

Recursive Descent Parsing (Top-down)

The parse tree is constructed (a) from the top, and (b) from left to right

Tokens are seen in order of appearance in the token stream

Example token stream (stream of terminals): t2 t5 t6 t8 t9

```
1 --> t2 3 9
3 --> 4 7
4 --> t5 t6
7 --> t8 t9
```

Consider the grammar

```
E --> T | T + E
T --> int | int * T | (E)
```

Token stream is: (int5)

Start with top-level non-terminal E (hence called top-down). Try rules for E *in order*. Walk through the input with an arrow left to right. Also try rules in order. As long as we are generating non-terminals, we do not know if we are on the right path or not. However if we hit a terminal, we can see if that matches our current token or not. If it does not, we can back-track and try another option (in order). In this case we go to T and then down to int, backtrack, then int * T, then backtrack, and then (E), and here '(' matches! So we might be on the right track, and we now need to advance the input pointer, and now we have to expand the non-terminal E. Again we will start with the first production T. And we will eventually match the input string with some production, and we are done.

General algorithm for Recursive Descent Parsing

Let TOKEN be the type of the tokens: e.g., INT, OPEN, CLOSE, PLUS, TIMES, ... Let next point to the next token in the input string (the arrow).

Define boolean functions that check for a match of:

- A given token terminal

```
bool term(TOKEN, tok) { return *next++ == tok; } //next pointer is incremented regardless!
Try the nth production of non-terminal S
bool Sn() { ... } //will succeed if input matches the nth production of S
Try all productions of S
bool S() { ... } //will succeed if input matches any production of S
```

Running on our example, functions for non-terminal E

- For production E --> T

```
bool E1() { return T(); }
```

- For production E --> T + E

```
bool E2() { return T() AND term(PLUS) AND E(); }
```

- For all productions of E (with backtracking)

```
bool E() {
    TOKEN *save = next;
    return (next = save, E1())
        || (next = save, E2());
}
```

The first rule that returns true will cause a return (we will not evaluate the remaining rules due to semantics of logical OR ||)

next pointer needs to be saved for backtracking

Functions for non-terminal T

- For production T --> int

```
bool T1() { return term(INT); }
```

- For production $E \rightarrow \text{int} * T$

```
bool T2() { return term(INT) AND term(TIMES) AND T(); }
```

- For production $E \rightarrow (E)$

```
bool T3() { return term(OPEN) AND E() AND term(CLOSE); }
```

```
bool T() {
    TOKEN *save = next;
    return (next = save, T1())
        || (next = save, T2())
        || (next = save, T3());
}
```

The first rule that returns true will cause a return (we will not evaluate the remaining rules due to semantics of logical OR ||)
next pointer needs to be saved for backtracking

To start the parser:

1. Initialize next to first token
2. Invoke E()
 - Show full execution of E() on the example string

Can implement by hand, e.g., TinyCC

Left Recursion : the main problem with Recursive Descent Parsing

Consider

$S \rightarrow S a$

```
bool S1() { return S() AND term(a); }
bool S() { return S1(); }
```

This will go into an infinite loop for parsing *any* input.

The problem is that this grammar is left recursive

$S \rightarrow^+ S \alpha$ for some α

If for some sequence of productions we end up with left recursion, recursive descent parsing will just not work

Consider the left-recursive grammar

$S \rightarrow S a \mid b$

This generates a string starting with b with any number of as following it. It produces the sentences right-to-left, and that's why it does not work with recursive descent parsing (which works left-to-right).

We can generate exactly the same language by producing strings left to right, which is also a solution to our problem.

$S \rightarrow b S'$
 $S' \rightarrow a S' \mid \epsilon$

In general, for any left recursive grammar

$S \rightarrow S a_1 \mid S a_2 \mid \dots \mid S a_n \mid b_1 \mid b_2 \mid \dots \mid b_m$

All strings starting with one of $b_1..b_m$ and continue with several instances of $a_1..a_n$
 Rewrite as

$S \rightarrow b_1 S' \mid b_2 S' \mid \dots \mid b_m S'$
 $S' \rightarrow a_1 S' \mid a_2 S' \mid \dots \mid a_n S' \mid \epsilon$

There are more general ways of encoding left recursion in a grammar

S \rightarrow A a | d
A \rightarrow S b

is also left-recursive because

S \rightarrow S b a

This left recursion can also be eliminated (see Dragon book for general algorithm).

Recursive descent

- Simple and general parsing strategy
- Left-recursion must be eliminated first
 - This can be done automatically
 - In practice however, this is done manually. The reason is that we also need to specify semantic actions with the productions used. Hence, people do elimination of left-recursion on their own, and this is not difficult to do.
- Popular strategy in production compilers. e.g., gcc's parser is a hand-written recursive-descent parser.

Predictive Parsing

- Like recursive-descent but parser can "predict" which production to use
 - by looking at the next few tokens
 - without backtracking
- Predictive parsing accepts LL(k) grammars
 - Left-to-right scanning of tokens
 - Leftmost derivation
 - k tokens to lookahead. In practice, k = 1
- In LL(1), at every step there is at most one choice for a possible production
 - At each step, only one choice of production

wAb \rightarrow w α b , next input: token t//only one choice for A at every step

Example

E \rightarrow T + E | T
T \rightarrow int | int * T | (E)

Here, with one lookahead, hard to predict because:

- we cannot choose the first two productions of T by one lookahead
- the same problem exists with E, because two productions start with non-terminal T. Here it is a non-terminal, but still there is an issue, because we do not know which production to use for one lookahead

Left factoring a grammar

Basic idea: Factor out common prefix into a single production

E \rightarrow TX
X \rightarrow + E | ϵ

Left factoring delays the decision on which production to use. Here we decide which production to use *after* we have seen T.

Similarly,

T \rightarrow int Y | (E)
Y \rightarrow ϵ | * T

Left-factored grammar

E \rightarrow TX X \rightarrow +E | ϵ T \rightarrow int Y | (E) Y \rightarrow ϵ | * T

LL(1) parsing table (assume it is somehow given, later we will discuss how to construct this):

	int	*	+	()	\$
E	TX			TX		
X			+E		ϵ	ϵ
T	int Y			(E)		
Y		*T	ϵ	ϵ	ϵ	ϵ

Empty entries represent error states

Algorithm: similar to recursive descent, except

- For the left-most non-terminal S , we look at the next input token a
- Choose the production rule shown at (S, a)

Instead of using recursive functions, maintain a stack of the frontier of the parse tree. The stack contains all the entries on the right of the top of the frontier. In other words, the top of the stack is the leftmost pending terminal or non-terminal. Terminals in the stack are yet to be matched in the input. Non-terminals in the stack are yet to be expanded.

Reject on reaching error state.

Accept on end of input or empty stack.

Pseudo-code:

```
initialize stack = <S $> and next
repeat
  case stack of
    <X, rest> : if T[X, *next] = Y1...Yn
                then stack <-- <Y1...Yn, rest>;
                else error();
    <t, rest> : if t == *next++
                then stack <-- <rest>;
                else error();
until stack == < >
```

Run this algorithm on the example

Constructing LL(1) Parsing Tables

Constructing First sets

Consider a non-terminal A , production $A \rightarrow \alpha$, and a token t

$T(A, t) = \alpha$ in two cases:

1. If $\alpha \rightarrow^* t \beta$
 - If α can derive t in the first position
 - We say that $t \in \text{First}(\alpha)$
2. If $\alpha \rightarrow^* \epsilon$ and $S \rightarrow^* \beta A \delta$
 - Useful if stack has A , input is t , and A cannot derive t
 - In this case, the only option is to get rid of A by deriving ϵ
 - Can work only if t can follow A in at least one derivation
 - We say that $t \in \text{Follow}(\alpha)$

Definition:

$\text{First}(X) = \{t \mid X \rightarrow^* t \alpha\} \cup \{\epsilon \mid X \rightarrow^* \epsilon\}$

Algorithm sketch:

1. $\text{First}(t) = \{t\}$
2. $\epsilon \in \text{First}(X)$
 - if $X \rightarrow \epsilon$
 - if $X \rightarrow A_1 A_2 \dots A_n$ AND $\epsilon \in \text{First}(A_i)$ for all i
3. $\text{First}(\alpha) \subseteq \text{First}(X)$ if $X \rightarrow A_1 \dots A_n \alpha$ and $\epsilon \in \text{First}(A_i)$ for all i

Example

```
E --> TX
X --> +E | \epsilon
T --> int Y | (E)
Y --> \epsilon | * T
```

First set for terminals

```
First(t) = {t}
First(*) = {*}
....
```

First set for non-terminals

```

First(E) \superset First(T)
First(T) = {'(', 'int')}
First(X) = {+, \epsilon}
First(Y) = {*, \epsilon}

```

Constructing Follow sets

Definition

$\text{Follow}(X) = \{t \mid S \xrightarrow{*} \beta X t \text{ } \delta\}$

Intuition

- if $X \rightarrow A B$, then $\text{First}(B) \subseteq \text{Follow}(A)$ and $\text{Follow}(X) \subseteq \text{Follow}(B)$
- if $B \rightarrow \epsilon$, $\text{Follow}(X) \subseteq \text{Follow}(A)$
- If S is the start symbol, then $\$ \in \text{Follow}(S)$

Algorithm sketch:

1. $\$ \in \text{Follow}(S)$
2. For each $A \rightarrow \alpha X \beta$,
 - $\text{First}(\beta) - \epsilon \subseteq \text{Follow}(X)$
3. For each $A \rightarrow \alpha X \beta$
 - $\text{Follow}(A) \subseteq \text{Follow}(X)$ if $\epsilon \in \text{First}(\beta)$

Example

```

E --> TX
X --> +E | \epsilon
T --> int Y | (E)
Y --> \epsilon | * T

```

Running the algo

```

Follow(E) = { $, ) }
Follow(E) \superset Follow(X)
Follow(X) \superset Follow(E)
(in other words Follow(X) = Follow(E))

```

Hence

$\text{Follow}(X) = \{ \$,) \}$

Looking at T, and the first production $E \rightarrow TX$

```

Follow(T) \superset First(X)
or
Follow(T) = { + } (ignoring \epsilon)
Follow(T) \superset Follow(E) because X --> * \epsilon
Follow(T) = { +, $ }

```

Looking at the second rule,

```

Follow(Y) \superset Follow(T)
Follow(T) \superset Follow(Y)

```

$\text{Follow}(Y) = \{ +, \$,) \}$

Looking at terminals

```

Follow('(') \superset First(E)
Follow('(') = { (, int }
Follow(')') \superset Follow(T)
Follow(')') = { +, $, ) }
Follow(+) \superset First(E)
Follow(+) = { (, int }
E does not go to \epsilon, so this is it.
Follow(*) \superset First(T)
Follow(*) = { (, int }
T does not go to \epsilon, so this is it.
Follow(int) \superset First(Y) - \epsilon
Follow(int) = { * }
Follow(int) \superset Follow(T) (because \epsilon \in First(Y))
Follow(int) = { *, +, $, ) }

```

LL(1) Parsing Tables

Goal: Construct a parsing table T for CFG G

- For each production $A \rightarrow \alpha$ in G , do:
 - For each terminal $t \in \text{First}(\alpha)$, do $T[A, t] = \alpha$
- If $\epsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do $T[A, t] = \alpha$ (for token t)
- If $\epsilon \in \text{First}(\alpha)$, for each $\$ \in \text{Follow}(A)$ do $T[A, \$] = \alpha$ (just a special case for $\$$)

Parsing table

	int	*	+	()	\$
E	TX			TX		
X			+E		ϵ	ϵ
T	int Y			(E)		
Y		*T	ϵ	ϵ	ϵ	ϵ

Fill some entries in this table

What happens if we try and build a parsing table for a grammar that is not LL(1). Example: $S \rightarrow Sa \mid b$

$\text{First}(S) = \{b\}$

$\text{Follow}(S) = \{\$, a\}$

In this case, both productions will appear at $T[S, b]$. If an entry is multiply defined (more than one productions appear in a table cell), G is not LL(1). Examples: a grammar that is ambiguous, not left-factored, left-recursive will all be not LL(1). But there are more grammars that may not be LL(1)

Most programming language CFGs are not LL(1). More powerful formulations for describing practical grammars that assemble some of these ideas in more sophisticated ways to develop more parsers