

Compiler

Parsing

Jyun-Ao Lin

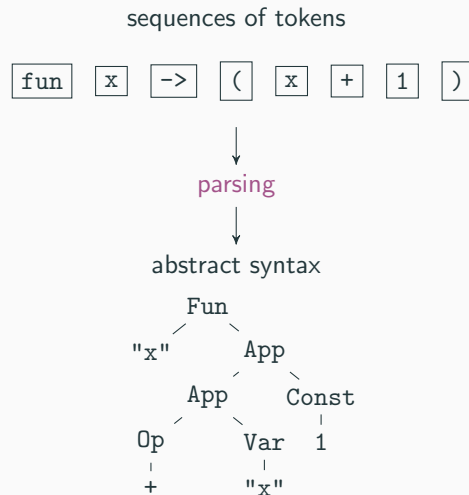
iFIRST & CSIE, NTUT
jalin@ntut.edu.tw

A large part of this course is based on the Compilation Course of J.-C. Filliâtre at ENS Ulm.

The goal of parsing is to recognize sentences that belong to the syntax of the language

Its input is the flow of tokens constructed by lexical analysis, its output is an abstract syntax tree

Lexical analysis: example



Syntax analysis must detect syntax errors and

- signal them with a position in the source
- explain them (most often limited to “syntax error” but also “unclosed parenthesis”, etc.)
- possibly resume the analysis to discover further errors

Which tool?

To implement syntax analysis, we are using

- a **context-free grammar** to define the syntax
- a **pushdown** automaton to recognize it

Similar to regular expressions / finite automata used in lexical analysis

Overview of the course

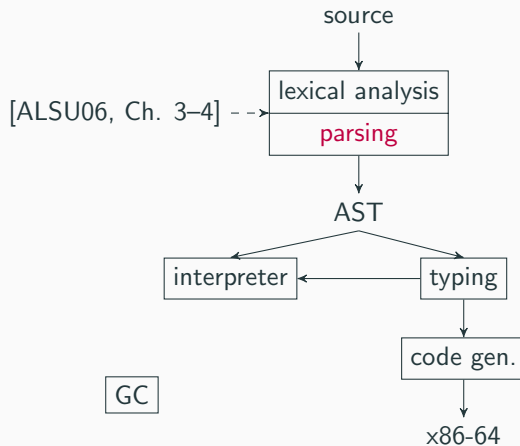


Table of contents

1. GRAMMARS
2. BOTTOM-UP PARSING
3. THE TOOL Menhir
4. CONSTRUCTION OF THE AUTOMATION AND THE TABLE
5. LOCALIZATION
6. ELEMENTARY PARSING

Grammars

Context-free grammar

Defn. (context-free grammar)

A context-free grammar is a tuple (N, T, S, R) where

- N is a finite set of non-terminal symbols
- T is a finite set of terminal symbols
- $S \in N$ is the start symbol (the axiom)
- $R \subseteq N \times (N \cup T)^*$ is a finite set of production rules



Example: arithmetic expressions

$N = \{E\}$, $T = \{+, *, (,), \text{int}\}$, $S = E$,
and $R = \{(E, E+E), (E, E*E), (E, (E)), (E, \text{int})\}$

In practice, we write production rules as follows:

$$\begin{array}{lcl} E & \rightarrow & E + E \\ & | & E * E \\ & | & (E) \\ & | & \text{int} \end{array}$$

The terminals are the tokens produced by the lexical analysis

Here `int` stands for an integer literal token (i.e. its nature, not its value)

Derivation

Defn. (derivation)

A word $u \in (N \cup T)^*$ **derives** to a word $v \in (N \cup T)^*$, denoted by $u \rightarrow v$, if there exists a decomposition

$$u = u_1 X u_2$$

with $X \in N, X \rightarrow \beta \in R$ and

$$v = u_1 \beta u_2$$



Example

$$\underbrace{E * (}_{u_1} \underbrace{E}_{X} \underbrace{)}_{u_2} \rightarrow E * (\underbrace{E + E}_{\beta})$$

Derivation

A sequence $w_1 \rightarrow w_2 \rightarrow \cdots \rightarrow w_n$ is called a derivation.

It is called **left derivation** (resp. **right**) if the reduced non-terminal is systemically the leftmost, i.e. $u_1 \in T^*$ (resp. the rightmost $u_2 \in T^*$)

We denote by \rightarrow^* the reflexive, transitive closure of \rightarrow .

Example

Example

Left derivation:

$$\begin{aligned}
 E &\rightarrow E * E \\
 &\rightarrow \text{int} * E \\
 &\rightarrow \text{int} * (E) \\
 &\rightarrow \text{int} * (E + E) \\
 &\rightarrow \text{int} * (\text{int} + E) \\
 &\rightarrow \text{int} * (\text{int} + \text{int})
 \end{aligned}$$

Then we have (in particular, but not uniquely)

$$E \xrightarrow{*} \text{int} * (\text{int} + \text{int})$$

Language

Defn. (language of grammar)

The language defined by a context-free grammar $G = (N, T, S, R)$ is the set of words of T^* derived from the axiom, i.e.

$$L(G) = \{w \in T^* \mid S \xrightarrow{*} w\}$$



Example

In the above example, we have

$$\text{int} * (\text{int} + \text{int}) \in L(G)$$

Derivation tree

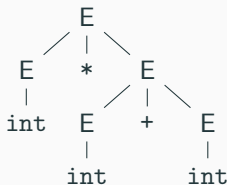
Defn. (derivation tree)

A **derivation tree** is a tree whose nodes are labeled with grammar symbols, such that

- the root is the axiom S
- any internal node X is a non-terminal whose subnodes are labeled by $\beta \in (N \cup T)^*$ with $X \rightarrow \beta$ a production rule
- leaves are terminal symbols



Example



Careful: this is **different** from the abstract syntax tree

Derivation tree

For a derivation tree in whose leaves form the word w in infix order, it is clear that we have $S \xrightarrow{*} w$

Conversely, to any derivation $S \xrightarrow{*} w$, we can associate a derivation tree whose leaves form the word w in infix order

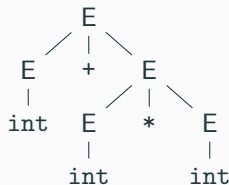
Idea: the derivation tree captures a whole set of derivations that we wish to identify

Example

The left derivation

$$E \rightarrow E + E \rightarrow \text{int} + E \rightarrow \text{int} + E * E \rightarrow \text{int} + \text{int} * E \rightarrow \text{int} + \text{int} * \text{int}$$

gives the derivation tree



but the right derivation

$$E \rightarrow E + E \rightarrow E + E * E \rightarrow E + E * \text{int} \rightarrow E + \text{int} * \text{int} \rightarrow \text{int} + \text{int} * \text{int}$$

too

Ambiguous

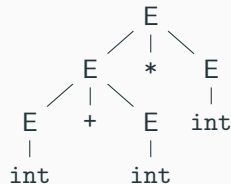
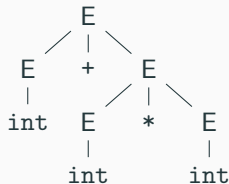
Defn. (ambiguous)

A grammar is called **ambiguous** if at least a word admits several derivation trees.



Example

The word `int + int * int` admits two derivation trees



and thus our grammar is ambiguous

Non-ambiguous grammar

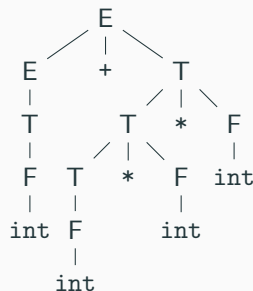
It is possible to propose another grammar, that is not ambiguous and that defines the same language

$$\begin{array}{lcl}
 E & \rightarrow & E + T \\
 & | & T \\
 T & \rightarrow & T * F \\
 & | & F \\
 F & \rightarrow & (E) \\
 & | & \text{int}
 \end{array}$$

This new grammar reflects the priority of multiplication over addition, and the choice of a left associativity for these two operations

Non-ambiguous grammar

Now, the word `int + int * int * int` has a single derivation tree,



corresponding to the left derivation

$$\begin{aligned}
 E &\rightarrow E + T \rightarrow T + T \rightarrow F + T \rightarrow \text{int} + T \rightarrow \text{int} + T * F \\
 &\rightarrow \text{int} + T * F * F \rightarrow \text{int} + F * F * F \\
 &\rightarrow \text{int} + \text{int} * F * F \rightarrow \text{int} + \text{int} * \text{int} * F \\
 &\rightarrow \text{int} + \text{int} * \text{int} * \text{int}
 \end{aligned}$$

Negative result

Whether a context-free grammar is ambiguous is **not decidable**

(reminder: decidable means that we can write a program that, for any input, terminates and outputs yes or no)

Approach

We are going to use **decidable sufficient criteria** to ensure that a grammar is not ambiguous, and for which we know how to decide membership efficiently (using a pushdown automaton)

The corresponding grammar classes are called LR(0), SLR(1), LALR(1), LR(1), LL(1), etc.

Bottom-up parsing

Main idea

- scan the input from left to right, and
- look for right-hand sides of production rules to build the derivation tree from bottom to top (*bottom-up parsing*)

Principle

The parser uses a stack that is a word of $(T \cup N)^*$

At each step, two actions can be performed

- a **shift** operation: we read a terminal from the input and we push it on the stack
- a **reduce** operation: the top of the stack is the right-hand side β of a production $X \rightarrow \beta$, and we replace β with X on the stack

Initially, the stack is empty

When no more action can be performed, the input is recognized if it was read entirely and if the stack is limited to the axiom S

Example

$$\begin{array}{lcl}
 E & \rightarrow & E + T \\
 & | & T \\
 T & \rightarrow & T * F \\
 & | & F \\
 F & \rightarrow & (E) \\
 & | & \text{int}
 \end{array}$$

| stack | input | action |
|------------------|-------------|-----------------------------------|
| ϵ | int+int*int | shift |
| int | +int*int | reduce $F \rightarrow \text{int}$ |
| F | +int*int | reduce $T \rightarrow F$ |
| T | +int*int | reduce $E \rightarrow T$ |
| E | +int*int | shift |
| $E+$ | int*int | shift |
| $E+\text{int}$ | *int | reduce $F \rightarrow \text{int}$ |
| $E+F$ | *int | reduce $T \rightarrow F$ |
| $E+T$ | *int | shift |
| $E+T*$ | int | shift |
| $E+T*\text{int}$ | | reduce $F \rightarrow \text{int}$ |
| $E+T*F$ | | reduce $T \rightarrow T*F$ |
| $E+T$ | | reduce $E \rightarrow E+T$ |
| E | | success |

LR parser (Knuth, 1965)

How to choose between shift and reduce?

Using an automaton and considering the first k tokens of the input;
this is called LR(k) analysis (LR means “**L**eft to right scanning, **R**ightmost derivation”)

In practice $k = 1$

i.e. we only consider the first token to take the decision

LR parsing

The stack looks like

$$s_0 x_1 s_1 x_2 \dots x_n s_n$$

where s_i is the state of the automation and $x_i \in T \cup N$ as before.

let a be the first token from the input; we look in the **action** table for state s_n and character a

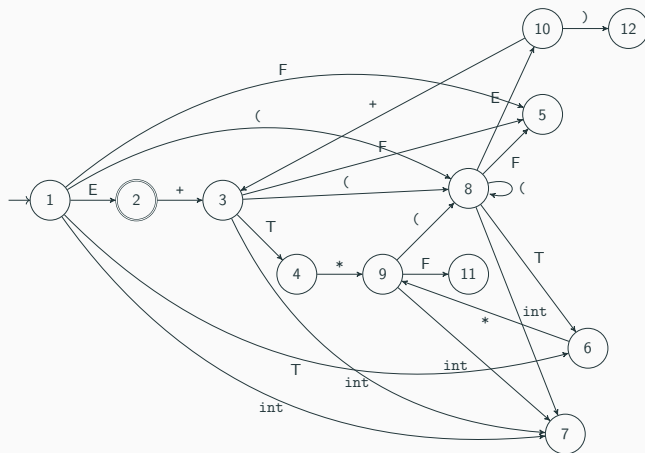
- if success or failure, we stop
- if shift, we push a and then the target state of the transition $s_n \xrightarrow{a} s$ into the stack
- if reduce rule $X \rightarrow \alpha$, with α of length p , then we have α on top of the stack

$$s_0 x_1 s_1 \dots x_{n-p} s_{n-p} \mid \alpha_1 s_{n-p+1} \dots \alpha_p s_n$$

we pop it and we push X s , where s is the target state of the transition $s_{n-p} \xrightarrow{X} s$, i.e.

$$s_0 x_1 s_1 \dots x_{n-p} s_{n-p} X s$$

Example

$$\begin{array}{lcl} E & \rightarrow & E + T \\ & | & T \\ T & \rightarrow & T * F \\ & | & F \\ F & \rightarrow & (E) \\ & | & \text{int} \end{array}$$


LR table

In practice, we do not work with automation but with two tables

- an **action** table with states as rows and terminals as columns; the box $\text{action}(s, a)$ indicates
 - shift s' for a shift and a new state s'
 - reduce $X \rightarrow \alpha$ for a reduction
 - a success
 - a failure
- a **transition** table having states for rows and non-terminals for columns; the box $\text{goto}(s, X)$ indicates the state resulting from a reduction of X

End of input

We also add a special token, denoted by #, which designates the end of the input

One can view it as adding it as a new non-terminal S (which becomes an axiom) and a new rule

$$\begin{array}{lcl} S & \rightarrow & E\# \\ E & \rightarrow & \dots \\ & \vdots & \end{array}$$

Example

```

1  E  →  E + T
2      |  T
3  T  →  T * F
4      |  F
5  F  →  ( E )
6      |  int

```

si = shift *i*

rj = reduce *j*

| | action | | | | | | goto | | |
|-------|--------|-----|----|----|-----|---------|------|---|----|
| state | (|) | + | * | int | # | E | T | F |
| 1 | s8 | | | | s7 | | 2 | 6 | 5 |
| 2 | | | s3 | | | success | | | |
| 3 | s8 | | | | s7 | | | 4 | 5 |
| 4 | | r1 | r1 | s9 | | r1 | | | |
| 5 | | r4 | r4 | r4 | | r4 | | | |
| 6 | | r2 | r2 | s9 | | r2 | | | |
| 7 | | r6 | r6 | r6 | | r6 | | | |
| 8 | s8 | | | | s7 | | 10 | 6 | 5 |
| 9 | s8 | | | | s7 | | | | 11 |
| 10 | | s12 | s3 | | | | | | |
| 11 | | r3 | r3 | r3 | | r3 | | | |
| 12 | | r5 | r5 | r5 | | r5 | | | |

Example of execution

| stack | input | action |
|------------------------------|--------------|---------------------------------|
| 1 | int+int*int# | s7 |
| 1 int 7 | +int*int# | $F \rightarrow \text{int}, g5$ |
| 1 F 5 | +int*int# | $T \rightarrow F, g6$ |
| 1 T 6 | +int*int# | $E \rightarrow T, g2$ |
| 1 E 2 | +int*int# | s3 |
| 1 E 2 + 3 | int*int# | s7 |
| 1 E 2 + 3 int 7 | *int# | $F \rightarrow \text{int}, g5$ |
| 1 E 2 + 3 F 5 | *int# | $T \rightarrow F, g4$ |
| 1 E 2 + 3 T 4 | *int# | s9 |
| 1 E 2 + 3 T 4 * 9 | int# | s7 |
| 1 E 2 + 3 T 4 * 9 int 7 | # | $F \rightarrow \text{int}, g11$ |
| 1 E 2 + 3 T 4 * 9 F 11 | # | $T \rightarrow T * F, g4$ |
| 1 E 2 + 3 T 4 | # | $E \rightarrow E + F, g2$ |
| 1 E 2 | # | success |

Automation

Bottom-up parsing is powerful but computing the tables is complex

We have tools to automate the process

This is the big family of yacc, bison, ocaml yacc, cup, menhir, . . .
(YACC means *Yet Another Compiler Compiler*)

The tool Menhir

Menhir

is a tool that transforms a grammar into an OCaml parser

Menhir is based on an LR(1) parsing

Each production rule of the grammar is accompanied by a semantic action

i.e. an OCaml code building a semantic value

(typically an abstract syntax tree)

Menhir is used in conjunction with a lexical analyzer

(typically `ocamllex`)

Structure

A Menhir file has the suffix `.mly` and has the following structure

```
%{  
    ... arbitrary OCaml code ...  
%}  
... declaration of tokens ...  
... declaration of precedents and associativity ...  
... declaration of entry points ...  
  
%%  
non-terminal-1:  
| production { action }  
| production { action }  
;  
non-terminal-2:  
| production { action }  
...  
%%  
    ... arbitrary OCaml code ...
```

Minimal example

```
%token PLUS LPAR RPAR EOF
```

```
%token <int> INT
```

```
%start <int> phrase
```

```
%%
```

```
phrase:
```

```
| e = expression; EOF { e }
```

```
;
```

```
expression:
```

```
| e1 = expression; PLUS; e2 = expression { e1 + e2 }
```

```
| LPAR; e = expression; RPAR { e }
```

```
| i = INT { i }
```

```
;
```

Example

We compile the `arith.mly` file in the following way

```
% menhir -v arith.mly
```

We obtain pure OCaml code in `arith.ml(i)`, which contains in particular

- the declaration of a token type

```
type token = RPAR | PLUS | LPAR | INT of int | EOF
```

- for each non-terminal declared with `%start`, a function of the type

```
val phrase: (Lexing.lexbuf -> token) -> Lexing.lexbuf -> int
```

as we can see, this function takes as an argument a lexical analyzer, of the type produced by `ocamllex`

ocamllex + menhir

When we combine ocamllex and menhir

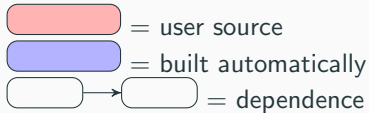
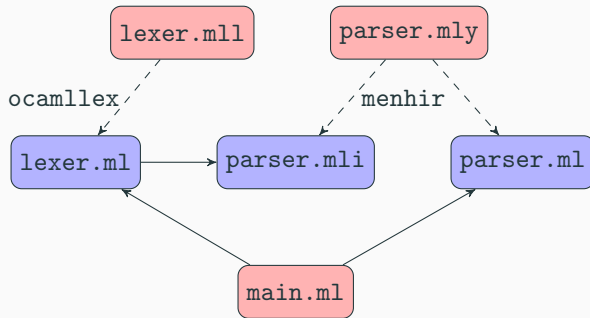
- `lexer.mll` refers to the tokens defined in `parser.mly`

```
{  
    open Parser  
}  
...
```

- The lexical analyzer and the parser are combined as follows:

```
let c = open_in file in  
let lb = Lexing.from_channel c in  
let e = Parser.phrase Lexer.token lb in  
...
```

ocamllex + menhir



Conflicts

When the grammar is not LR(1), Menhir shows the **conflicts** to the user

- the file `.automaton` contains the LR(1) automaton (more later), with conflicts listed
- the file `.conflicts` contains an explanation for each conflict, as a sequence of tokens leading to two distinct derivation trees

Example

On the above grammar, Menhir reports a conflict

```
% menhir -v arith.mly
Warning: one state has shift/reduce conflicts.
Warning: one shift/reduce conflict was arbitrarily resolved.
```

The file `arith.automaton` contains in particular

```
State 6:
expression -> expression . PLUS expression [ RPAR PLUS EOF ]
expression -> expression PLUS expression . [ RPAR PLUS EOF ]
-- On PLUS shift to state 5
-- On RPAR reduce production expression -> expression PLUS expressi
-- On PLUS reduce production expression -> expression PLUS expressi
-- On EOF reduce production expression -> expression PLUS expressio
** Conflict on PLUS
```

Example

The file `arith.conflicts` contains a clear explanation

```

** Conflict (shift/reduce) in state 6.
** Token involved: PLUS
** This state is reached from phrase after reading:

expression PLUS expression

** In state 6, looking ahead at PLUS, shifting is permitted
** because of the following sub-derivation:

expression PLUS expression
      expression . PLUS expression

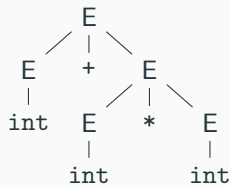
** In state 6, looking ahead at PLUS, reducing production
** expression -> expression PLUS expression
** is permitted because of the following sub-derivation:

expression PLUS expression // lookahead token appears

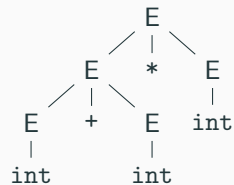
```

Example

In other word, the question is to choose between



and



Conflict resolution

One way to resolve conflicts is to tell Menhir how to choose between shift and reduction

For this, we can give **priorities** to tokens and productions, and rules of **associativity**

By default, the priority of a production is that of its rightmost token (but it can be specified explicitly)

Conflict resolution

If the priority of the production is higher than that of the tokens to be read,
then reduction is favored

Conversely, if the priority of the tokens is higher,
then reading is favored

In case of equality, associativity is consulted: a left-associative tokens favors reduction and a right-associative tokens favors shift

Example

In our example, it is enough to indicate for example that PLUS is left-associative

```
%token PLUS LPAR RPAR EOF
%token <int> INT
%left PLUS
%start <int> phrase
%%
phrase:
| e = expression; EOF           { e }
;
expression:
| e1 = expression; PLUS; e2 = expression { e1 + e2 }
| LPAR; e = expression; RPAR         { e }
| i = INT                             { i }
;
```

Priorities

To associate priorities with tokens, we use the following convention:

- the order of declaration of the associativity sets the priorities (the first tokens have the lowest priorities)
- Several tokens can appear on the same line, thus having the same priority

Example:

```
%left PLUS MINUS  
%left TIMES DIV
```

A great classic

The following grammar contains a conflict

```
expression:
| IF e1 = expression; THEN; e2 = expression
{ ... }
| IF e1 = expression; THEN; e2 = expression;
ELSE; e3 = expression
{ ... }
| i = INT
{ ... }
| ...
```

(known in English as dangling else)

Explanation, resolution

It corresponds to the situation

```
IF a THEN IF b THEN c ELSE d
```

To associate the ELSE with the closest THEN, it is necessary to favor shift

```
%nonassoc THEN
```

```
%nonassoc ELSE
```

Menhir assets

offers many advantages over traditional tools such as `ocaml yacc` :

- non-terminals parameterized by (non-)terminals
- in particular, easy to write regular expressions ($E?$, E^* , E^+) and lists with separator
- explanation of conflicts
- interactive mode
- LR(1) parsing, where most tools only offer LALR(1)

read the Menhir manual! (accessible from the course page)

Localization

For the following phases of the analysis (typically typing) can **localize** error messages, it is advisable to keep a localization information in the abstract syntax tree

Menhir provides this information in `$startpos` and `$endpos`, two value of the type `Lexing.position`; this information was transmitted to it by the lexical analyzer

Be careful: `ocamllex` automatically maintains only the absolute position in the file; to have the line and column numbers up to date, it is necessary to call `Lexing.new_line` for each carriage return

Localization

One way to keep location information in the abstract syntax tree is as follows

```
type expression =  
  { desc: desc;  
    loc : Lexing.position * Lexing.position }  
  
and desc =  
  | Econst of int  
  | Eplus of expression * expression  
  | Eneg of expression  
  | ...
```

Each node is thus decorated by a localization

Localization

The grammar can therefore look like this

expression:

```
| d = desc { { desc = d; loc = $startpos, $endpos } }  
;
```

desc:

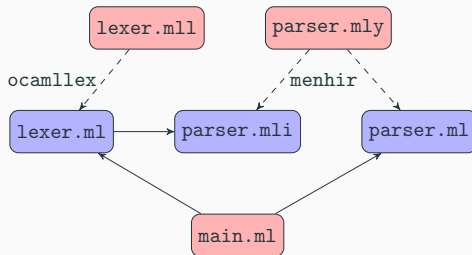
```
| i = INT { Econst i }  
| e1 = expression; PLUS; e2 = expression { Eplus (e1, e2) }  
| ...
```


Compilation and dependencies

As in the case of `ocamllex`, it is necessary to ensure the application of `menhir` before calculating the dependencies

If we use `dune`, we indicate the presence of a `menhir` file:

```
(ocamllex
(modules lexer))
(menhir
(flags --explain --dump)
(modules parser))
(executable
(name minilang)
...)
```



Construction of the automation and the table

Definitions

Defn. (NULL)

Let $\alpha \in (N \cup T)^*$. $\text{NULL}(\alpha)$ holds if and only if we can derive ϵ from α , ie, $\alpha \xrightarrow{*} \epsilon$



Defn. (FIRST)

Let $\alpha \in (N \cup T)^*$. $\text{FIRST}(\alpha)$ is the set of all terminals starting words derived from α , ie,
 $\{a \in T \mid \exists w. \alpha \xrightarrow{*} aw\}$



Defn. (FOLLOW)

Let $X \in N$. $\text{FOLLOW}(X)$ is the set of all terminals that may appear after X in a derivation, ie,
 $\{a \in T \mid \exists u, w. S \xrightarrow{*} uXaw\}$



Computing NULL, FIRST and FOLLOW

To compute $\text{NULL}(\alpha)$, we simply need to compute $\text{NULL}(X)$ for $X \in N$

$\text{NULL}(X)$ holds if and only if

- there exists a production $X \rightarrow \epsilon$
- or there exists a production $X \rightarrow Y_1 \dots Y_m$ where $\text{NULL}(T_i)$ for all i

Issue: this is a set of mutually recursive equations

In other words, if $N = \{X_1, \dots, X_n\}$ and if

$\vec{V} = (\text{NULL}(X_1), \dots, \text{NULL}(X_n))$, we look for the least fixpoint to an equation such as

$$\vec{V} = F(\vec{V})$$

Fixpoint computation

Thm. (existence of a least fixpoint (Tarski))

Let A be a finite set with an order relation \leq and a least element ϵ . Any monotonically increasing function $f : A \rightarrow A$, ie, such that $\forall x, y. x \leq y \implies f(x) \leq f(y)$, admits a *least fixpoint*. □

Proof.

As ϵ is the smallest element, we have $\epsilon \leq f(\epsilon)$. The function f being increasing, therefore we have $f^k(\epsilon) \leq f^{k+1}(\epsilon)$ for every k . Since A is a finite set, there exists a smallest k_0 such that $f^{k_0}(\epsilon) = f^{k_0+1}(\epsilon)$. We have just found a fixed point $a_0 := f^{k_0}(\epsilon)$ of f .

Let b be another fixed point of f . We have $\epsilon \leq b$ and hence $f^k(\epsilon) \leq f^k(b)$ for every k . In particular, $a_0 = f^{k_0}(\epsilon) \leq f^{k_0}(b) = b$. Thus the fixed point a_0 is the smallest fixed point of f . ■

Computing NULL

To compute NULL, we have

$A = \text{BOOL} \times \cdots \times \text{BOOL}$ with $\text{BOOL} = \{\text{true}, \text{false}\}$.

We can equip BOOL with the order $\text{false} \leq \text{true}$ and A with point-wise order

$$(x_1, \dots, x_n) \leq (y_1, \dots, y_n) \quad \text{if and only if} \quad \forall i. x_i \leq y_i.$$

The theorem then applies by taking

$$\epsilon = (\text{false}, \dots, \text{false}).$$

since the function calculating $\text{NULL}(X)$ from $\text{NULL}(X_i)$ is increasing.

Computing NULL

To compute $\text{NULL}(X_i)$, we thus start with

$$\text{NULL}(X_1) = \text{false}, \dots, \text{NULL}(X_n) = \text{false}$$

and we use the equations until we get a fixpoint
i.e. until the values $\text{NULL}(X_i)$ do not change anymore

Example

$$\begin{aligned}
 E &\rightarrow TE' \\
 E' &\rightarrow + TE' \\
 &\mid \epsilon \\
 T &\rightarrow FT' \\
 T' &\rightarrow * FT' \\
 &\mid \epsilon \\
 F &\rightarrow (E) \\
 &\mid \text{int}
 \end{aligned}$$

| E | E' | T | T' | F |
|-------|-------|-------|-------|-------|
| false | false | false | false | false |
| false | true | false | true | false |
| false | true | false | true | false |

Explanation

Why do we seek for a **least** fixpoint?

- ⇒ by induction on the number of steps of the fixpoint computation, we show that if $\text{NULL}(X) = \text{true}$ then $X \xrightarrow{*} \epsilon$
- ⇐ by induction on the number of steps of derivation $X \xrightarrow{*} \epsilon$, we show that $\text{NULL}(X) = \text{true}$ in the previous computation

Computing FIRST

Similarly, the equations defining FIRST are mutually recursive

$$\text{FIRST}(X) = \bigcup_{X \rightarrow \beta} \text{FIRST}(\beta)$$

and

$$\begin{aligned} \text{FIRST}(\epsilon) &= \emptyset \\ \text{FIRST}(a\beta) &= \{a\} \quad \text{if } a \in T \\ \text{FIRST}(X\beta) &= \begin{cases} \text{FIRST}(X) & \text{if } \neg \text{NULL}(X) \\ \text{FIRST}(X) \cup \text{FIRST}(\beta) & \text{if } \text{NULL}(X) \end{cases} \end{aligned}$$

Again, we compute a least fixpoint using Tarski's theorem, with

$$A = \mathcal{P}(X) \times \cdots \times \mathcal{P}(X)$$

equipped with point-wise ordered with \subseteq and with $\epsilon = (\emptyset, \dots, \emptyset)$

Example

$$\begin{aligned}
 E &\rightarrow TE' \\
 E' &\rightarrow + TE' \\
 &\mid \epsilon \\
 T &\rightarrow FT' \\
 T' &\rightarrow * FT' \\
 &\mid \epsilon \\
 F &\rightarrow (E) \\
 &\mid \text{int}
 \end{aligned}$$

NULL

| E | E' | T | T' | F |
|-------|------|-------|------|-------|
| false | true | false | true | false |

FIRST

| E | E' | T | T' | F |
|--------------------|-------------|--------------------|-------------|--------------------|
| \emptyset | \emptyset | \emptyset | \emptyset | \emptyset |
| \emptyset | $\{+\}$ | \emptyset | $\{*\}$ | $\{(\text{int})\}$ |
| \emptyset | $\{+\}$ | $\{(\text{int})\}$ | $\{*\}$ | $\{(\text{int})\}$ |
| $\{(\text{int})\}$ | $\{+\}$ | $\{(\text{int})\}$ | $\{*\}$ | $\{(\text{int})\}$ |
| $\{(\text{int})\}$ | $\{+\}$ | $\{(\text{int})\}$ | $\{*\}$ | $\{(\text{int})\}$ |

Computing FOLLOW

Again, the equations defining FOLLOW are mutually recursive

$$\text{FOLLOW}(X) = \bigcup_{Y \rightarrow \alpha X \beta} \text{FIRST}(\beta) \cup \bigcup_{Y \rightarrow \alpha X \beta, \text{NULL}(\beta)} \text{FOLLOW}(Y)$$

We compute a least fixpoint, using the same domain as for first

Note: we add a special symbol $\#$ in $\text{FOLLOW}(S)$
(which we can do directly, or by adding a rule $S' \rightarrow S\#$)

Example

NULL

| E | E' | T | T' | F |
|-------|------|-------|------|-------|
| false | true | false | true | false |

$$\begin{aligned}
 E &\rightarrow TE' \\
 E' &\rightarrow + TE' \\
 &\mid \epsilon \\
 T &\rightarrow FT' \\
 T' &\rightarrow * FT' \\
 &\mid \epsilon \\
 F &\rightarrow (E) \\
 &\mid \text{int}
 \end{aligned}$$

FIRST

| E | E' | T | T' | F |
|-----------------------|-----------|-----------------------|-----------|-----------------------|
| $\{ (, \text{int} \}$ | $\{ + \}$ | $\{ (, \text{int} \}$ | $\{ * \}$ | $\{ (, \text{int} \}$ |

FOLLOW

| E | E' | T | T' | F |
|---------------|---------------|------------------|------------------|---------------------|
| $\{ \# \}$ | \emptyset | \emptyset | \emptyset | \emptyset |
| $\{ \#,) \}$ | $\{ \# \}$ | $\{ +, \# \}$ | \emptyset | $\{ * \}$ |
| $\{ \#,) \}$ | $\{ \#,) \}$ | $\{ +, \#,) \}$ | $\{ +, \# \}$ | $\{ *, +, \# \}$ |
| $\{ \#,) \}$ | $\{ \#,) \}$ | $\{ +, \#,) \}$ | $\{ +, \#,) \}$ | $\{ *, +, \#,) \}$ |
| $\{ \#,) \}$ | $\{ \#,) \}$ | $\{ +, \#,) \}$ | $\{ +, \#,) \}$ | $\{ *, +, \#,) \}$ |

Exercise

Compute NULL, FIRST, FOLLOW for the grammar LISP

$$\begin{array}{lcl}
 S & \rightarrow & E \# \\
 E & \rightarrow & \text{sym} \\
 & & | (L) \\
 L & \rightarrow & \epsilon \\
 & & | E L
 \end{array}$$

Exercise

Compute NULL, FIRST, FOLLOW for the grammar of the terms of λ -calculus

$$\begin{array}{lcl} E & \rightarrow & T \# \\ T & \rightarrow & \text{nat} \\ & | & (T \ T) \\ & | & \text{lam } T \end{array}$$

LR(0) automation

Let us use $k = 0$ for the moment

We start by constructing an asynchronous automata

that is,

LR(0) automation

The **states** are of the form

$$[X \rightarrow \alpha \bullet \beta]$$

where $X \rightarrow \alpha\beta$ is a production of the grammar; the intuition is

“We want to recognize X , we have already seen α and we still have to see β ”

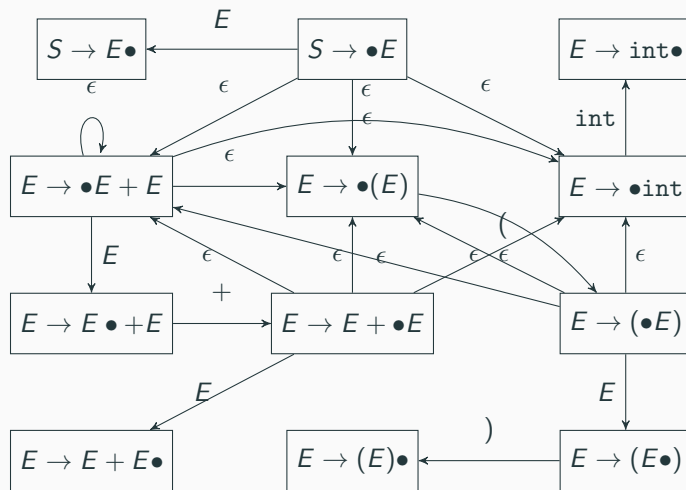
The **transitions** are labeled by $N \cup T$ and are as follows

$$[Y \rightarrow \alpha \bullet a\beta] \xrightarrow{a} [Y \rightarrow \alpha a \bullet \beta]$$

$$[Y \rightarrow \alpha \bullet X\beta] \xrightarrow{X} [Y \rightarrow \alpha X \bullet \beta]$$

$$[Y \rightarrow \alpha \bullet X\beta] \xrightarrow{\epsilon} [Y \rightarrow \bullet \gamma] \quad \text{for all production } X \rightarrow \gamma$$

Example

$$\begin{array}{lcl}
 S & \rightarrow & E \\
 E & \rightarrow & E + E \\
 & | & (E) \\
 & | & \text{int}
 \end{array}$$


Deterministic LR(0) automation

Let us **determinize** the LR(0) automation

For this, we group the states connected by ϵ -transitions

The states of the deterministic automaton are therefore sets of items, such that

| | | |
|-----|---------------|------------------|
| S | \rightarrow | $\bullet E$ |
| E | \rightarrow | $\bullet E + E$ |
| E | \rightarrow | $\bullet (E)$ |
| E | \rightarrow | <code>int</code> |

Deterministic LR(0) automation

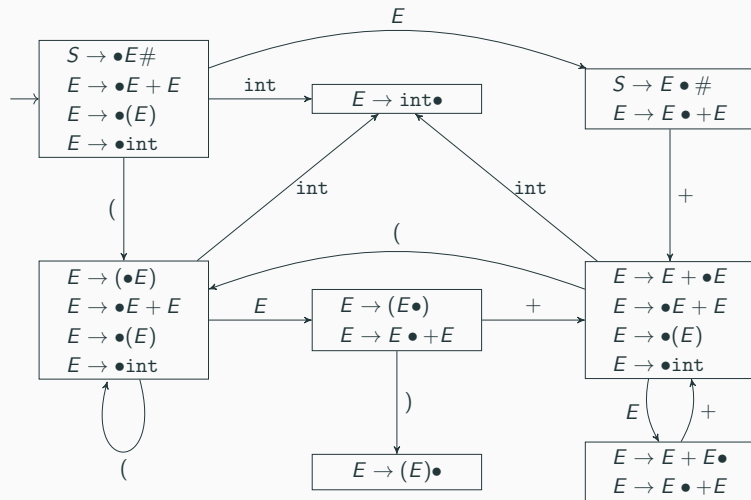
By construction, each state is **saturated** by the property

if $Y \rightarrow \alpha \bullet X \beta \in s$
and if $X \rightarrow \gamma$ is a production
then $X \rightarrow \bullet \gamma \in s$

The initial state is the one containing $S \rightarrow \bullet E \#$

Example

$S \rightarrow E$
 $E \rightarrow E + E$
 $\quad |$
 $\quad (E)$
 $\quad |$
 $\quad \text{int}$



Building the LR(0) table

For the table action, we set

- $\text{action}(s, \#) = \text{success}$ if $[S \rightarrow E \bullet \#] \in s$;
- $\text{action}(s, a) = \text{shift } s'$ if there is a transition $s \xrightarrow{a} s'$;
- $\text{action}(s, a) = \text{reduce } X \rightarrow \beta$ if $[S \rightarrow \beta \bullet] \in s$ for each a ;
- fail in all the other cases.

For the table goto, we set

- $\text{goto}(s, X) = s'$ if and only if there is a transition $s \xrightarrow{X} s'$.

Example

In our example, the table is of follows:

| | action | | | | | goto |
|-------|-----------------------------------|---------|---------|---------|---------|------|
| state | (|) | + | int | # | E |
| 1 | shift 4 | | | shift 2 | | 3 |
| 2 | reduce $E \rightarrow \text{int}$ | | | | | |
| 3 | | | shift 6 | | success | |
| 4 | shift 4 | | | shift 2 | | 5 |
| 5 | | shift 7 | shift 6 | | | |
| 6 | shift 4 | | | shift 2 | | 8 |
| 7 | reduce $E \rightarrow (E)$ | | | | | |
| 8 | | | shift 6 | | | |
| | reduce $E \rightarrow E + E$ | | | | | |

Conflicts

The LR(0) table can contain several possible actions in the same box. We call this a **conflict**. There are two kinds of conflicts:

- a conflict **shift/reduce** if in a state s we can perform a shift but also a reduction;
- a conflict **reduce/reduce**, if in a state s two different reductions are possible.

Defn. (LR(0) grammar)

A grammar is said to be LR(0) if the table thus constructed does not contain any conflicts. □

Conflicts

In our example, there is a conflict in state 8

$$\begin{array}{l} E \rightarrow E + E \bullet \\ E \rightarrow E \bullet + E \end{array}$$

This conflict precisely illustrates the ambiguity of grammar on a word such as `int + int + int`, after reading `int + int`.

We can resolve the conflict in two ways:

- either we favor the **shift**, then translating an associativity on the right;
- either we favor the **reduction**, then translating a left associativity.

Exercise

Show that the grammar for the terms of the λ -calculus is LR(0)

SLR(1) parsing

The construction of LR(0) generates conflicts too easily

We will therefore seek to limit the reductions

A very simple idea is to set $\text{action}(s, a) = \text{reduce } X \rightarrow \beta$ if and only if

$$[X \rightarrow \beta \bullet] \in s \quad \text{and} \quad a \in \text{FOLLOW}(X).$$

Defn. (the class SLR(1))

A grammar is called SLR(1) if the constructed table does not contain any conflict. (SLR stands for *simple LR*.)



Example

The grammar

$$\begin{array}{lcl}
 S & \rightarrow & E\# \\
 E & \rightarrow & E + T \\
 & | & T \\
 T & \rightarrow & T * F \\
 & | & F \\
 F & \rightarrow & (E) \\
 & | & \text{int}
 \end{array}$$

is SLR(1)

Exercise: Verify it (the automata contains 12 states).

Exercise: Show that the grammar for LISP is SLR(1).

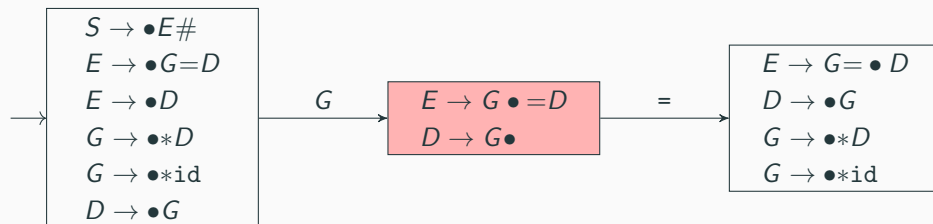
Limitation of SLR(1) parsing

In practice, SLR(1) grammars are not powerful enough

Example

$$\begin{aligned} S &\rightarrow E\# \\ E &\rightarrow G=D \\ &\quad | \quad D \\ G &\rightarrow *D \\ &\quad | \quad \text{id} \\ D &\rightarrow G \end{aligned}$$

| | | |
|---|-------------------------------------|-----|
| | = | |
| 1 | ... | ... |
| 2 | shift 3 reduce $D \rightarrow G$ | ... |
| 3 | : | .. |



LR(1) parsing

We introduce a larger class of grammars, LR(1), with larger tables

items now look like

$$[X \rightarrow \alpha \bullet \beta, a]$$

and the meaning is “we want to recognize X , we have already seen α , we still need to see β and then to check that the next token is a ”

LR(1) parsing

The non-deterministic LR(1) automaton has transitions

$$[Y \rightarrow \alpha \bullet a\beta, b] \xrightarrow{a} [Y \rightarrow \alpha a \bullet \beta, b]$$

$$[Y \rightarrow \alpha \bullet X\beta, b] \xrightarrow{X} [Y \rightarrow \alpha X \bullet \beta, b]$$

$$[Y \rightarrow \alpha \bullet X\beta, b] \xrightarrow{\epsilon} [Y \rightarrow \bullet \gamma, c] \quad \text{for each } c \in \text{FIRST}(\beta b)$$

the initial state is that containing $[S \rightarrow \bullet \alpha, \#]$

As before, we can determinize the automaton and construct the corresponding table; we introduce a reduction action for (s, a) only when s contains an item of the form $[X \rightarrow \alpha \bullet, a]$

Defn. (the class LR(1))

A grammar is said to be LR(1) if the resulting table contains no conflict.

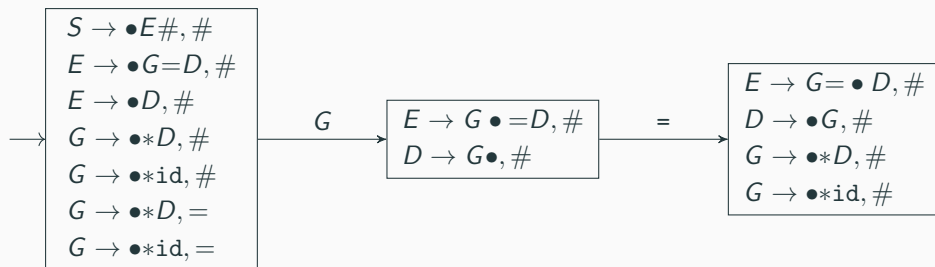


Example

Example

$S \rightarrow E\#$
 $E \rightarrow G=D$
 $\quad | \quad D$
 $G \rightarrow *D$
 $\quad | \quad id$
 $D \rightarrow G$

| | # | = | |
|---|--------------------------|----------|----------|
| 1 | ... | ... | ... |
| 2 | reduce $D \rightarrow G$ | shift 3 | ... |
| 3 | \vdots | \vdots | \ddots |



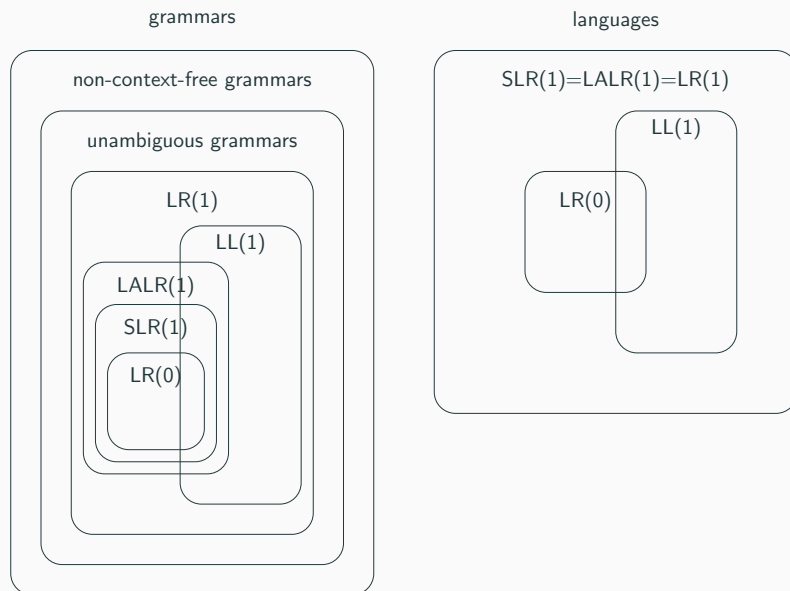
LALR(1) parsing

The LR(1) tables can be large, so we introduced approximations

The class LALR(1) (*lookahead LR*) is such an approximation, used in tools of the yacc family

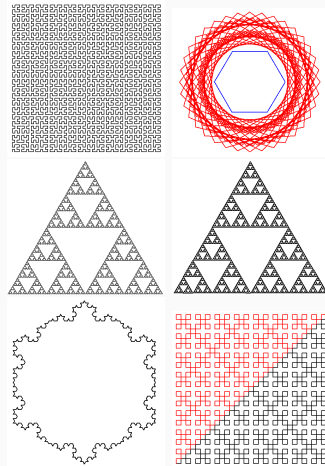
For more details, see Compilers (“the dragon book”) by A. Aho, R. Sethi, J. Ullman, [ALSU06, section 4.7]

Grammar hierarchies



Next

- HW4
 - syntax analysis of mini-Turtle
 - Java or OCaml
- more parsing



Localization

Localization

The `ocamllex` tool maintains, in the `Lexing.lexbuf` type structure, the current position in the source text that is analyzed (file name, line, column)

```
type position = ...
```

We can obtain the location of the last string recognized by `ocamllex`

```
val lexeme_start_p: lexbuf -> position  
val lexeme_end_p : lexbuf -> position
```

Localization

This information is used to locate a syntax error

```
let lb = Lexing.from_channel c in
try
  let f = Parser.file Lexer.next_token lb in
  ...
with
| Parser.Error ->
  let pos = Lexing.lexeme_start_p lb in
  fprintf "%d: syntax error" pos.pos_lnum;
  exit 1
```

but also possibly lexical errors like a string or unclosed comments

Localization in Menhir

The `menhir` tool retrieves this information and provides it in two values `$startpos` and `$endpos` of the `Lexing.position` type

In a semantic action, they correspond to the beginning and the end of the text that has been recognized by the grammar rule

This information can be stored in the abstract syntax tree

```
expression:  
| e1 = expression; PLUS; e2 = expression  
  { Add ($startpos, $endpos, e1, e2) }
```

(anyway we have seen an even more elegant solution)

Elementary parsing

A good exercise

Let us write, in the most elementary way possible, a parser for arithmetic expressions including

- constants
- addition
- multiplication
- parentheses

More precisely

Starting point: a lexical analyzer (e.g. written with `ocamllex`)

```
type token =  
  | CONST of int  
  | PLUS  
  | TIMES  
  | LEFTPAR  
  | RIGHTPAR  
  | EOF
```

Goal: an abstract syntax tree

```
type expr =  
  | Const of int  
  | Add of expr * expr  
  | Mul of expr * expr
```

A good advice

Start by writing a **print function** (pretty-printer)

```
let rec print fmt = function
  | Add (e1, e2) -> fprintf fmt "%a +@ %a" print e1 print e2
  | e             -> print2 fmt e
and print2 fmt = function
  | Mul (e1, e2) -> fprintf fmt "%a *@ %a" print2 e1 print2 e2
  | e             -> print3 fmt e
and print3 fmt = function
  | Const n       -> fprintf fmt "%d" n
  | e             -> fprintf fmt "(@[%a@])" print e
```

(here we use the library Format)

References i



Alfred V. Aho, Monica S. Lam, Ravi Sethi, and Jeffrey D. Ullman.
Compilers: Principles, Techniques, and Tools (2nd Edition).
Addison-Wesley Longman Publishing Co., Inc., USA, 2006.

Questions?