

# AGP: Algorithmique et programmation

tanguy.risset@insa-lyon.fr

Lab CITI, INSA de Lyon

Version du July 22, 2016

Tanguy Risset

July 22, 2016

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Tanguy Risset

AGP: Algorithmique et programmation

1

Expressions régulières Analyse lexicale Grammaire Analyse syntaxique Yacc (bison) Annexe: quelques précisions

## Langages, Grammaires, et Compilateurs

- Objectif de cette partie du cours: comprendre le mécanisme de compilation, le lien avec la notion de grammaire et savoir créer un parseur.
- A quoi est dû le succès des ordinateurs?
  - Progrès technologique pour l'intégration de transistors
  - Augmentation de la productivité des programmeurs
    - Langage de haut niveau
    - Compilateurs rapides, codes portables
- La notion de compilation n'est pas limitée aux langages de programmation.
- L'informaticien produit en permanence des programmes qui font passer les données d'une représentation à une autre.

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Tanguy Risset

AGP: Algorithmique et programmation

2

# Plan

- 1 Expressions régulières
- 2 Analyse lexicale
- 3 Grammaire
- 4 Analyse syntaxique
- 5 Yacc (bison)
- 6 Annexe: quelques précisions

Sources:

- Cours Compilation T. Risset,
- “Engineering a compiler, Cooper & Torczon”
- <http://ds9a.n1/lex-yacc/>

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

Tanguy Risset

AGP: Algorithmique et programmation

3

Expressions régulières Analyse lexicale Grammaire Analyse syntaxique Yacc (bison) Annexe: quelques précisions

# Table of Contents

- 1 Expressions régulières
- 2 Analyse lexicale
- 3 Grammaire
- 4 Analyse syntaxique
- 5 Yacc (bison)
- 6 Annexe: quelques précisions

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

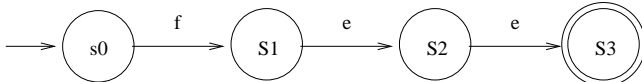
Tanguy Risset

AGP: Algorithmique et programmation

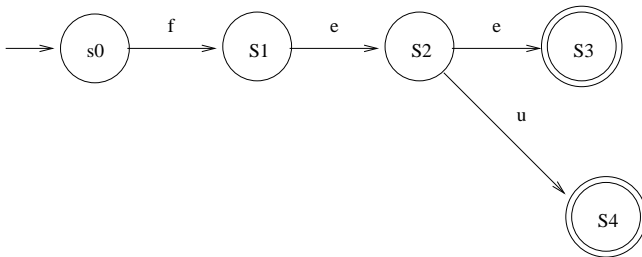
4

# Langage régulier et automate

- Formellement, un *langage* (au sens *syntaxe*) est simplement un ensemble de mots formés à partir d'un alphabet fini.
- Exemples de langage
  - l'ensemble des mots du dictionnaire (alphabet de a à z)
  - l'ensemble des mots constitués de deux 'a' suivis d'un nombre quelconque de 'b':  $\mathcal{L} = \{aa, aab, aabb, aabbb, \dots\}$
- Un langage est dit *régulier* s'il est reconnu par un automate.
- Automate reconnaissant le mot fee:



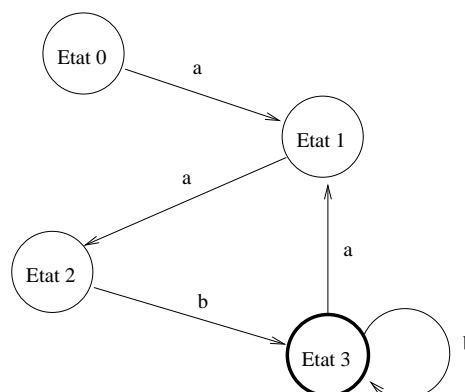
- Automate reconnaissant le langage {fee, feu}



Navigation: ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

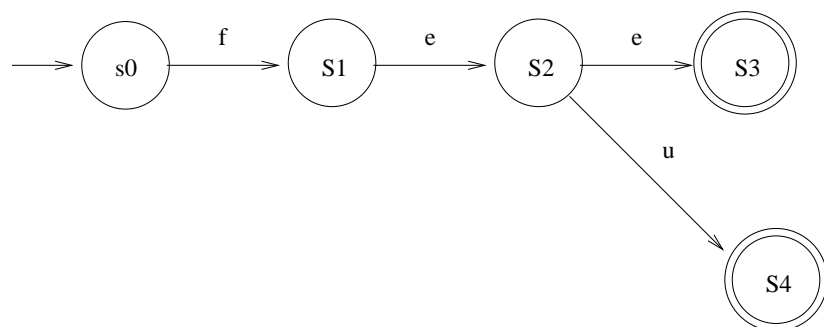
## Notion d'automate

- Un automate est une collection de K états numérotés de 0 à K-1, ainsi qu'une collection de transitions
- Un état particulier est l'état initial.
- Tous les états sont soit des états d'acceptation et soit des états de refus
- Les transitions, sont des triplets (état 1, lettre x, état 2) qui signifient : lorsque je suis dans l'état 1 et que je lis la lettre x, alors je vais dans l'état 2.



Navigation: ◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

# Notion de mot reconnu



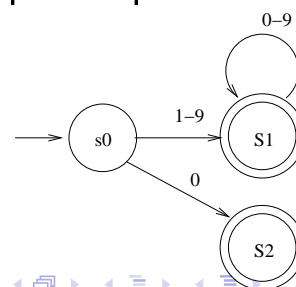
- $fee \rightarrow$  reconnu
- $feu \rightarrow$  reconnu
- $fei \rightarrow$  non reconnu (impossible de lire 'i')
- $fe \rightarrow$  non reconnu (arrêt dans un état non final)

## Automate: définition formelle

- Un automate fini déterministe est donné par  $(S, \Sigma, \delta, s_0, S_F)$  ou:
  - $S$  est un ensemble fini d'états;
  - $\Sigma$  est un alphabet fini;
  - $\delta : S \times \Sigma \rightarrow S$  est la fonction de transition;
  - $s_0$  est l'état initial;
  - $S_F$  est l'ensemble des états finaux
- pour l'automate reconnaissant "fee", on a  $S = \{s_0, s_1, s_2, s_3\}$ ,  
 $\Sigma = \{f, e\}$  (ou tout l'alphabet)  
 $\delta = \{\delta(s_0, f) \rightarrow s_1, \delta(s_1, e) \rightarrow s_2, \delta(s_2, e) \rightarrow s_3\}$ .

- Il y a en fait un état implicite d'erreur ( $s_e$ ) vers lequel vont toutes les transitions qui ne sont pas définies.
- Un automate *accepte* une chaîne de caractère si et seulement si en démarrant de  $s_0$ , il s'arrête dans un état final

automate qui reconnaît  
n'importe quel entier:



# Expressions régulières

- Les expressions régulières sont couramment employées en ligne de commande, par exemple: `ls *.c`
- Une *expression régulière*  $X$  basée sur un alphabet décrit un *langage*, c'est à dire un ensemble de mots sur cet alphabet, on note ce langage  $E(X)$ .
- Exemples:

Expression régulière	langage reconnu
$a^*.b$	mots constitués d'un nombre quelconque de $a$ suivi d'un $b$
$a.(a+b)^*.a + b.(a+b)^*.b$	mots constitués des lettres $a$ et $b$ , commençant et finissant par la même lettre
$a.(a+b)^*.a+a$	mots commençant par $a$ et finissant par $a$ (alphabet $\{a,b\}$ )

Navigation icons: back, forward, search, etc.

## Expressions régulières: Définition formelle

- Individuellement, chaque lettre est une expression régulière:  
 $E(a) = \{a\}$
- Si  $X_1$  et  $X_2$  sont deux expressions régulières, alors  $X_1.X_2$  est une expression régulière (concaténations des mots des deux langages).
- Si  $X_1$  et  $X_2$  sont deux expressions régulières, alors  $X_1+X_2$  est une expression régulière (union des mots des deux langages:  
 $E(X_1+X_2) = E(X_1) \cup E(X_2)$ )
- Si  $X$  est une expression régulière alors  $X^*$  est une expression régulière qui décrit l'ensemble des mots construits en répétant autant de fois que l'on veut (éventuellement zéro fois) un mot décrit par  $X$
- $*$  plus prioritaire que  $.$  qui est plus prioritaire que  $+$

Navigation icons: back, forward, search, etc.

# Expression régulière et analyse lexicale

- Les expressions régulières définissent une classe de langage simple (langages réguliers) reconnue par des automates finis.
- On va utiliser cela en compilation pour identifier les *tokens* (ou *lexèmes*) du langage, c'est à dire les unités lexicales élémentaires:
  - les mots-clés (do, while, for, etc.)
  - les identificateurs (noms de variable, de fonction, etc.)
  - les constantes (entier, flottant, chaîne de caractère)
- C'est la première phase de la compilation: ***l'analyse lexicale***
- Il existe depuis longtemps des outils (historiquement: `lex`) qui, à partir des expressions régulières décrivant les lexèmes, génèrent automatiquement le programme qui reconnaît les lexèmes.
- C'est la notion de compilateur de compilateur.

## Table of Contents

1 Expressions régulières

2 Analyse lexicale

3 Grammaire

4 Analyse syntaxique

5 Yacc (bison)

6































## Exemple d'analyse syntaxique ascendante (LR(1))

Action	Règle utilisée	Etat
		$\uparrow x - 2 \times y$
shift	Identifier	Identifier $\uparrow - 2 \times y$
reduce	$Factor \rightarrow Identifier$	Factor $\uparrow - 2 \times y$
reduce	$Term \rightarrow Factor$	Term $\uparrow - 2 \times y$
reduce	$Expr \rightarrow Term$	Expr $\uparrow - 2 \times y$
shift	—	Expr $- \uparrow 2 \times y$
shift	Number	Expr $- Number \uparrow \times y$
reduce	$Factor \rightarrow Number$	Expr $- Factor \uparrow \times y$
reduce	$Term \rightarrow Factor$	Expr $- Term \uparrow \times y$
shift	$\times$	Expr $- Term \times \uparrow y$
shift	Identifier	Expr $- Term \times Identifier \uparrow (EOF)$
reduce	$Factor \rightarrow Identifier$	Expr $- Term \times Factor \uparrow (EOF)$
reduce	$Term \rightarrow Term \times Factor$	Expr $- Term \uparrow (EOF)$
reduce	$Expr \rightarrow Expr + Term$	Expr $\uparrow (EOF)$
success		

## Comment ça marche

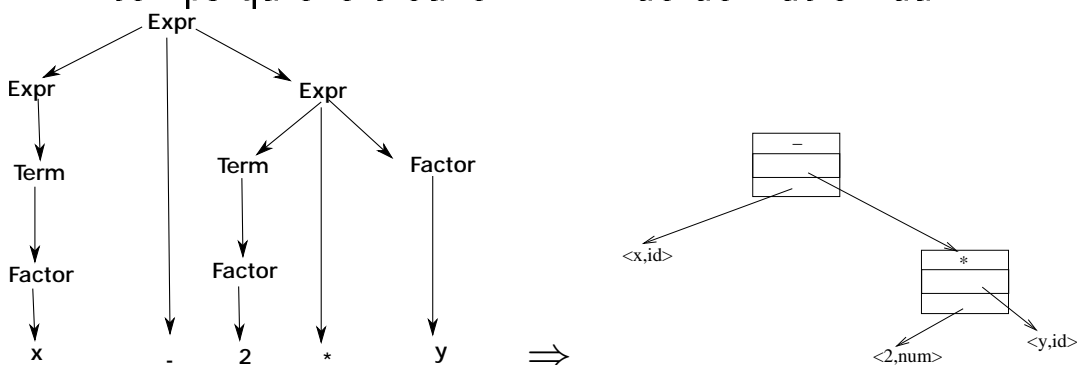
- Nous ne détaillons pas ici la théorie de l'analyse syntaxique.
- Cette théorie prouve que l'on peut choisir la bonne règle à chaque fois, en regardant simplement le prochain caractère à lire.
- Ce qu'il faut savoir pour utiliser yacc:
  - La technique utilisée est dite *LR(1)* (ou *left-to-right, reduce as needed*). Ces parseurs lisent de gauche à droite et construisent (à l'envers) une dérivation la plus à droite en regardant au plus un symbole sur l'entrée, d'où leur nom: *Left-to-right can, Reverse-ightmost derivation with 1 symbol lookahead*.
  - La grammaire ne doit pas être ambiguë
  - Lors de chaque réduction, yacc réalise des *action* qui permettent de construire, en *avant*, la représentation intermédiaire du programme.

## Exemple d'action possible

<i>S/R</i>	<i>Règle utilisée</i>	<i>Action</i>	<i>Etatdel' analyse</i>
<i>shift</i>	<i>Identifier</i>	Creer nœud Id	$\uparrow x - 2 \times y$ <i>Identifier</i> $\uparrow - 2 \times$
<i>reduce</i>	<i>Factor</i> $\rightarrow$ <i>Identifier</i>		<i>Factor</i> $\uparrow - 2 \times y$
<i>reduce</i>	<i>Term</i> $\rightarrow$ <i>Factor</i>		<i>Term</i> $\uparrow - 2 \times y$
<i>reduce</i>	<i>Expr</i> $\rightarrow$ <i>Term</i>		<i>Expr</i> $\uparrow - 2 \times y$
<i>shift</i>	—		<i>Expr</i> $- \uparrow 2 \times y$
<i>shift</i>	<i>Number</i>	Creer nœud Num	<i>Expr</i> $-$ <i>Number</i>
<i>reduce</i>	<i>Factor</i> $\rightarrow$ <i>Number</i>		<i>Expr</i> $-$ <i>Factor</i> $\uparrow$
<i>reduce</i>	<i>Term</i> $\rightarrow$ <i>Factor</i>		<i>Expr</i> $-$ <i>Term</i> $\uparrow >$
<i>shift</i>	$\times$		<i>Expr</i> $-$ <i>Term</i> $\times \uparrow$
<i>shift</i>	<i>Identifier</i>	Creer nœud Id	<i>Expr</i> $-$ <i>Term</i> $\times$ <i>I</i>
<i>reduce</i>	<i>Factor</i> $\rightarrow$ <i>Identifier</i>		<i>Expr</i> $-$ <i>Term</i> $\times$ <i>I</i>
<i>reduce</i>	<i>Term</i> $\rightarrow$ <i>Term</i> $\times$ <i>Factor</i>	Creer nœud $\times$	<i>Expr</i> $-$ <i>Term</i> $\uparrow ($
<i>reduce</i>	<i>Expr</i> $\rightarrow$ <i>Expr</i> $+$ <i>Term</i>	Creer nœud $-$	<i>Expr</i> $\uparrow (EOF)$
<i>success</i>			

## Résultat du l'analyse syntaxique

L'analyse syntaxique construit un AST (Abstract Syntax Tree), en même temps qu'elle trouve l'arbre de dérivation du mot.



# Table of Contents

- 1 Expressions régulières
- 2 Analyse lexicale
- 3 Grammaire
- 4 Analyse syntaxique
- 5 Yacc (bison)**
- 6 Annexe: quelques précisions

## Yacc (bison)

- `flex` est l'implémentation de `lex` par Vern Paxson et `bison` est la version GNU de `yacc`, dans la suite on parlera toujours de `lex` et `yacc` pour désigner `flex` et `bison`
- `yacc` signifie *Yet Another Compiler Compiler*
- `yacc` peut parser des flots de *tokens*, il doit donc être utilisé avec un front-end qui transforme un flot de caractères en un flot de tokens (par exemple `lex`)

## Premier exemple simple: un petit additionneur

- On doit lire une suite d'addition et afficher le résultat.
  - entrée du parser: 3+4+6
  - affichage du parseur: résultat: 13
- Proposition de grammaire élémentaire:
  - Grammaire (S,T,N,P),
  - $S = \{ \text{ OM } \}$ ,
  - $N = \{ \text{ EXPR, OM } \}$
  - $T = \{ \text{ ENTIER, '+' } \}$
  - $$P = \left\{ \begin{array}{l} \text{ SOM} ::= \text{ EXPR} \\ \text{ EXPR} ::= \text{ NUMBER} \mid \\ \qquad \qquad \text{ EXPR '+' NUMBER} \end{array} \right\}$$

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

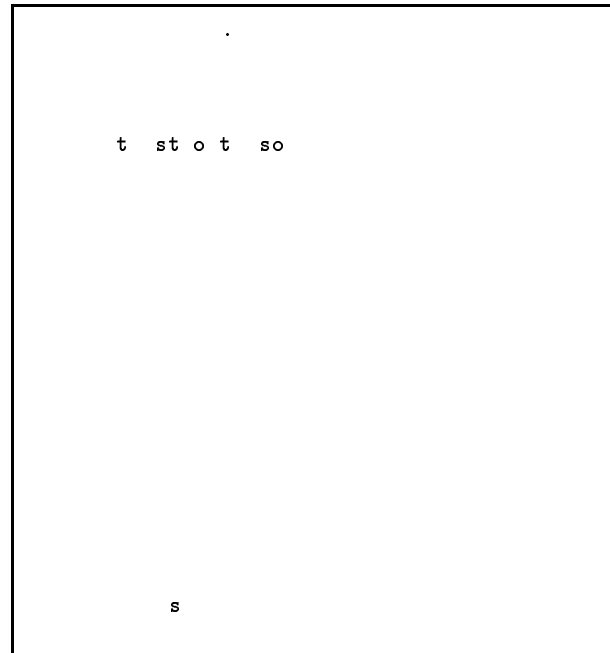
## Additionneur simple: fichier adder.l

```
%{
#include <stdio.h>
#include "adder.tab.h"
}%
%%
[0-9]+      yyl al=atoi(yytext)  return NUMBER
+          return PLU
n          /* ignore linebreak */
[ t]+      /* ignore whitespace */
%%
```

- adder.tab.h sera généré par la commande yacc, à inclure dans le fichier lex.yy.c
- yyval et yytext sont des variables partagées par lex et yacc: quand lex reconnaît un lexème, il le met dans la variable yytext et le transmet à yacc. Quand c'est une valeur, yacc s'attend à le trouver dans yyval

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ 🔍 ↻

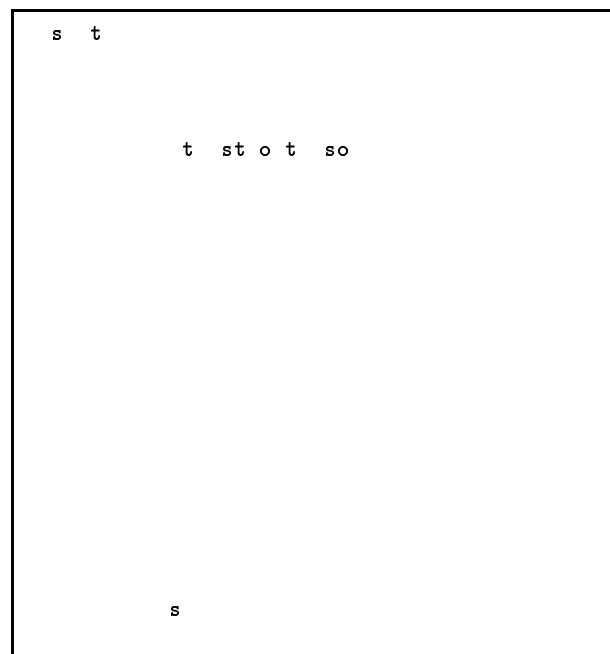
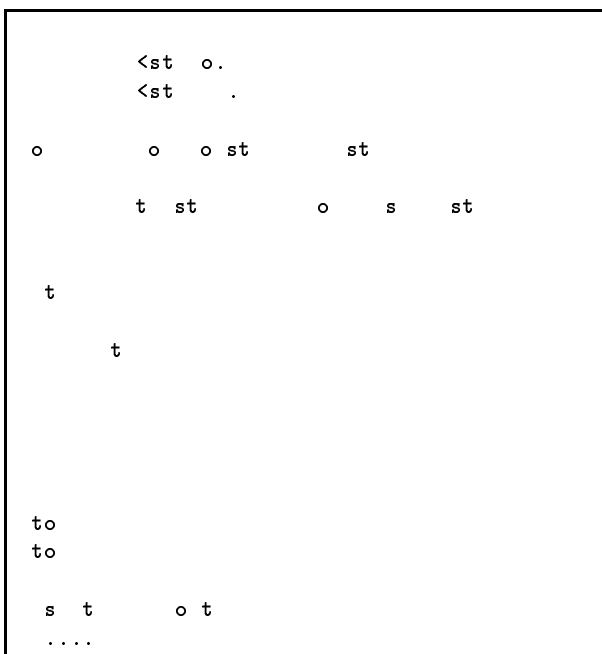
## Additionneur simple: grammaire $\Rightarrow$ adder.y



- A chaque réduction de règles est associée une action
- chaque symbole renvoie un objet
- \$1, \$2,... correspondent aux objets renvoyés par la partie droite de la règles. \$\$ correspond à l'objet renvoyé par la partie gauche

Navigation icons: back, forward, search, etc.

## Additionneur simple: fichier adder.y complet



- A chaque réduction de règles est associé une action
- chaque symbole renvoie un objets
- \$1, \$2,... correspondent aux objets renvoyés par la partie droite de la règles. \$\$ correspond à l'objet renvoyé par la partie gauche

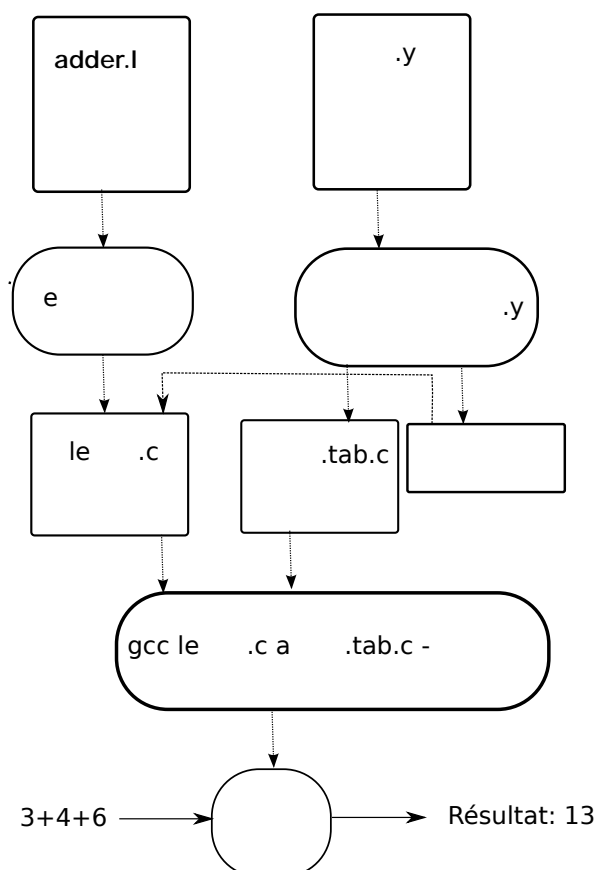
Navigation icons: back, forward, search, etc.

## Structure d'un fichier yacc/bison

```
... t o ...
.... s...
... o ...
```

- Définitions:
  - Code C copié au début du fichier .c généré (entre %{ et %})
  - On doit notamment définir yyerror et yywrap
  - définition des tokens (type, associativité éventuellement)
- Règles:
  - définition des règles de grammaire avec les actions associées aux réductions
- Code:
  - Le code que vous rajoutez, au minimum un `ain` qui appelle `yyparse` comme ici

## Adder: Processus de compilation global





## Compilation et utilisation du programme adder

compilation:

```
trisset@fania:~/$ ake
bison -d adder.y
gcc -c -o adder.tab.o adder.tab.c
flex adder.l
gcc -c -o lex.yy.o lex.yy.c
gcc adder.tab.o lex.yy.o -o adder
trisset@fania:~/$
```

exécution:

```
trisset@fania:~/$ ./adder
2+3
^D
so e:
trisset@fania:~/$ ./adder
1+2
+

+9
^D
so e: 19
```

Navigation icons: back, forward, search, etc.

## Le Langage TC-LOGO

- Le langage TC-LOGO (version 1.0) est hérité du formalisme LOGO ([http://en.wikipedia.org/wiki/Logo\\_programming\\_language](http://en.wikipedia.org/wiki/Logo_programming_language)) destiné à dessiner un graphique.
- Il sert à décrire le chemin que va suivre le crayon pour dessiner ce graphique.
- Il comporte 4 instructions:
  - FORWARD  $n$
  - LEFT  $T$
  - RIGHT  $T$
  - REPEAT [ *programme TC-LOGO* ]
- Il n'y a pas de point-virgule en TC-LOGO, les séparateurs valides sont l'espace, le retour chariot ou la tabulation

# Le Langage TC-LOGO

- Exemple de programme TC-LOGO:

```
FORWARD 100 FORWARD 10
REPEAT 10
  [FORWARD 10 FORWARD 100]
```

```
REPEAT 360
[FORWARD 1
LEFT
1]
```

- Grammaire (S,T,N,P),
- $S = \{ \text{PROG} \}$ ,  $N = \{ \langle \text{PROG} \rangle, \langle \text{INST} \rangle \}$
- $T = \{ \text{FORWARD}, \text{LEFT}, \text{RIGHT}, \text{REPEAT}, '[', ']', \text{ENTIER}, \text{SEPARATEUR} \}$
- ENTIER et SEPARATEUR sont des lexèmes qui seront identifiés par l'analyse lexicale

## Grammaire possible pour TC-LOGO

$P = \{$

$\langle \text{PROG} \rangle ::= \langle \text{INST} \rangle \mid \langle \text{PROG} \rangle \langle \text{INST} \rangle$

$\langle \text{INST} \rangle ::= \text{FORWARD ENTIER} \mid$

$\text{LEFT ENTIER} \mid$

$\text{RIGHT ENTIER} \mid$

$\text{REPEAT ENTIER } '[' \langle \text{PROG} \rangle ']'$

$\}$

# Lex et Yacc en résumé

- Outils pour produire des parseurs
- Utiles pour traiter les fichiers de données ou pour analyser des formats simples.
- Outils open source (gnu) extrêmement solides et portables (produisent du C travaillant sur les E/S standard).

## Table of Contents

- 1 Expressions régulières
- 2 Analyse lexicale
- 3 Grammaire
- 4 Analyse syntaxique
- 5 Yacc (bison)
- 6 Annexe: quelques précisions

## Exemple de grammaire ambiguë

- Une grammaire est ambiguë lorsque deux arbres de dérivations différents peuvent donner le même mot du langage.
- Exemple: grammaire ambiguë pour *if-then-else*

1		<i>Instruction</i>	→	<u>if</u>	<i>Expr</i>	<u>then</u>	<i>Instruction</i>	<u>else</u>	<i>Instruction</i>
2				<u>if</u>	<i>Expr</i>	<u>then</u>	<i>Instruction</i>		
3							<i>Assiguation</i>		
4							<i>AutreInstructions....</i>		

- Avec cette grammaire, le code

*if Expr<sub>1</sub> then if Expr<sub>2</sub> then Ass<sub>1</sub> else Ass<sub>2</sub>*

- peut être compris de deux manières:

*if Expr<sub>1</sub>*  
     *then if Expr<sub>2</sub>*  
         *then Ass<sub>1</sub>*  
         *else Ass<sub>2</sub>*

*if Expr<sub>1</sub>*  
     *then if Expr<sub>2</sub>*  
         *then Ass<sub>1</sub>*  
     *else Ass<sub>2</sub>*

## Solution: desambiguitisation (desambiguation ?)

- Il faut changer la grammaire

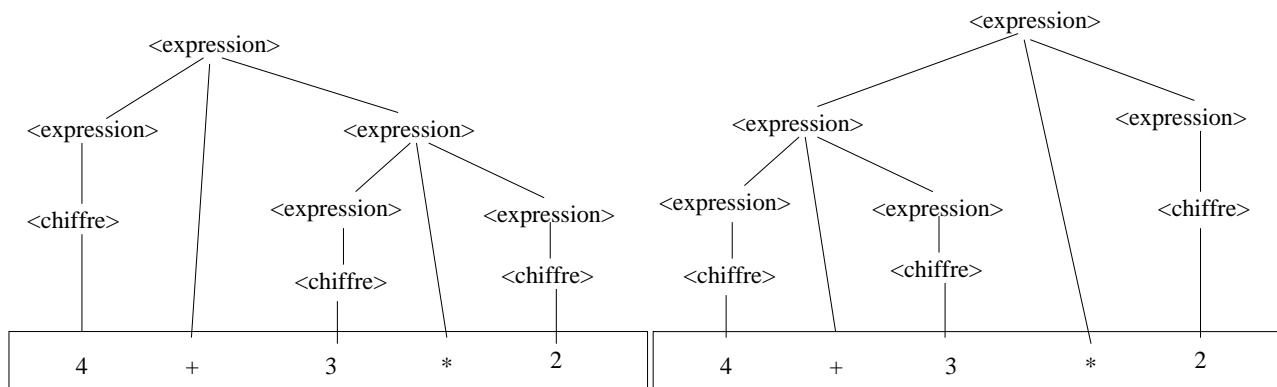
1		<i>Instruction</i>	→	<i>AvecElse</i>
2				<i>DernierElse</i>
3		<i>AvecElse</i>	→	<u>if</u> <i>Expr</i> <u>then</u> <i>AvecElse</i> <u>else</u> <i>AvecElse</i>
4				<i>Assiguation</i>
5				<i>AutreInstructions....</i>
6		<i>DernierElse</i>	→	<u>if</u> <i>Expr</i> <u>then</u> <i>Instruction</i>
7				<u>if</u> <i>Expr</i> <u>then</u> <i>AvecElse</i> <u>else</u> <i>DernierElse</i>
8				<i>AutreInstructions....</i>

## Retour sur la grammaire des expression ambiguë

- Soit la grammaire suivante:

```

<expression> ::= <expression> + <expression> |
                <expression> * <expression> |
                <number>
<number>      ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  
```



- Deux dérivations sont équivalentes si elles ont le même arbre de dérivation (seul l'ordre dans lequel on a choisit les règles peut changer).

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ ↻

## Solution: modifier la grammaire

- ... Sans modifier le langage:

```

<expression> ::= <expression> + <expression> |
                <terme>
<terme>      ::= <terme> * <terme> | <number>
<number>     ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  
```

- $\Rightarrow$  un seul arbre de dérivation possible.
- On peut faire d'autres améliorations pour l'associativité des opérateurs et les parenthèses. Une grammaire couramment utilisée pour les expressions arithmétiques est la suivante:

```

<expression> ::= <expression> + <terme> | <terme>
<terme>      ::= <terme> * <facteur> | <facteur>
<facteur>    ::= ( <expression> ) | <number>
<number>     ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
  
```

◀ ◻ ▶ ◀ ◻ ▶ ◀ ≡ ▶ ◀ ≡ ▶ ≡ ↺ ↻

