

Non-recursive predictive parsing

Observation:

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

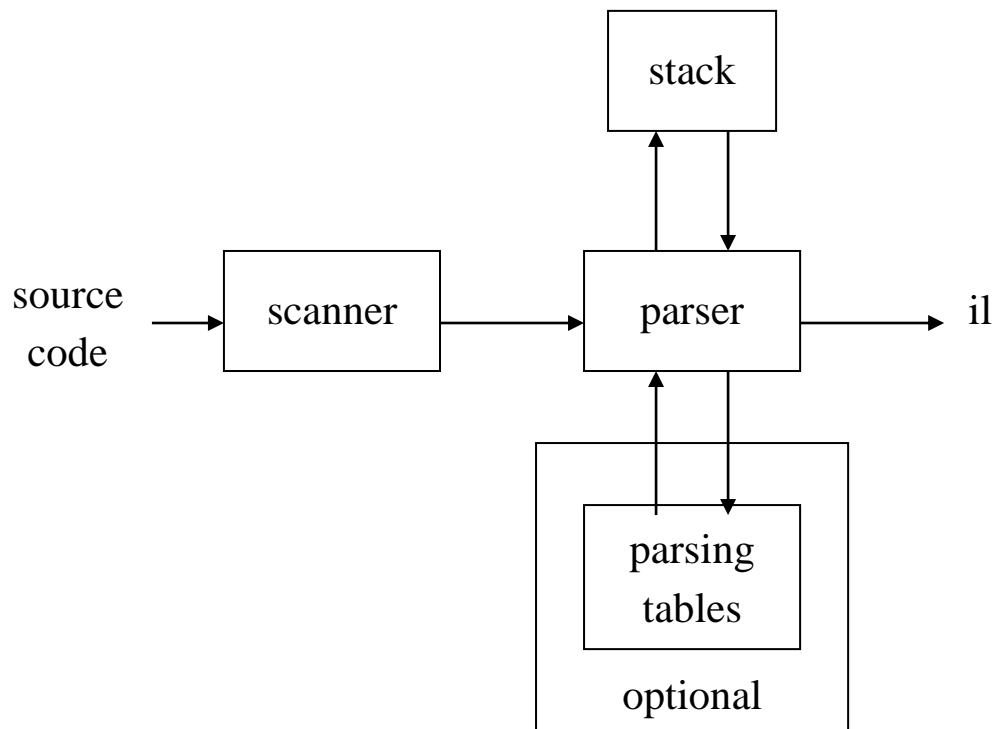
Using recursive procedure calls to implement a stack abstraction may not particularly efficient.

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:

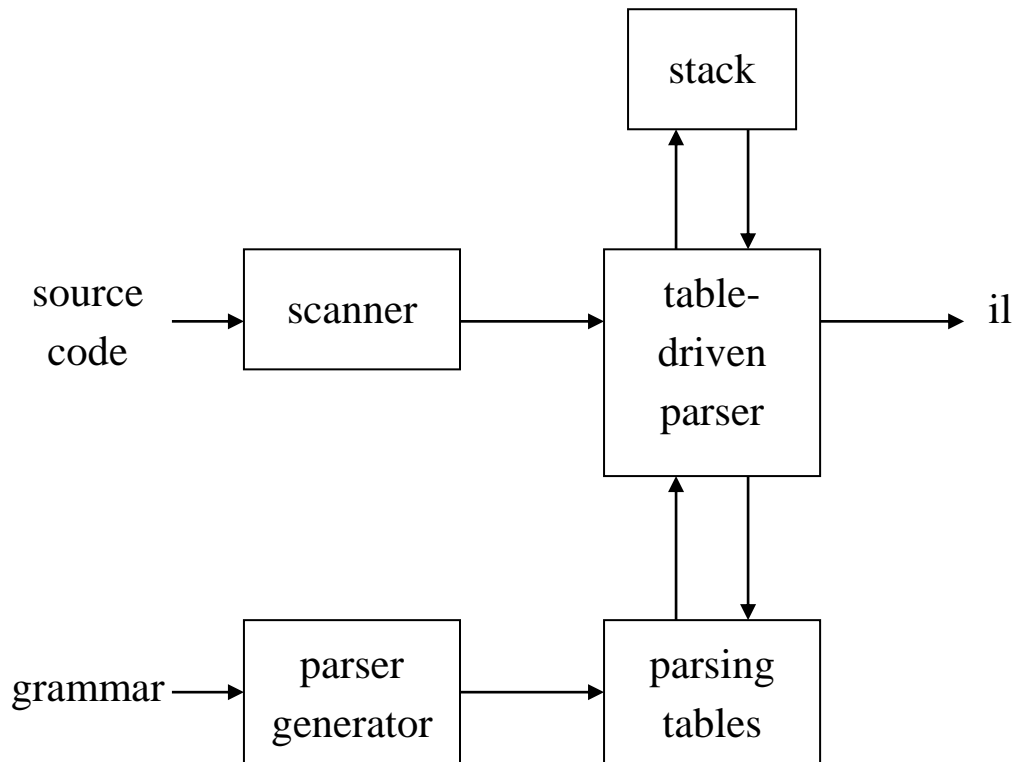


Rather than writing code, we build tables.

Building tables can be automated!

Table-driven parsers

A parser generator system often looks like:



This is true for both top down and bottom up parsers

LL(1): left to right scan, leftmost derivation, one lookahead symbol

LR(1): left to right scan, rightmost derivation, one lookahead symbol

Table-driven parsing algorithm

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

$\text{tos} \leftarrow 0$

$\text{Stack}[\text{tos}++] \leftarrow \text{eof}$

$\text{Stack}[\text{tos}++] \leftarrow \text{Start Symbol}$

$\text{token} \leftarrow \text{next_token}()$

$X \leftarrow \text{Stack}[\text{tos}]$

repeat

 if X is a terminal or eof then

 if $X = \text{token}$ then

 pop X

$\text{token} \leftarrow \text{next_token}()$

 else error()

 else /* X is a non-terminal */

 if $M[X, \text{token}] = X \rightarrow Y_1 Y_2 \dots Y_k$ then

 pop X

 push Y_k, Y_{k-1}, \dots, Y_1

 else error()

$X \leftarrow \text{Stack}[\text{tos}]$

until $X = \text{eof}$

Aho, Sethi, and Ullman, Algorithm 4.34

The grammar and its table

Our long-suffering expression grammar

$\langle \text{goal} \rangle ::= \langle \text{expr} \rangle$
 $\langle \text{expr} \rangle ::= \langle \text{term} \rangle \langle \text{expr}' \rangle$
 $\langle \text{expr}' \rangle ::= + \langle \text{term} \rangle \langle \text{expr}' \rangle$
 $| - \langle \text{term} \rangle \langle \text{expr}' \rangle$
 $| \varepsilon$
 $\langle \text{term} \rangle ::= \langle \text{factor} \rangle \langle \text{term}' \rangle$
 $\langle \text{term}' \rangle ::= * \langle \text{term}' \rangle$
 $| / \langle \text{term}' \rangle$
 $| \varepsilon$
 $\langle \text{factor} \rangle ::= \text{number}$
 $| \text{id}$

LL(1) parsing table

	id	num	+	−	*	/	eof
$\langle \text{goal} \rangle$	$g \rightarrow e$	$g \rightarrow e$	-	-	-	-	-
$\langle \text{expr} \rangle$	$e \rightarrow te'$	$e \rightarrow te'$	-	-	-	-	-
$\langle \text{expr}' \rangle$	-	-	$e' \rightarrow +e$	$e' \rightarrow -e$	-	-	$e' \rightarrow \varepsilon$
$\langle \text{term} \rangle$	$t \rightarrow ft'$	$t \rightarrow ft'$	-	-	-	-	-
$\langle \text{term}' \rangle$	-	-	$t' \rightarrow \varepsilon$	$t' \rightarrow \varepsilon$	$t' \rightarrow *t$	$t' \rightarrow /t$	$t' \rightarrow \varepsilon$
$\langle \text{factor} \rangle$	$f \rightarrow \text{id}$	$f \rightarrow \text{num}$					

The FIRST set

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

- the set of terminal symbols that begin strings derived from α
- if $\alpha \Rightarrow^* \varepsilon$, then $\varepsilon \in \text{FIRST}(\alpha)$

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

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

1. if X is a terminal, $\text{FIRST}(X)$ is $\{X\}$
2. if $X ::= \varepsilon$, then $\varepsilon \in \text{FIRST}(X)$
3. if $X ::= Y_1 Y_2 \dots Y_k$ then put $\text{FIRST}(Y_1)$ in $\text{FIRST}(X)$
4. if X is a non-terminal and $X ::= Y_1 Y_2 \dots Y_k$, then
 $a \in \text{FIRST}(X)$ if $a \in \text{FIRST}(Y_i)$ and $\varepsilon \in \text{FIRST}(Y_j)$ for all
 $1 \leq j < i$
 (If $\varepsilon \notin \text{FIRST}(Y_1)$, then $\text{FIRST}(Y_i)$ is irrelevant, for $1 < i$)

The FIRST construction

rule	1	2	3	4	FIRST
goal	-	-	num, id	-	{num,id}
expr	-	-	num, id	-	{num,id}
expr'	-	ε	+, -	-	{ ε ,+, -}
term	-	-	num, id	-	{num,id}
term'	-	ε	*, /	-	{ ε ,*, /}
factor	-	-	num, id	-	{num,id}
num	num	-	-	-	{num}
id	id	-	-	-	{id}
+	+	-	-	-	{+}
-	-	-	-	-	{-}
*	*	-	-	-	{*}
/	/	-	-	-	{/}

The FOLLOW set

For a non-terminal A, define 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 legally appear after it

A terminal symbol has no FOLLOW set

To build FOLLOW(X)

1. place eof in FOLLOW(<goal>)
2. if $A ::= \alpha B \beta$, then put $\{\text{FIRST}(\beta) - \varepsilon\}$ in FOLLOW(B)
3. if $A ::= \alpha B$ then put FOLLOW(A) in FOLLOW(B)
4. if $A ::= \alpha B \beta$ and $\varepsilon \in \text{FIRST}(\beta)$, then put FOLLOW(A) in FOLLOW(B)

The FOLLOW construction

rule	1	2	3	4	FOLLOW
goal	eof	-	-	-	{eof}
expr	-	-	eof	-	{eof}
expr'	-	-	eof	-	{eof}
term	-	+,-	-	eof	{eof, +,-}
term'	-	-	eof, +,-	-	{ ϵ ,*, /}
factor	-	*,/	-	eof, +,-	{ eof, +,-,*,/}

LL(1) parsing table construction

Input: a grammar G

Method

1. \forall production $A ::= \alpha$, perform steps 2-4
2. \forall terminal a in $\text{FIRST}(\alpha)$, add $A ::= \alpha$ to $M[A, a]$
3. if $\varepsilon \in \text{FIRST}(\alpha)$, for any terminal $b \in \text{FOLLOW}(A)$
add $A ::= \alpha$ to $M[A, b]$
4. if $\varepsilon \in \text{FIRST}(\alpha)$, and $\text{eof} \in \text{FOLLOW}(A)$,
add $A ::= \alpha$ to $M[A, \text{eof}]$
5. set each undefined entry of M to error

If this fails, the grammar is not LL(1)

Aho, Sethi, and Ullman, Algorithm 4.31

LL(1) parsing table for example

	id	num	+	−	*	/	eof
<goal>	$g \rightarrow e$	$g \rightarrow e$	-	-	-	-	-
<expr>	$e \rightarrow te'$	$e \rightarrow te'$	-	-	-	-	-
<expr'>	-	-	$e' \rightarrow +e$	$e' \rightarrow -e$	-	-	$e' \rightarrow e$
<term>	$t \rightarrow ft'$	$t \rightarrow ft'$	-	-	-	-	-
<term'>	-	-	$t' \rightarrow \varepsilon$	$t' \rightarrow \varepsilon$	$t' \rightarrow *t$	$t' \rightarrow /t$	$t' \rightarrow \varepsilon$
<factor>	$f \rightarrow id$	$f \rightarrow num$					

Symbol	FIRST	FOLLOW
<goal>	{id, number}	{eof}
<expr>	{id, number}	{eof}
<expr'>	{ ε , +, −}	{eof}
<term>	{id, number}	{eof, +, −}
<term'>	{ ε , *, /}	{eof, +, −}
<factor>	{id, number}	{eof, +, −, *, /}
+	{+}	-
−	{−}	-
*	{*}	-
/	{/}	-
id	{id}	-
number	{number}	-

Building the tree

Insert some code at the appropriate points

```
tos  $\leftarrow$  0
```

```
Stack[tos++]  $\leftarrow$  eof
```

```
Stack[tos++]  $\leftarrow$  root node
```

```
Stack[tos++]  $\leftarrow$  Start Symbol
```

```
token  $\leftarrow$  next_token()
```

```
X  $\leftarrow$  Stack[tos]
```

```
repeat
```

```
  if X is a terminal or eof then
```

```
    if X = token then
```

```
      pop X
```

```
      token  $\leftarrow$  next_token()
```

```
      pop and fill in node
```

```
    else error()
```

```
  else /* X is a non-terminal */
```

```
    if  $M[X, \text{token}] = X \rightarrow Y_1 Y_2 \dots Y_k$  then
```

```
      pop X
```

```
      pop node for X
```

```
      build node for each child and make it a child of node for X
```

```
      push  $n_k, Y_k, n_{k-1}, Y_{k-1}, \dots, n_1, Y_1$ 
```

```
    else error()
```

```
  X  $\leftarrow$  Stack[tos]
```

```
until X = eof
```

LL(1) grammars

Features

- input parsed from left to right
- leftmost derivation
- one token lookahead

Definition

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

A grammar G is LL(1) if and only if for each set of productions $A ::= \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$

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

If G is ε -free, condition 1 is sufficient.

LL(1) grammars

Provable facts about LL(1) grammars:

- no left recursive grammar is LL(1)
- no ambiguous grammar is LL(1)
- LL(1) parsers operate in linear time
- an ε -free grammar where each alternative expansion for A begins with a distinct terminal is a simple LL(1) grammar

Not all grammars are LL(1)

- $S ::= aS \mid a$
is not LL(1)
 $\text{FIRST}(aS) = \text{FIRST}(a) = \{a\}$

- $S ::= aS'$
 $S' ::= aS' \mid \varepsilon$
accepts the same language and is LL(1)

LL grammars

LL(1) grammars

- may need to rewrite grammar
(left recursion, left factoring)
- resulting grammar larger, less maintainable

LL(k) grammars

- k-token lookahead
- more powerful than LL(1) grammars
- example:
 $S ::= ac \mid abc$ is LL(2)

Not all grammars are LL(k)

- example:
 $S ::= a^i b^j$ where $i \geq j$
- equivalent to dangling else problem
- problem – must choose production after k tokens of lookahead

Bottom-up parsers avoid this problem