

# Non-recursive predictive parsing

---

Observation:

*Our recursive descent parser encodes state information in its run-time stack, or call stack.*

Using recursive procedure calls to implement a stack abstraction may not be particularly efficient.

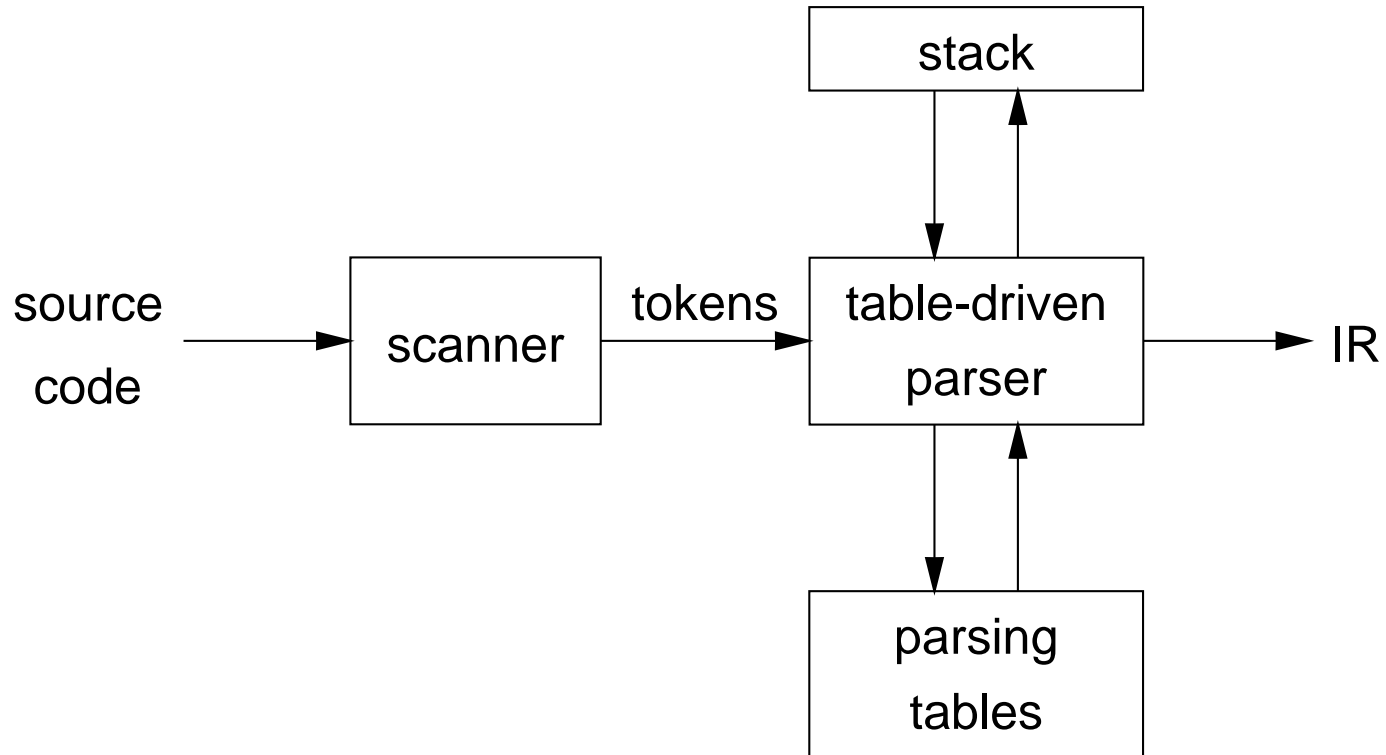
This suggests other implementation methods:

- explicit stack, hand-coded parser
- stack-based, table-driven parser

# Non-recursive predictive parsing

---

Now, a predictive parser looks like:



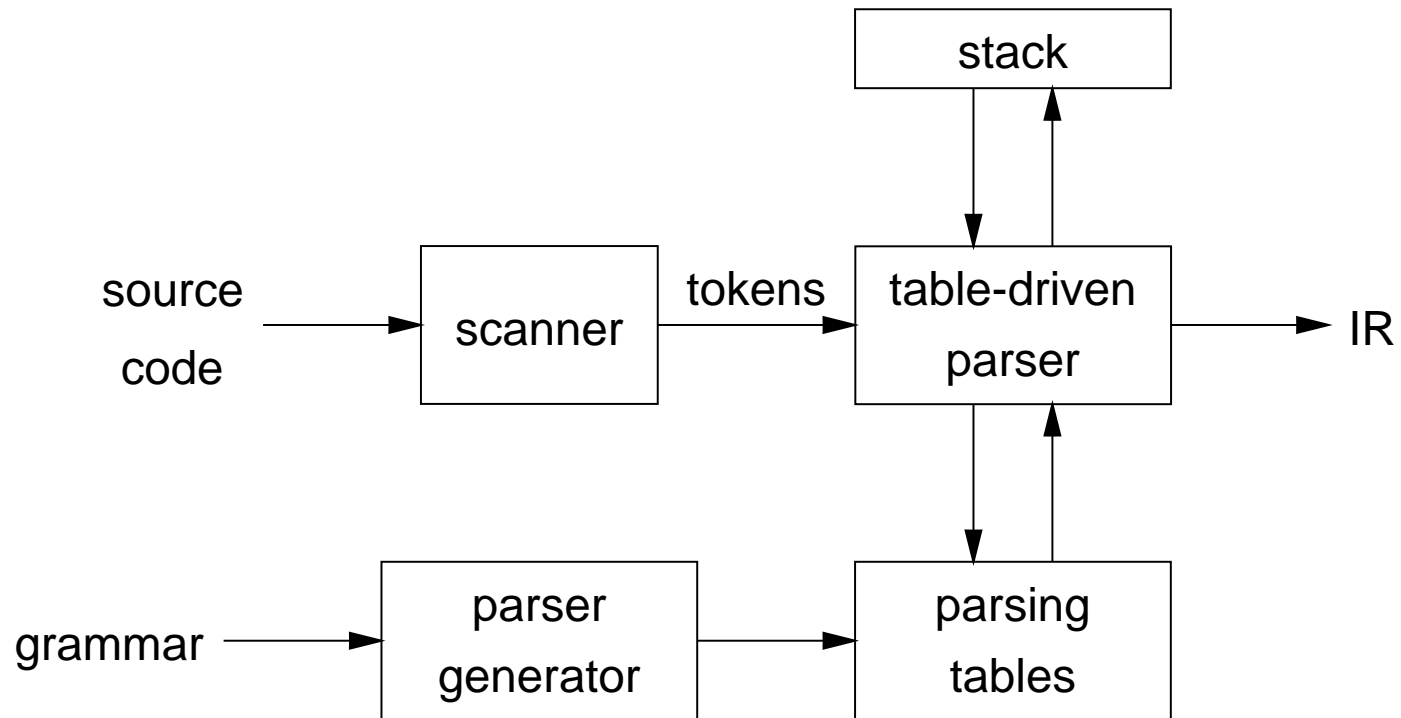
Rather than writing code, we build tables.

*Building tables can be automated!*

# Table-driven parsers

---

A parser generator system often looks like:



This is true for both top-down (LL) and bottom-up (LR) parsers

# Non-recursive predictive parsing

---

*Input:* a string  $w$  and a parsing table  $M$  for  $G$

```
tos ← 0
Stack[tos] ← EOF
Stack[++tos] ← Start Symbol
token ← next_token()
repeat
  X ← Stack[tos]
  if X is a terminal or EOF then
    if X = token then
      pop X
      token ← next_token()
    else error()
  else /* X is a non-terminal */
    if  $M[X, token] = X \rightarrow Y_1 Y_2 \dots Y_k$  then
      pop X
      push  $Y_k, Y_{k-1}, \dots, Y_1$ 
    else error()
until X = EOF
```

# Non-recursive predictive parsing

---

What we need now is a parsing table  $M$ .

Our expression grammar:

1	$\langle \text{goal} \rangle$	$::=$	$\langle \text{expr} \rangle$
2	$\langle \text{expr} \rangle$	$::=$	$\langle \text{term} \rangle \langle \text{expr}' \rangle$
3	$\langle \text{expr}' \rangle$	$::=$	$+\langle \text{expr} \rangle$
4			$-\langle \text{expr} \rangle$
5			$\epsilon$
6	$\langle \text{term} \rangle$	$::=$	$\langle \text{factor} \rangle \langle \text{term}' \rangle$
7	$\langle \text{term}' \rangle$	$::=$	$*\langle \text{term} \rangle$
8			$/\langle \text{term} \rangle$
9			$\epsilon$
10	$\langle \text{factor} \rangle$	$::=$	num
11			id

Its parse table:

	id	num	+	-	*	/	$\$^\dagger$
$\langle \text{goal} \rangle$	1	1	-	-	-	-	-
$\langle \text{expr} \rangle$	2	2	-	-	-	-	-
$\langle \text{expr}' \rangle$	-	-	3	4	-	-	5
$\langle \text{term} \rangle$	6	6	-	-	-	-	-
$\langle \text{term}' \rangle$	-	-	9	9	7	8	9
$\langle \text{factor} \rangle$	11	10	-	-	-	-	-

$\dagger$  we use  $\$$  to represent EOF

# FIRST

---

For a string of grammar symbols  $\alpha$ , define  $\text{FIRST}(\alpha)$  as:

- the set of terminal symbols that begin strings derived from  $\alpha$ :  
 $\{a \in V_t \mid \alpha \Rightarrow^* a\beta\}$
- If  $\alpha \Rightarrow^* \varepsilon$  then  $\varepsilon \in \text{FIRST}(\alpha)$

$\text{FIRST}(\alpha)$  contains the set of tokens valid in the initial position in  $\alpha$

To build  $\text{FIRST}(X)$ :

1. If  $X \in V_t$  then  $\text{FIRST}(X)$  is  $\{X\}$
2. If  $X \rightarrow \varepsilon$  then add  $\varepsilon$  to  $\text{FIRST}(X)$
3. If  $X \rightarrow Y_1 Y_2 \cdots Y_k$ :
  - (a) Put  $\text{FIRST}(Y_1) - \{\varepsilon\}$  in  $\text{FIRST}(X)$
  - (b)  $\forall i : 1 < i \leq k$ , if  $\varepsilon \in \text{FIRST}(Y_1) \cap \cdots \cap \text{FIRST}(Y_{i-1})$   
(i.e.,  $Y_1 \cdots Y_{i-1} \Rightarrow^* \varepsilon$ )  
then put  $\text{FIRST}(Y_i) - \{\varepsilon\}$  in  $\text{FIRST}(X)$
  - (c) If  $\varepsilon \in \text{FIRST}(Y_1) \cap \cdots \cap \text{FIRST}(Y_k)$  then put  $\varepsilon$  in  $\text{FIRST}(X)$

Repeat until no more additions can be made.

# FOLLOW

---

For a non-terminal  $A$ , define  $\text{FOLLOW}(A)$  as

the set of terminals that can appear immediately to the right of  $A$  in some sentential form

Thus, a non-terminal's FOLLOW set specifies the tokens that can legally appear after it.

A terminal symbol has no FOLLOW set.

To build  $\text{FOLLOW}(A)$ :

1. Put  $\$$  in  $\text{FOLLOW}(\langle \text{goal} \rangle)$
2. If  $A \rightarrow \alpha B \beta$ :
  - (a) Put  $\text{FIRST}(\beta) - \{\epsilon\}$  in  $\text{FOLLOW}(B)$
  - (b) If  $\beta = \epsilon$  (i.e.,  $A \rightarrow \alpha B$ ) or  $\epsilon \in \text{FIRST}(\beta)$  (i.e.,  $\beta \Rightarrow^* \epsilon$ ) then put  $\text{FOLLOW}(A)$  in  $\text{FOLLOW}(B)$

Repeat until no more additions can be made

# LL(1) grammars

---

## *Previous definition*

A grammar  $G$  is LL(1) iff. for all non-terminals  $A$ , each distinct pair of productions  $A \rightarrow \beta$  and  $A \rightarrow \gamma$  satisfy the condition  $\text{FIRST}(\beta) \cap \text{FIRST}(\gamma) = \phi$ .

What if  $A \Rightarrow^* \epsilon$ ?

## *Revised definition*

A grammar  $G$  is LL(1) iff. for each set of productions

$A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ :

1.  $\text{FIRST}(\alpha_1), \text{FIRST}(\alpha_2), \dots, \text{FIRST}(\alpha_n)$  are all pairwise disjoint
2. If  $\alpha_i \Rightarrow^* \epsilon$  then  $\text{FIRST}(\alpha_j) \cap \text{FOLLOW}(A) = \phi, \forall 1 \leq j \leq n, i \neq j$ .

If  $G$  is  $\epsilon$ -free, condition 1 is sufficient.



# LL(1) grammars

---

Provable facts about LL(1) grammars:

1. No left-recursive grammar is LL(1)
2. No ambiguous grammar is LL(1)
3. Some languages have no LL(1) grammar
4. A  $\epsilon$ -free grammar where each alternative expansion for  $A$  begins with a distinct terminal is a *simple* LL(1) grammar.

Example

- $S \rightarrow aS \mid a$  is not LL(1) because  $\text{FIRST}(aS) = \text{FIRST}(a) = \{a\}$
- $S \rightarrow aS'$   
 $S' \rightarrow aS' \mid \epsilon$   
accepts the same language and is LL(1)

# LL(1) parse table construction

---

*Input:* Grammar  $G$

*Output:* Parsing table  $M$

*Method:*

1.  $\forall$  productions  $A \rightarrow \alpha$ :

(a)  $\forall a \in \text{FIRST}(\alpha)$ , add  $A \rightarrow \alpha$  to  $M[A, a]$

(b) If  $\epsilon \in \text{FIRST}(\alpha)$ :

i.  $\forall b \in \text{FOLLOW}(A)$ , add  $A \rightarrow \alpha$  to  $M[A, b]$

ii. If  $\$ \in \text{FOLLOW}(A)$  then add  $A \rightarrow \alpha$  to  $M[A, \$]$

2. Set each undefined entry of  $M$  to error

If  $\exists M[A, a]$  with multiple entries then grammar is not LL(1).

Note: recall  $a, b \in V_t$ , so  $a, b \neq \epsilon$

# Example

---

Our long-suffering expression grammar:

$$\begin{array}{l|l}
 S \rightarrow E & T \rightarrow FT' \\
 E \rightarrow TE' & T' \rightarrow *T \mid /T \mid \varepsilon \\
 E' \rightarrow +E \mid -E \mid \varepsilon & F \rightarrow \text{id} \mid \text{num}
 \end{array}$$

	FIRST	FOLLOW
$S$	{num, id}	{ $\$$ }
$E$	{num, id}	{ $\$$ }
$E'$	{ $\varepsilon$ , +, -}	{ $\$$ }
$T$	{num, id}	{+, -, $\$$ }
$T'$	{ $\varepsilon$ , *, /}	{+, -, $\$$ }
$F$	{num, id}	{+, -, *, /, $\$$ }
id	{id}	-
num	{num}	-
*	{*}	-
/	{/}	-
+	{+}	-
-	{-}	-

	id	num	+	-	*	/	$\$$
$S$	$S \rightarrow E$	$S \rightarrow E$	-	-	-	-	-
$E$	$E \rightarrow TE'$	$E \rightarrow TE'$	-	-	-	-	-
$E'$	-	-	$E' \rightarrow +E$	$E' \rightarrow -E$	-	-	$E' \rightarrow \varepsilon$
$T$	$T \rightarrow FT'$	$T \rightarrow FT'$	-	-	-	-	-
$T'$	-	-	$T' \rightarrow \varepsilon$	$T' \rightarrow \varepsilon$	$T' \rightarrow *T$	$T' \rightarrow /T$	$T' \rightarrow \varepsilon$
$F$	$F \rightarrow \text{id}$	$F \rightarrow \text{num}$	-	-	-	-	-

# Building the tree

---

Again, we insert code at the right points:

```
tos ← 0
Stack[tos] ← EOF
Stack[++tos] ← root node
Stack[++tos] ← Start Symbol
token ← next_token()
repeat
  X ← Stack[tos]
  if X is a terminal or EOF then
    if X = token then
      pop X
      token ← next_token()
      pop and fill in node
    else error()
  else /* X is a non-terminal */
    if  $M[X,token] = X \rightarrow Y_1Y_2\cdots Y_k$  then
      pop X
      pop node for X
      build node for each child and
      make it a child of node for X
      push  $n_k, Y_k, n_{k-1}, Y_{k-1}, \cdots, n_1, Y_1$ 
    else error()
until X = EOF
```

# A grammar that is not LL(1)

---

$$\begin{aligned}\langle \text{stmt} \rangle & ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \\ & \quad | \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \text{ else } \langle \text{stmt} \rangle \\ & \quad | \dots\end{aligned}$$

Left-factored:

$$\begin{aligned}\langle \text{stmt} \rangle & ::= \text{if } \langle \text{expr} \rangle \text{ then } \langle \text{stmt} \rangle \langle \text{stmt}' \rangle | \dots \\ \langle \text{stmt}' \rangle & ::= \text{else } \langle \text{stmt} \rangle | \varepsilon\end{aligned}$$

Now,  $\text{FIRST}(\langle \text{stmt}' \rangle) = \{\varepsilon, \text{else}\}$

Also,  $\text{FOLLOW}(\langle \text{stmt}' \rangle) = \{\text{else}, \$\}$

But,  $\text{FIRST}(\langle \text{stmt}' \rangle) \cap \text{FOLLOW}(\langle \text{stmt}' \rangle) = \{\text{else}\} \neq \emptyset$

On seeing `else`, conflict between choosing

$$\langle \text{stmt}' \rangle ::= \text{else } \langle \text{stmt} \rangle \quad \text{and} \quad \langle \text{stmt}' \rangle ::= \varepsilon$$

$\Rightarrow$  grammar is not LL(1)!

The fix:

Put priority on  $\langle \text{stmt}' \rangle ::= \text{else } \langle \text{stmt} \rangle$  to associate `else` with closest previous `then`.

# Error recovery

---

Key notion:

- For each non-terminal, construct a set of terminals on which the parser can synchronize
- When an error occurs looking for  $A$ , scan until an element of  $\text{SYNCH}(A)$  is found

Building  $\text{SYNCH}$ :

1.  $a \in \text{FOLLOW}(A) \Rightarrow a \in \text{SYNCH}(A)$
2. place keywords that start statements in  $\text{SYNCH}(A)$
3. add symbols in  $\text{FIRST}(A)$  to  $\text{SYNCH}(A)$

If we can't match a terminal on top of stack:

1. pop the terminal
2. print a message saying the terminal was inserted
3. continue the parse

(i.e.,  $\text{SYNCH}(a) = V_t - \{a\}$ )