
=	relop	EQ
<>	relop	NE
>	relop	GT
>=	relop	GE

Figure 3.12: Tokens, their patterns, and attribute values

Reconnaissance de jetons

- La reconnaissance des jetons s'effectue via un automate NFA (non-déterministe) ou DFA (déterministe)

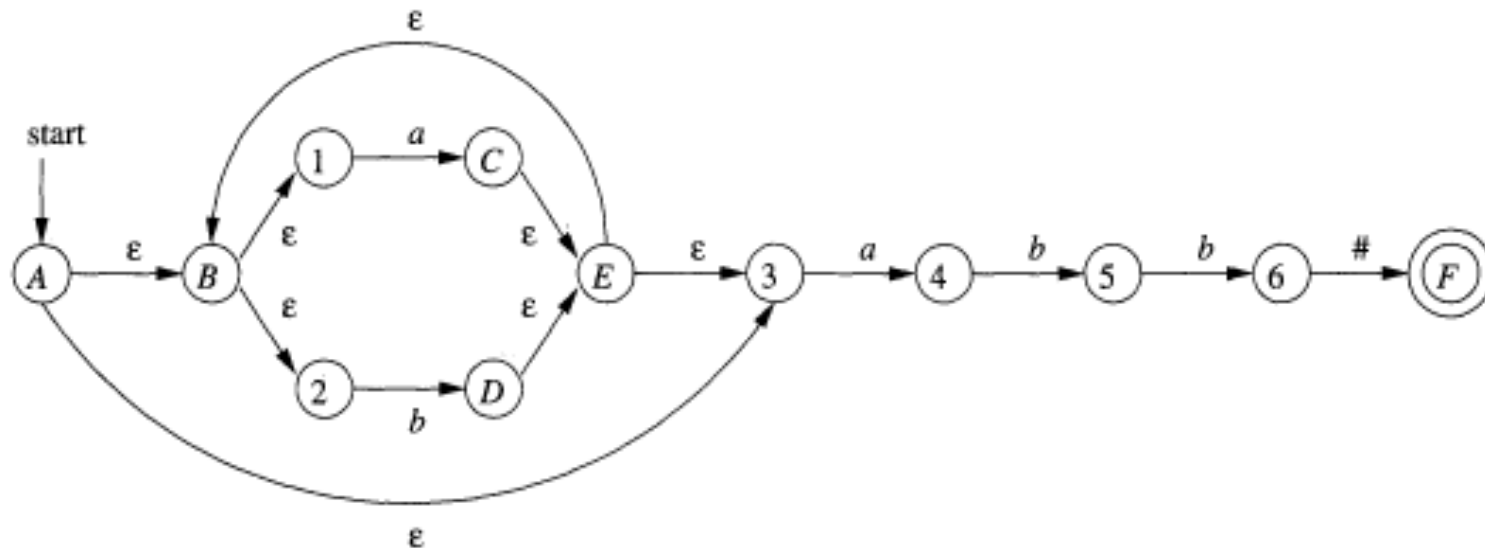


Figure 3.57: NFA constructed by Algorithm 3.23 for $(a|b)^*abb\#$

Grammaires

Grammaires

- ▶ Une grammaire décrit la structure hiérarchique d'un langage.
- ▶ Exemple:
 - ▶ **if** (expression) statement **else** statement
- ▶ Règle de production:
 - ▶ *stmt* → **if** (*expr*) *stmt* **else** *stmt*

Grammaires libres de contexte (CFG)

- ▶ $CFG = (N, T, P, S)$:
 - ▶ Symboles non terminaux (N)
 - ▶ Symboles terminaux (T)
 - ▶ Règles de productions P du type
 - ▶ $a \rightarrow b, \quad a \in N \quad \text{et} \quad b \in (N \cup T)$
 - ▶ Exemple:
 - $E \rightarrow E + E \mid E * E \mid (E) \mid -E \mid id$
 - ▶ Un symbole de départ $S \in N$

Grammaires libres de contexte (CFG)

- ▶ **Exemple** : Grammaire qui reconnaît une chaîne de "a", suivie d'une addition (+), suivie d'une chaîne de "b".
- ▶ $CFG = (N, T, P, S)$:
 - ▶ Symboles non terminaux (N) : $N = \{S, A, B\}$
 - ▶ Symboles terminaux (T) : $T = \{a, b, +\}$
 - ▶ Règles de productions P :
 - ▶ $S \rightarrow A + B$
 - ▶ $A \rightarrow aA \mid a$
 - ▶ $B \rightarrow bB \mid b$
 - ▶ Un symbole de départ $S \in N$
- ▶ **Question** : Combien de règles de productions ?

Exercice 1.1

- ▶ En utilisant les symboles suivants, écrivez une grammaire permettant de reconnaître la multiplication et la division d'un nombre aléatoire de chiffres.

Par exemple, la chaîne $8*2*4/2*9/1/2$ doit être reconnue.

- ▶ Symboles non terminaux: *list*, *digit*
- ▶ Symboles terminaux: $*$, $/$, les chiffres de 0 à 9

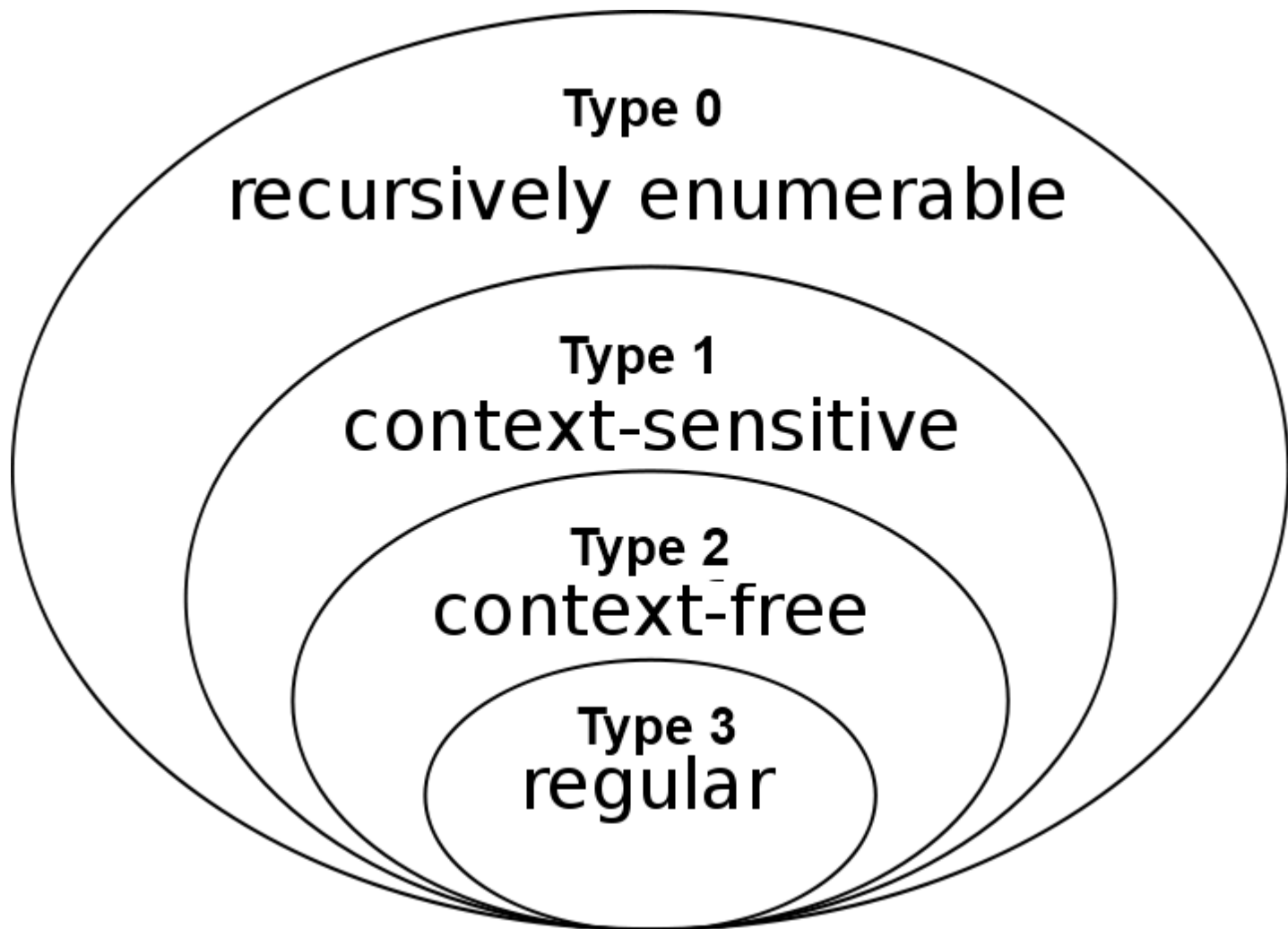
Représentation de CFG

- ▶ Backus-Naur Form (BNF)
- ▶ Extended BNF (EBNF)
 - ▶ Groupage: (...), 1 occurrence
 - ▶ Parties optionnelles: [...], 0 ou 1 occurrence
 - ▶ {...}+, 1 ou plus occurrences
 - ▶ {...}n, 1 à n occurrences
 - ▶ {...}*, 0 à plus occurrences
- ▶ JavaCC permet d'utiliser la notation EBNF

Langages libres de contexte

- ▶ Un langage libre de contexte (CFL) est généré par une CFG.
 - ▶ Un langage est généré par une grammaire.
- ▶ Il existe des langages non libres de contexte.
 - ▶ Par exemple, C/C++ !

Types de grammaires



Types de grammaires

► Grammaires régulières (type 3)

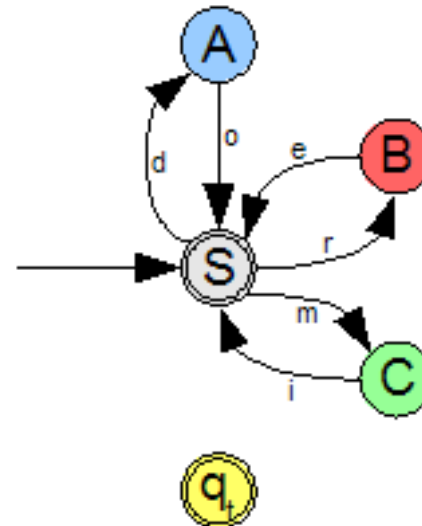
- Ne peuvent comprendre que des règles de production ayant la forme $A \rightarrow a$ ou $A \rightarrow aB$
- Il s'agit de tous les langages reconnaissables par un automate à fini (*Finite State Machine*).
- Tous les langages réguliers peuvent être obtenus via une expression régulière.
- Exemple:

$$S \rightarrow dA | rB | mC | \epsilon$$

$$A \rightarrow oS$$

$$B \rightarrow eS$$

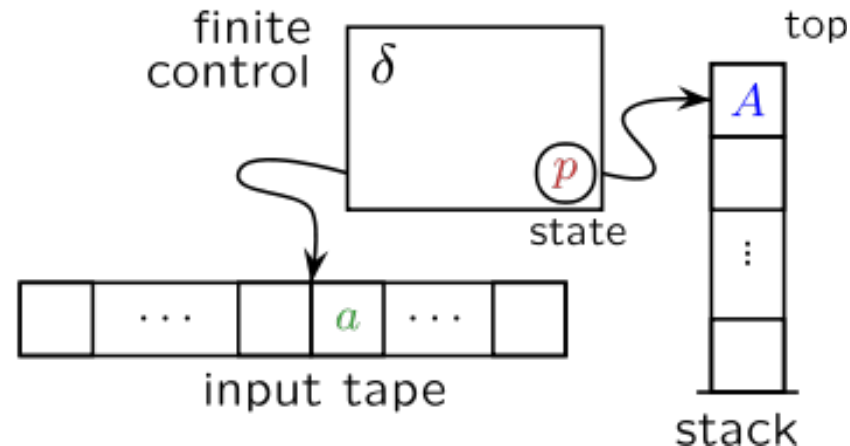
$$C \rightarrow iS$$



Types de grammaires

► Grammaires libres de contexte (type 2)

- Règles de la forme $A \rightarrow \gamma$, où A est un non terminal et γ est un ensemble de terminaux et de non terminaux.
- Il s'agit de tous les langages reconnaissables par un automate à pile (*pushdown automaton* ou parseur ascendant).
- La plupart des langages de programmation sont libres de contexte.



Types de grammaires

▶ **Grammaires contextuelles (type I)**

- ▶ Règles de la forme $\alpha A \beta \rightarrow \alpha \gamma \beta$ où A est un non terminal et α , β et γ sont un ensemble de terminaux et de non terminaux. α et β peuvent être vides, mais γ est non-vide.
- ▶ Il s'agit de tous les langages reconnaissables par un automate linéairement borné (forme restreinte d'une machine de Turing non déterministe).

▶ **Grammaires syntagmatiques (ou non restreintes) (type 0)**

- ▶ Aucune restriction sur les règles.
- ▶ Génèrent tous les langages reconnaissables par une machine de Turing

Types de grammaires

- ▶ **Grammaires régulières (Type 3)**
 - ▶ Reconnues par l'analyse lexicale (automate)
- ▶ **Grammaires libres de contexte (Type 2)**
 - ▶ Reconnues par l'analyse syntaxique
- ▶ **Grammaires contextuelles (Type 1)**
 - ▶ Partiellement reconnues par l'analyse sémantique (table de symboles)
- ▶ **Grammaires non restreintes (Type 0)**
 - ▶ Machine de Turing

Exemples

Langages non libres de contexte

- ▶ $L1 = \{wcw \mid w \in (a|b)^*\}$
 - ▶ Problème de vérifier la déclaration d'un identificateur de longueur arbitraire avant son utilisation
 - ▶ Solution: table des symboles et routines associées
- ▶ $L2 = \{a^n b^m c^n d^m \mid n \geq 1, m \geq 1\}$
 - ▶ Problème de vérifier la correspondance entre paramètres formels et paramètres actuels dans les définitions et les appels à procédures
 - ▶ Solution: analyse sémantique

Arbre d'analyse syntaxique

- ▶ L'arbre d'analyse syntaxique (ou arbre de parsage) montre comment le symbole de départ d'une grammaire dérive une chaîne du langage.
- ▶ L'arbre de parsage représente la structure syntaxique d'une chaîne en entrée par rapport à une grammaire.
- ▶ **Exemple :** 9-5+2 pour la grammaire suivante
 - ▶ $list \rightarrow list + digit$
 - ▶ $list \rightarrow list - digit$
 - ▶ $list \rightarrow digit$
 - ▶ $digit \rightarrow 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

Exercice 1.2

- ▶ En utilisant les symboles suivants, écrivez une grammaire permettant de reconnaître la multiplication et la division d'un nombre aléatoire de chiffres.

Donnez un arbre de parsage pour la chaîne $8*2*4/2*9/1/2$.

- ▶ Symboles non terminaux: *list*, *digit*
- ▶ Symboles terminaux: $*$, $/$, les chiffres de 0 à 9

Exercice 2

► Soit la grammaire CFG suivante:

► $S \rightarrow AB$ $A \rightarrow AA \mid Aa \mid aa$ $B \rightarrow BB \mid Bb \mid bb$

1. Identifiez les éléments N, T, P, S de la définition de la grammaire $G = (N, T, P, S)$

2. Dites si les chaînes suivantes appartiennent au langage et, pour celles qui y appartiennent, dessinez l'arbre d'analyse syntaxique:

<i>a</i>	<i>a a b b</i>
<i>a a</i>	<i>a a a b b</i>
<i>a a a</i>	<i>a b b b</i>
<i>a b</i>	<i>a a b b b</i>
<i>a a b</i>	<i>a a a b b b</i>
<i>a a a b</i>	<i>a b a b</i>
<i>a b b</i>	<i>a a b b a a b b</i>

Grammaires ambiguës

- ▶ Une grammaire est ambiguë si elle peut produire plus qu'un arbre d'analyse syntaxique pour la même chaîne de caractères en entrée (c'est-à-dire plus que une dérivation gauche ou droite).
- ▶ Le problème de déterminer si une CFG est ambiguë est indécidable (NP-complet).
 - ▶ Par contre, il ne suffit que d'un exemple pour démontrer qu'elle est ambiguë.
- ▶ Lire les sections 2.2.4 et 4.2.5 du livre.

Éliminer l'ambiguïté

- ▶ Pour un parseur, il est préférable qu'une grammaire ne soit pas ambiguë.
- ▶ Solutions à l'ambiguïté:
 - ▶ Modifier la grammaire
 - ▶ Exemple: $A \rightarrow A + A \mid A - A \mid a$
 - ▶ Il y a deux dérivations possibles pour la chaîne « $a + a + a$ »
 - ▶ Grammaire non ambiguë pour le même langage:
 $A \rightarrow A + a \mid A - a \mid a$
 - ▶ Modifier le langage
 - ▶ Ajouter des règles d'élimination de l'ambiguïté
 - ▶ Élimination des arbres de parsage incorrects

Éliminer l'ambiguïté

Exemple

- ▶ $stmt \rightarrow$ **if** $expr$ **then** $stmt$ |
 if $expr$ **then** $stmt$ **else** $stmt$ |
 other
- ▶ Ambiguë, car la chaîne suivante a deux arbres de parsage:

If E_1 then if E_2 then S_1 else S_2

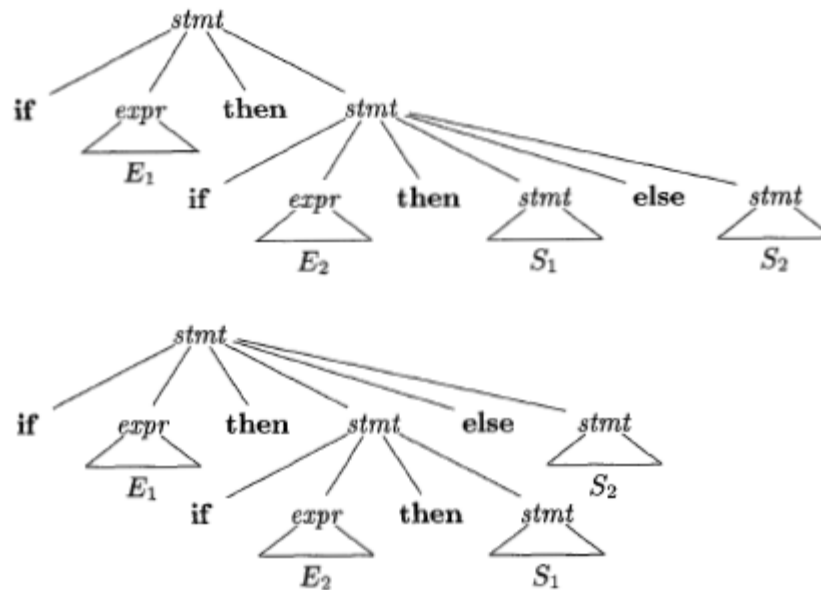


Figure 4.9: Two parse trees for an ambiguous sentence

Éliminer l'ambiguïté

Exemple

- Solution en modifiant la grammaire:

```
stmt    →  matched_stmt  
         |  open_stmt  
matched_stmt → if expr then matched_stmt else matched_stmt  
              |  other  
open_stmt  → if expr then stmt  
              |  if expr then matched_stmt else open_stmt
```

- Solution en modifiant le langage:

- $stmt \rightarrow \text{if expr then } stmt \textbf{fi} \mid$
 $\text{if expr then } stmt \textbf{else } stmt \textbf{fi} \mid$
 other
- $stmt \rightarrow \text{if expr then } \{ stmt \} \mid \text{if expr then } \{ stmt \} \textbf{else } \{ stmt \} \mid \text{other}$

- Solution avec des règles d'élimination:

- Toujours associer les **else** avec le **then** non associé le plus près.

Exercice 3

► Soit la grammaire de l'exercice 2:

► $S \rightarrow AB$ $A \rightarrow AA \mid Aa \mid aa$ $B \rightarrow BB \mid Bb \mid bb$

1. Montrez que la grammaire est ambiguë.
2. Éliminez l'ambiguïté en réécrivant la grammaire pour en obtenir une équivalente, mais sans ambiguïté.
3. Expliquez brièvement pourquoi votre grammaire est équivalente.

Associativité des opérateurs

- ▶ Une grammaire se doit de représenter l'associativité des opérateurs.
- ▶ **Exemple :**
 - ▶ $9+5+2$ est équivalent à $(9+5)+2$ parce que le $+$ est associatif à gauche.
 - ▶ $list \rightarrow list + digit \mid digit$
 - ▶ Dans un langage de programmation, $a=b=c$ est équivalent à $a=(b=c)$ parce que l'opérateur $=$ est associatif à droite.
 - ▶ $list \rightarrow id = list \mid id$
- ▶ Écrivons les grammaires correspondantes.

Précédence des opérateurs

- ▶ En arithmétique, les opérateurs $*$ et $/$ ont une précedence plus élevée que les opérateurs $+$ et $-$.
- ▶ **Exemple :**
 - ▶ La grammaire doit donc reconnaître que $9+5*2$ est équivalent à $9+(5*2)$.
- ▶ Quels sont les problèmes avec la grammaire suivante:
 - ▶ $E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$

Précédence des opérateurs

- ▶ La grammaire doit donc reconnaître que $9+5*2$ est équivalent à $9+(5*2)$.
- ▶ Quels sont les problèmes avec la grammaire suivante:
 - ▶ $E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$
 1. La grammaire est ambiguë.
 2. Plusieurs arbres de parsage sont possibles et ne garantissent pas la priorité des opérations.

Précédence des opérateurs

- Solution

- ▶ Quels sont les problèmes avec la grammaire suivante:
 - ▶ $E \rightarrow E + E \mid E * E \mid (E) \mid \text{id}$
- ▶ Une grammaire qui reconnaît la priorité des opérateurs et l'associativité à gauche:
 - ▶ $E \rightarrow E + T \mid T$
 - ▶ $T \rightarrow T * F \mid F$
 - ▶ $F \rightarrow (E) \mid \text{id}$
- ▶ Lire le livre du Dragon section 4.1.2

TP1 – Grammaire et parseur

- ▶ Écrire la grammaire pour générer un parseur pour le « langage des expressions ».
- ▶ JavaCC est un générateur de parseurs descendants-récurrents (Semaine 2).
 - ▶ *Java Compiler Compiler.*
- ▶ Les parseurs LL (descendants-récurrents) ne supportent pas la récursivité à gauche.

Réversivité à gauche

► Quels sont les problèmes avec la grammaire suivante:

► $E \rightarrow E + T \mid T$

► $T \rightarrow T * F \mid F$

► $F \rightarrow (E) \mid \text{id}$

1. La grammaire est réversible à gauche.

Réversivité à gauche

- Solution

- ▶ Quels sont les problèmes avec la grammaire suivante:

- ▶ $E \rightarrow E + T \mid T$

- ▶ $T \rightarrow T * F \mid F$

- ▶ $F \rightarrow (E) \mid \text{id}$

- ▶ Il faut éliminer la réversivité à gauche:

$$E \rightarrow T E'$$

$$E' \rightarrow + T E' \mid \varepsilon$$

$$T \rightarrow F T'$$

$$T' \rightarrow * F T' \mid \varepsilon$$

$$F \rightarrow (E) \mid \text{id}$$