

# Compilation Avanc  e (MII90) 2014

## Cours I : Introduction & Analyseurs

Carlos Agon  
agonc@ircam.fr



### C'est quoi la compilation ?

**Traduire** : un programme vers du code en langage machine

- 50's assembler  
du code machine textuel pour des langages de haut niveau (Fortran)
- 60's Les langages   voluent (la r  cursion) le langage machine suit (utilisation d'une pile),  
mais pas tant que   a.
- 80's Repr  sentation automatique des donn  es, le langage machine ne suit plus..

Il y a de plus en plus un   cart entre l'expressivit   des langages des haut niveau et  
celle du langage machine. C'est grave ?

... le compilateur doit s'occuper de cette fracture.

## Autres tâches

**Link** : liaison des fichiers des unités de compilation (module, class, interface, package, etc.)

**Runtime** : utilisation des bibliothèques d'exécution (i/o, gc, appels de méthodes, etc.)

**VM** : production de *bytecode* et interprétation/compilation (jit) vers du code machine

**Transformations source-à-source** : macros – interopérabilité (C)

**Optimisations** : temps, espace, énergie.

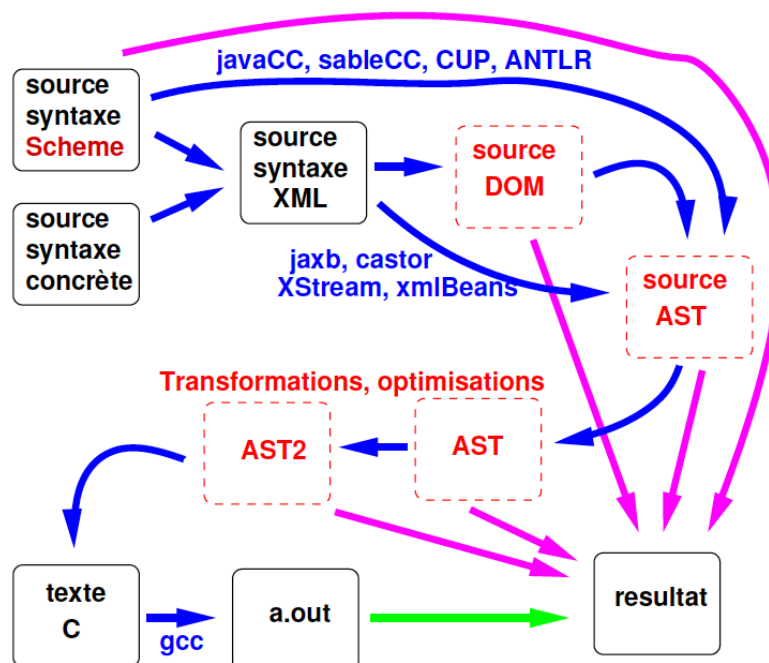
**Sûreté** : typage.

...

... la compilation est une chaîne de processus.

3

## Chaîne de compilation (Souvenir d'ILP)



4

## Chaîne de compilation classique

**Code source**

**Analyse lexicale**

**Analyse syntaxique**

**Analyses statiques**

**Génération**

5

## Informations

### **Première partie**

Cours et TD/TME par Carlos Agon

- Cours 1 : Introduction, analyseurs
- Cours 2-3 : Machines abstraites, virtuelles et bibliothèques d'exécution
- Cours 4 : Modèles mémoire - Garbage collection
- Cours 5 : Contrôle de haut niveau

### **Deuxième partie**

Cours et TD/TME par Karine Heydemann

- Mémoire cache et optimisation pour la hiérarchie mémoire
- Architecture matérielle basée sur le parallélisme d'instructions
- Ordonnancement et optimisation

### **Evaluation**

- Un examen 60 %
- Un projet par partie ( 20 % + 20 %)

6

## **Cours I**

**Code source**

**Analyse lexicale**

**Analyse syntaxique**

**Analyses statiques**

**Génération**

7

**Langages formels**

8

## Définitions

**Symbole** : signe primitif (lettre, chiffre, point, ligne, ...)

**Alphabet** : ensemble fini de symboles ( $\Sigma$ )

**Chaîne** : séquence finie de symboles (la chaîne vide  $\epsilon$ )

**Langages formel** : ensemble de chaînes définies sur un alphabet  $\Sigma$

### Problèmes :

- Comment définir un langage en intension ?
- Comment calculer si une chaîne  $w$  appartient ou non à  $L$  ?

9

## Hiérarchie des langages

Caractérisation des langages en fonction de la difficulté (calculabilité, espace, temps) de reconnaissance d'une chaîne.

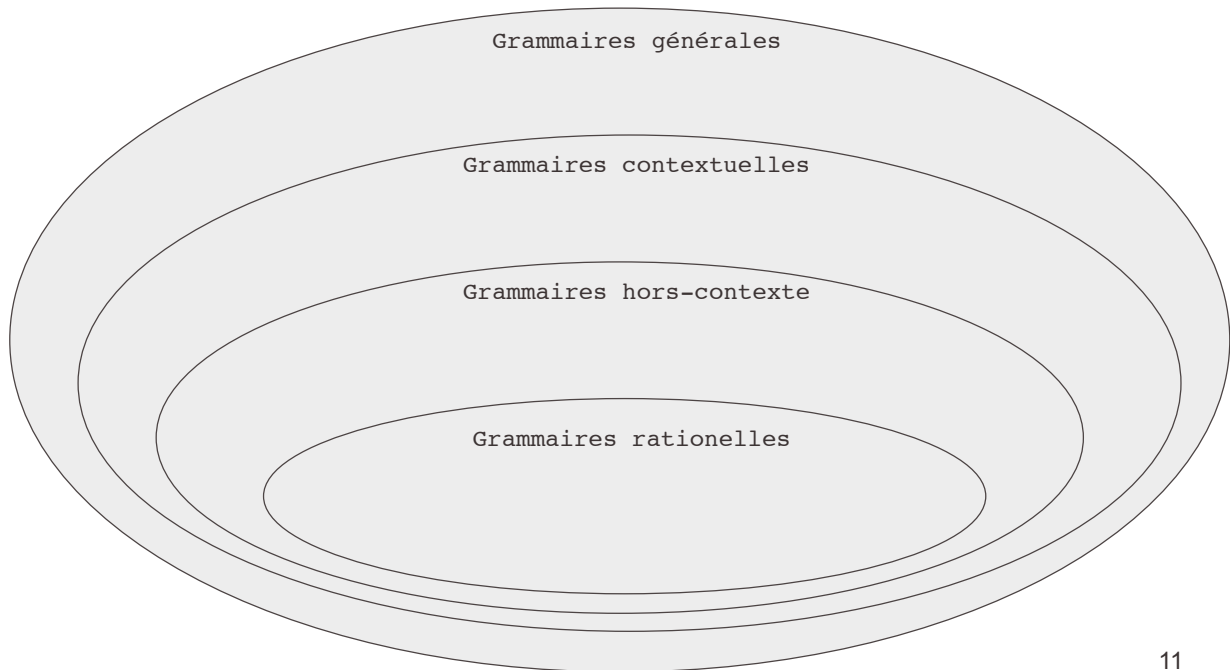
**Langages finis** : plutôt facile.

**Langages infinis**: hiérarchie de Chomsky :

- rationnels ou réguliers,
- indépendants du contexte,
- dépendants et
- récursivement énumérables.

10

## Hiérarchie de Chomsky



11

## Grammaires formelles

**T** : symboles terminaux

**N** : symboles non terminaux

**R** :  $\{\alpha \rightarrow \beta\}$  avec  $\alpha \in (\mathbf{T} \cup \mathbf{N})^*$  et  $\beta \in (\mathbf{T} \cup \mathbf{N})^*$

**S** : un symbole de départ e **N**

### Exemple :

- **T** = {a,b,c}
- **N** = {S}
- **S** = S
- **R** =  $S \rightarrow aSb$   
 $aSb \rightarrow aaSbb$   
 $aSb \rightarrow c$

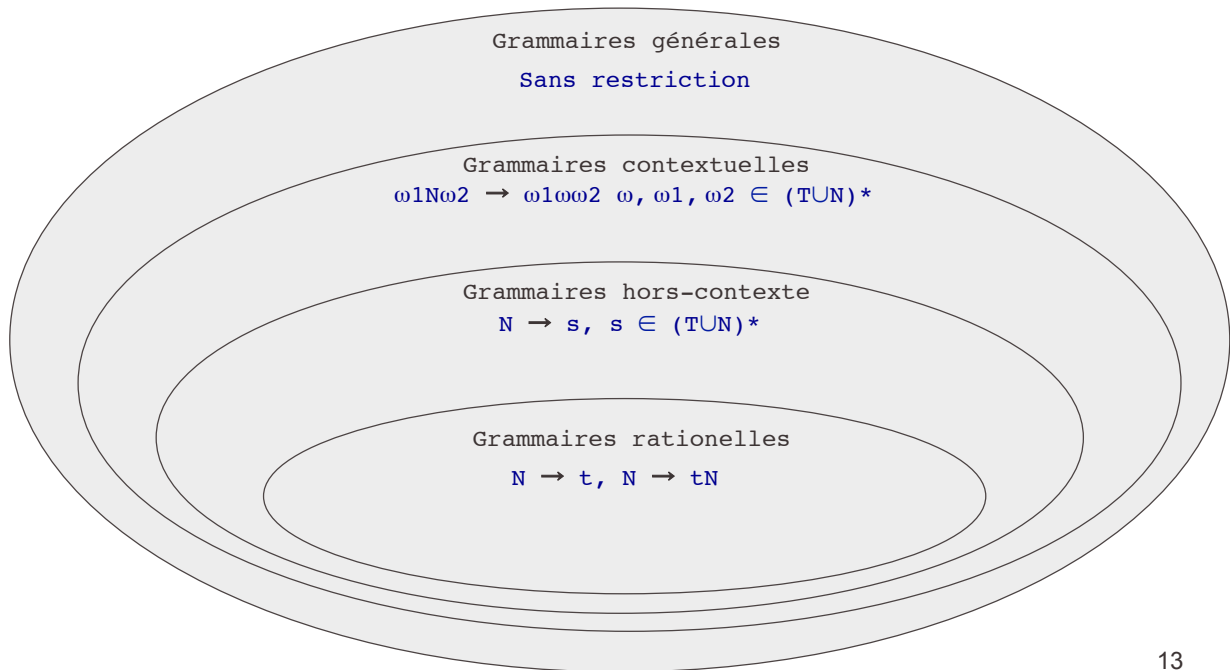
**Acceptées** : c, aacbb, acb,...

**Non-acceptées** : acab,...

Langage :  $a^n c b^n$

12

## Hiérarchie de Chomsky



13

## Langages réguliers

Soit  $\Sigma = \{a_1, \dots, a_n\}$

**Les langages triviaux** :  $L_i = \{a_i\}$

Et trois opérations :

**union** :  $L \cup L' = \{\omega \mid \omega \in L \text{ ou } \omega \in L'\}$

**produit** :  $L.L' = \{v.\omega \mid v \in L \text{ et } \omega \in L'\}$

**Itération** :  $L^* = \bigcup_{i \geq 0} L^i$

$L^0 = \{\varepsilon\}$  ;  $L^1 = L$  ;  $L^2 = L.L^1$  ; ... ;  $L^{i+1} = L.L^i$

**Langage régulier** :  $L \subseteq \Sigma^*$

ssi  $L$  peut être obtenu à partir des langages triviaux sur  $\Sigma$  uniquement par l'application des opérations d'union, de produit et d'itération.

14

## automates

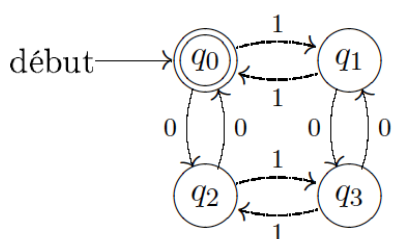
Comment calculer si une chaîne  $\omega$  appartient ou non à  $L$  ?

**Automate fini déterministe (AFD)** :  $(Q, \Sigma, \delta, q_0, F)$  tq.

- $Q = \{ q_0, \dots, q_k \}$
- $\Sigma = \text{alphabet } \{x_1, x_2, \dots, x_n\}$
- $\delta : Q \times \Sigma \rightarrow Q$
- $q_0 \in Q$  état initial
- $F \subseteq Q$  états finaux

15

## AFD exemple



États	Entrées	
	0	1
$q_0$	$q_2$	$q_1$
$q_1$	$q_3$	$q_0$
$q_2$	$q_0$	$q_3$
$q_3$	$q_1$	$q_2$

✓ 1010

✗ 2014

✗ 101

✗ 000

✓ 0000

**Langage accepté** : soit  $A$  un automate fini  
 $L(A) = \{ \omega \mid \omega \text{ est acceptée par } A \}$

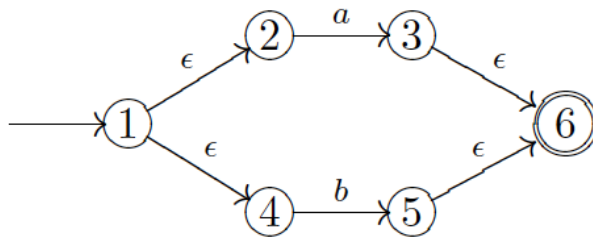
16



## AFN

Pareil que un AFD sauf

- $\Sigma = \text{alphabet } \{\epsilon, x_1, x_2, \dots, x_n\}$
- $\delta \subseteq Q \times \Sigma \times Q \times \Sigma$



17

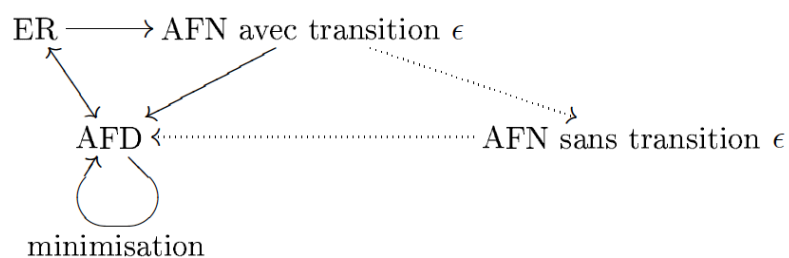
## Expressions régulières

Une autre façon de définir un langage en intension.

$(a|b)^*abb$

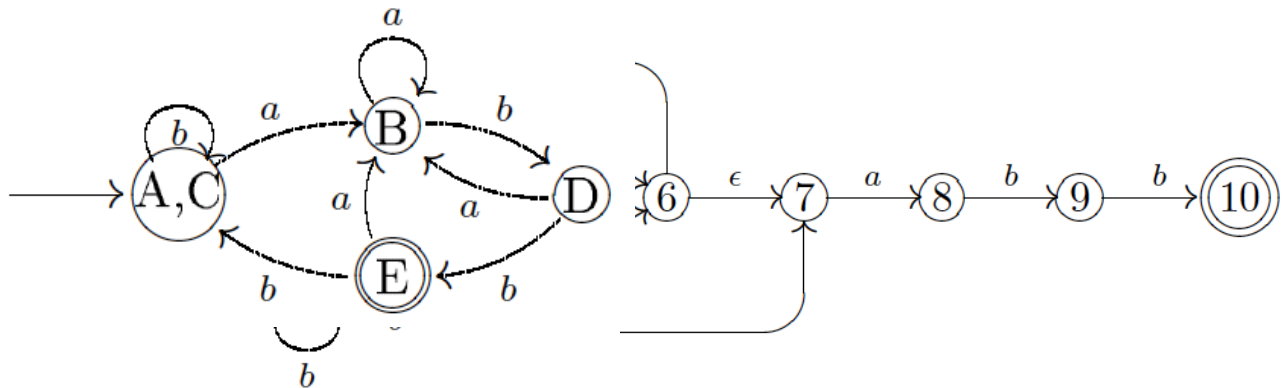
✓ aaaabb ; babb ; abb...

✗ Bb ; aaaba ; bbbbbb...



18

## Conversions



19

## Langages non-contextuels $N \rightarrow s, s \in (TUN)^*$

L'automate à pile est aux grammaires non-contextuelles ce que les automates finis sont aux grammaires régulières.

**Automate à pile (AP) :**  $(Q, \Sigma, \Gamma, \delta, q_0, Z_0, F)$  tq.

- $Q = \{ q_0, \dots, q_k \}$
- $\Sigma$  = alphabet fini d'entrée
- $\Gamma$  = alphabet fini de la pile
- $\delta : (Q \times \Sigma \cup \{\epsilon\}) \times \Gamma \rightarrow Q \times \Gamma^*$
- $q_0 \in Q$  état initial
- $Z_0$  symbole initial dans la pile
- $F \subseteq Q$  états finaux

20

## Automates à pile

### Comportement d'acceptation :

$\omega$  est accepté par l'automate AP si toute la chaîne  $\omega$  est consommée, on arrive dans un état final et la pile est vide.

### Comment construire l'automate à pile à partir d'une grammaire ?

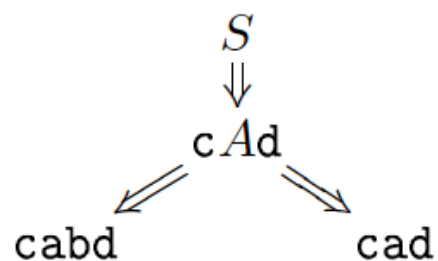
- Méthode universelle (Cocke, Younger et Kasami)  $O(n^3)$
- Analyse descendante LL(k)
- Analyse ascendante LR(k)

21

## Analyse descendante

$$S \rightarrow cAd$$

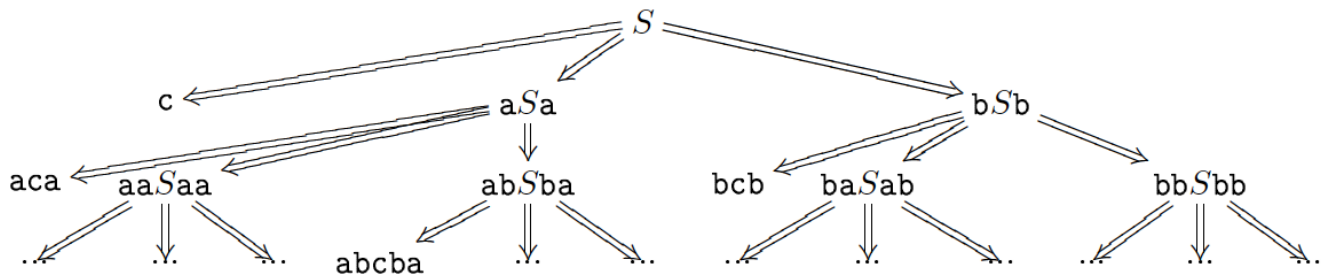
$$A \rightarrow ab \mid a$$



22

## Analyse descendante

$$S \rightarrow c \mid aSa \mid bSb$$



① **Descente récursive non prédictive**

② **Descente récursive prédictive LL(k)**

**k** est le nombre de lexème à tester pour prendre une décision

23

## LL(1)

- Un seul *token* suffit pour prendre la bonne décision
- Bon avec de grammaires petites
- Restriction :
  - Récursion gauche interdite
  - Nécessité de factorisation à gauche

24

## Factorisation à gauche

$S \longrightarrow \text{if } E \text{ then } S \mid \text{if } E \text{ then } S \text{ else } S$   
 $E \longrightarrow \dots$

Trouver le plus long préfixe  $\alpha$  commun aux deux alternatives  
Si  $\alpha \neq \epsilon$  remplacer  $A \rightarrow \alpha\beta_1 \mid \alpha\beta_2 \mid \dots \mid \alpha\beta_n \mid \omega$   
par :

$A \rightarrow \alpha A' \mid \omega$   
 $A' \rightarrow \beta_1 \mid \beta_2 \mid \dots \mid \beta_n$

$S \longrightarrow \text{if } E \text{ then } S S'$   
 $S' \longrightarrow \epsilon \mid \text{else } S$   
 $E \longrightarrow \dots$

25

## Elimination des récursivités à gauche

$E \longrightarrow E + E$

$E \longrightarrow E + E + E + E + E \dots$

$S \longrightarrow A\alpha$   
 $A \longrightarrow S$

$S \longrightarrow A\alpha \longrightarrow A\alpha\alpha \longrightarrow A\alpha\alpha\alpha \longrightarrow A\alpha\alpha\alpha\alpha \dots$

26

## Elimination des récursivités directes

$$A \longrightarrow A\alpha_1 \mid \dots \mid A\alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$

$$\begin{array}{l} A \longrightarrow \beta_1 A' \mid \dots \mid \beta_n A' \\ A' \longrightarrow \alpha_1 A' \mid \dots \mid \alpha_m A' \mid \epsilon \end{array}$$

$$\begin{array}{l} E \longrightarrow E + T \mid T \\ T \longrightarrow T * F \mid F \\ F \longrightarrow (E) \mid id \end{array}$$

$$\begin{array}{l} E \longrightarrow TE' \\ E' \longrightarrow + TE' \mid \epsilon \\ T \longrightarrow FT' \\ T' \longrightarrow * FT' \mid \epsilon \\ F \longrightarrow (E) \mid id \end{array}$$

27

## Elimination des récursivités

1. On ordonne les non-terminaux  $A_1, A_2, \dots, A_n$ .
2. **pour**  $i := 1$  jusqu'à  $n$  **faire**
3.     **pour**  $j := 1$  jusqu'à  $i - 1$  **faire**
4.         remplacer toutes les productions  $A_i \rightarrow A_j \gamma$  par les productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ , où  $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  sont toutes les  $A_j$ -productions courantes.
- fin**
5. éliminer les récursivités à gauche immédiates des  $A_i$ -productions à l'aide de la transformation suivante :
 
$$A_i \rightarrow A_i \alpha_1 \mid \dots \mid A_i \alpha_m \mid \beta_1 \mid \dots \mid \beta_n$$
 devient :
 
$$\begin{array}{l} A_i \rightarrow \beta_1 A'_i \mid \dots \mid \beta_n A'_i \\ A'_i \rightarrow \alpha_1 A'_i \mid \dots \mid \alpha_m A'_i \mid \epsilon \end{array}$$
- fin**

28

## Elimination des récursivités

$$\begin{aligned} S &\rightarrow Aa \mid b \\ A &\rightarrow bdA' \\ A' &\rightarrow cA' \mid adA' \mid \varepsilon \end{aligned}$$

29

## Analyse ascendante

LR(I) est l'approche utilisée presque partout

Plus souple car  $LL(I) \subset LR(I)$

$$\begin{aligned} \text{EXPR} &::= \textit{integer} & (R1) \\ &| \text{EXPR EXPR} + & (R2) \\ &| \text{EXPR EXPR} * & (R3) \end{aligned}$$

$$((1 + 2) * 3) + 4 \quad \equiv \quad 1 \ 2 + 3 * 4 + \$$$

30

## Analyse ascendante

ACTION	INPUT	STACK
	<u>1</u> 2+3*4+\$	[]
Shift	2+3*4+\$	[1]
Reduce (R1)	2+3*4+\$	[EXPR]
Shift	<u>+</u> 3*4+\$	[2EXPR]
Reduce (R1)	+3*4+\$	[EXPR EXPR]
Shift, Reduce (R2)	<u>3</u> *4+\$	[EXPR]
Shift, Reduce (R1)	*4+\$	[EXPR EXPR]
Shift, Reduce (R3)	<u>4</u> +\$	[EXPR]
Shift, Reduce (R1)	+\$	[EXPR EXPR]
Shift, Reduce (R2)	<u>\$</u>	[EXPR]

$EXPR ::= integer \quad (R1)$   
 $\quad \mid EXPR EXPR + \quad (R2)$   
 $\quad \mid EXPR EXPR * \quad (R3)$

31

## Analyse ascendante (ambigüités)

$E0 ::= integer \quad (R1)$   
 $\quad \mid E0 + E0 \quad (R2)$

ACTION	INPUT	STACK
⋮	<u>+</u> ...	[E0 + E0 ...]
⋮		

Conflit **shift/reduce** : on empile le **+** ou on réduit **E0+E0** ?

32



## Analyse ascendante (ambiguïtés)

$$\begin{array}{lcl} E1 & ::= & integer \quad (R1) \\ & | & E1 + E1 \quad (R2) \\ & | & E1 * E1 \quad (R3) \end{array}$$

$(integer + integer) * integer \neq integer + (integer * integer)$

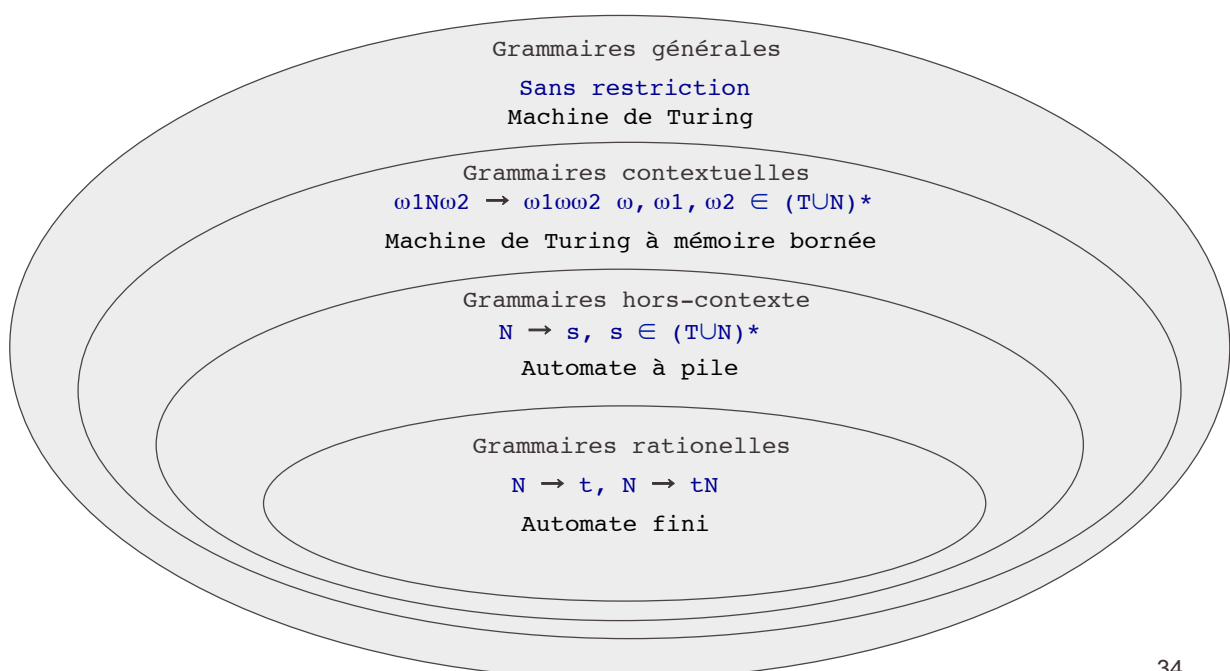
On donne des priorités dans la grammaire...

... mais aussi peut transformer la grammaire.

$$\begin{array}{lcl} E & ::= & E + T \quad (R1) \\ & | & T \quad (R2) \\ T & ::= & T * F \quad (R3) \\ & | & F \quad (R4) \\ F & ::= & integer \quad (R5) \end{array}$$

33

## Hiérarchie de Chomsky



34

## Pour les langages de programmation

Caractérisation des langages en fonction de la difficulté de reconnaissance d'une chaîne.

- **Forth, LISP**
- **ML, Pascal, C, etc**
- ...
- **C++ Java**

Plus une grammaire est difficile, plus

- Augmente la complexité (temps / espace)
- Lisibilité des automates
- Possibilité d'ambiguïté
- Messages d'erreur des parseurs lisibles.

- Lexique simple (textuel).
- Grammaire contextuelle dans la plus part des cas.

35

## Backus-Naur Form (BNF)

<terme> ::= <nombre signé> | <nombre signé> <op> <terme>

<nombre signé> ::= <nombre> | « + » <nombre> | « - » <nombre>

<nombre> ::= <entier> | <nombre fractionnaire>

<nombre fractionnaire> ::= <entier> | <entier> « / » <entier>

<entier> ::= <chiffre> | <chiffre><entier>

<chiffre> ::= « 0 » | « 1 » | ... | « 8 » | « 9 »

<op> ::= « \* » | « / » | ... | « + » | « - »

méta-symboles + non terminaux + terminaux

::= , | , ( , )

<nom>

« mots »

36

## Exemple d'un langage

	Lexique	$\{1, 0\}$
Syntaxe	$B ::= BC \mid C$ $C ::= 0 \mid 1$	
		$B_v : B \rightarrow \text{Nat}$ $B_v[[BC]] = (B_v[[B]] * 2) + C_v[[C]]$ $B_v[[C]] = C_v[[C]]$
Sémantique		$C_v : C \rightarrow \text{Nat}$ $C_v[[0]] = 0$ $C_v[[1]] = 1$

37

## Exemple d'un langage

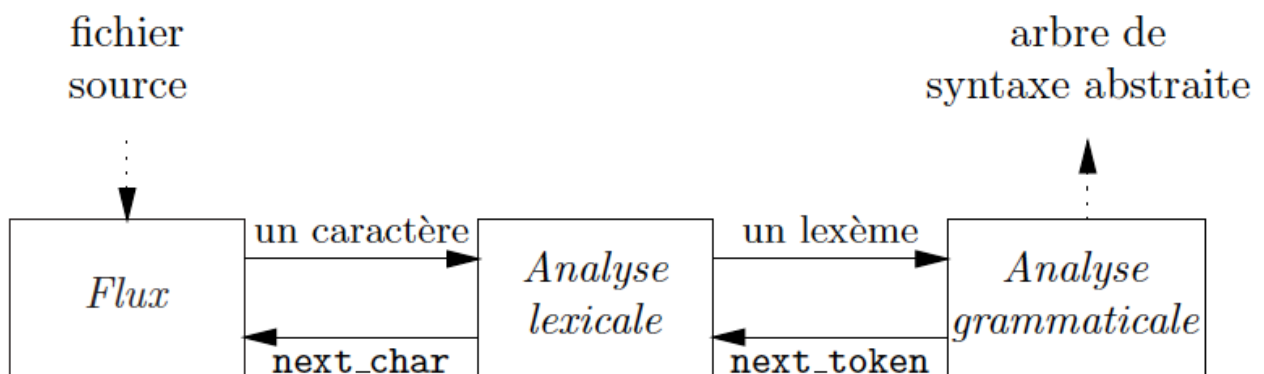
$$\begin{aligned} B_v[[101]] &= (B_v[[10]] * 2) + C_v[[1]] \\ &= ((B_v[[1]] * 2) * 2 + C_v[[0]]) + 1 \\ &= ((C_v[[1]] * 2) * 2 + 0) + 1 \\ &= (1 * 2) * 2 + 0 + 1 \\ &= 5 \end{aligned}$$

38

## Outils d'analyse

39

## Rapport entre les analyseurs



40

## Analyse lexicale (Lex)

Transformer une suite de caractères en une suite de mots (lexèmes / *tokens*). Peut se faire au même temps que l'analyse grammaticale, mais...

Le but de l'analyseur lexical est donc de 'consommer' des symboles et de les fournir à l'analyseur syntaxique.

Un fichier de description pour Lex est formé de trois parties (optionnelles)

déclarations

%%

productions

%%

code additionnel

41

## Lex (déclarations)

- ① Code dans le langage cible `%{ %}`
- ② Expressions régulières *non\_terminal* *expression\_regulière*

```
%{  
  
#include "calc.h"  
  
#include <stdio.h>  
#include <stdlib.h>  
  
%}
```

```
blancs    [\t\n ]+  
lettre    [A-Za-z]  
chiffre   [0-9]  
identificateur {lettre}(_|{lettre}|{chiffre})*  
entier    {chiffre}+
```

42

## Lex (productions)

- ① Code dans le langage cible `%{ %}` (placé au début de `yylex()` )
- ② Productions *expression\_régulière* *action*

```
%{  
printf ("hello");  
%}
```

```
[a-z]          printf ("%c", ((yytext[0]-'a'+13)%26)+'a');
```

```
[\\r\\n]|(\\r\\n)  {printf ("%s", yytext) ; fflush (stdout);  
                  printf ("hola");  
                  printf ("%c", yytext [0]);}
```

43

## Lex (code additionnel)

```
/* Code par default*/
```

```
main() {  
    yylex();  
}
```

### Compil :

```
Flex myfile.lex && gcc lex.yy.c -lfl -o exe
```

44

## Analyse grammaticale Yacc (Yet Another Compiler Compiler)

Produire le texte source d'un analyseur syntaxique du langage engendré par une grammaire du type LR(1)

Il est aussi possible, en plus de la vérification de la syntaxe de la grammaire, de lui faire effectuer des actions sémantiques.

Pareil que pour Lex, un fichier Yacc se compose de trois parties :

```
déclarations
%%
productions
%%
code additionnel
```

45

## Yacc (déclarations)

- ① Code dans le langage cible `%{ %}`
- ② Déclaration des terminaux pouvant être rencontrés, grâce au mot-clé `%token`
- ③ Le type de donnée du terminal courant, avec le mot-clé `%union`.
- ④ Des informations donnant la priorité et l'associativité des opérateurs.
- ⑤ L'axiome de la grammaire, avec le mot-clé `%start`.

46

## Yacc (productions)

- ① Déclarations et/ou définitions encadrées par `%{` et `%}`
- ② Productions de la grammaire

```
notion_non_terminale:  
corps_1 { action_semantique_1 }  
| corps_2 { action_semantique_2 }  
| ...  
| corps_n { action_semantique_n }  
;
```

Les `corps_i` peuvent être des notions terminales ou non terminales du langage.

47

## Yacc (code additionnel)

Cette partie, qui comporte le code additionnel, devra obligatoirement comporter une déclaration du `main()` (qui devra appeler la fonction `yyparse()`), et de la fonction `yyerror(char *message)`, appelée lorsqu'une erreur de syntaxe est trouvée.

48



## Yacc (exemple)

```
%{
#include "global.h"
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
%}

%token NOMBRE
%token PLUS MOINS FOIS DIVISE PUISSANCE
%token PARENTHESE_GAUCHE PARENTHESE_DROITE
%token FIN

%left PLUS MOINS
%left FOIS DIVISE
%left NEG
%right PUISSANCE

%start Input
%%
Input:
    /*Vide */
    | Input Ligne
    ;

Ligne:
    FIN
    | Expression FIN { printf("Resultat : %f\n",$1); } ;
```

```
Expression:
    NOMBRE { $$=$1; }
    | Expression PLUS Expression { $$=$1+$3; }
    | Expression MOINS Expression { $$=$1-$3; }
    | Expression FOIS Expression { $$=$1*$3; }
    | Expression DIVISE Expression { $$=$1/$3; }
    | MOINS Expression %prec NEG { $$=-$2; }
    | Expression PUISSANCE Expression { $$=pow($1,$3); }
    | PARENTHESE_GAUCHE Expression
    PARENTHESE_DROITE { $$=$2; }
    ;

%%

int yyerror(char *s) {
    printf("%s\n",s);
}

int main(void) {
    yyparse();
}
```

49

## Yacc + Lex

```
%{
#include "global.h"
#include "calc.h"
#include <stdlib.h>
%}

blancs [ \t]+
chiffre [0-9]
entier {chiffre}+
exposant [eE][+-]?{entier}
reel {entier}("."{entier})?{exposant}?

%%
{blancs} { /* On ignore */ }
{reel} {
    yylval=atof(yytext);
    return(NOMBRE);
}
```

```
"+" return(PLUS);
"-" return(MOINS);
"*" return(FOIS);
"/" return(DIVISE);

"^" return(PUISSANCE);

"(" return(PARENTHESE_GAUCHE);
")" return(PARENTHESE_DROITE);

"\n" return(FIN);
```

50

## Compilation (Demo)

```
>bison -d calc.y
>mv calc.tab.h calc.h
>mv calc.tab.c calc.y.c
>flex calc.lex
>mv lex.yy.c calc.lex.c
>gcc -c calc.lex.c -o calc.lex.o
>gcc -c calc.y.c -o calc.y.o
>gcc -o calc calc.lex.o calc.y.o -ll
```

L'appel à bison (le Yacc de GNU) calc.tab.h pour définir les terminaux. On appelle flex (le Lex de GNU) et on donne des noms un peu plus convenables pour le tout. Il suffit alors de compiler, sans oublier la librairie spéciale Lex ("-ll").

```
>calc
1+2*3
Resultat : 7.000000
2.5*(3.2-4.1^2)
Resultat : -34.025000
```

51

**ocamllex et ocaml yacc**

52

## Ocamllex format du fichier

```
{  
Code Ocaml exécuté avant  
}
```

Let <nom> = <expreg>

```
Rule <nom> = parse  
<expreg>  
| <expreg>      { <code action> }  
...  
| <expreg>      { <code action> }
```

```
{  
Code Ocaml exécuté après  
}
```

53

## Ocamllex exemple

```
{open Printf open Char }  
  
rule rot13 = parse  
| ['a'-'z' ] as c  
    { printf "%c" (chr (((code c - code 'a' + 13) mod 26) + code 'a')) }  
| ['n' 'r'] | "\r\n" as s  
    { printf "%s%!" s }  
| eof    { raise Exit }  
| _ as c { printf "%c" c }  
  
{  
try let chan = Lexing.from_channel stdin in  
  while true do  
    rot13 chan  
  done  
with Exit -> ()  
}
```

**Compil** : ocamllex file.mll && Ocamlpt file.ml -o exe

54

## Ocamllex avec un parseur

```
{  
  open Parser  
  (* Parser doit définir type token = INT of int | OP of string | END *)  
}
```

```
rule expr= parse  
| [ ' ' '\t' ]      { expr lexbuf }  
| eof              { END }  
| [ '0'-'9' ]+ as s { INT (int_of_string s) }  
| [ '+' '-' '*' '/' ] as s { OP s }
```

55

## Ocamllex Demo

56

## Références

- « Automates et langages » Support de cours de J. Malenfant
- Slides Compilation Avancée (MI190) Benjamin Canu
- « Compilers Principles, Techniques and Tools » Aho, Sethi et Ulman
- « Développement d'applications avec Ocaml » E.CHAILLOUX, P. MANOURY, B.PAGANO
- « Tutorial de Lex/Yacc » Etienne Bernard