

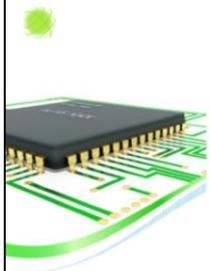
Definition

In a typical language processor implementation, syntactic analysis is the attempt to recognize string of tokens as a sentence in a given language.

Parser is the general term for the component that performs the syntactic analysis.

In a typical language processor implementation, syntactic analysis is the attempt to recognize string of tokens as a sentence in a given language. Parser is the general term for the component that performs the syntactic analysis.

Typical implementations receive the sequence of the tokens generated by the lexical analyzer. In such cases, the sequence of tokens is fixed for a given input. There may be other approaches, especially in parsing strategies using top-down processing, where the tokens are recognized as the parsing proceeds. For such cases, back-tracking and alternative sequence of tokens may be observed.



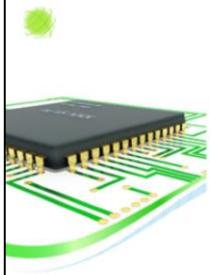
Syntactic Validity

Yes, valid program is a sentence.

Usually, validity transcends syntax.

- **Declarative Integrity**
- **Parameter Matching**
- **Type Conformance**
- **Operator Applicability**
- **Statement Applicability**
- ...

Thinking of our experiences with the programming languages which we are accustomed to, even without consulting to the programming language definitions we can safely state that syntactic compliance is necessary to validate a program but not enough. Validity of a given program is governed by the rules that transcend the grammatical constraints. Can you give examples of context dependent properties of a programming language that you know?



Grammars

CFG and (E)BNF

- Context Free Grammar (Type 2: Grammars)

$$G = (V, \Sigma, R, S)$$

V: Set of non-terminal (Variables)

Σ : Set of terminals (Reported by lexical analyzer)

R: Rule set (Production rules, productions)

S: Start variable

- Backus-Naur Form (Notation)

$$\begin{array}{ll} rexpr & \rightarrow rexpr + rterm \mid rterm \\ rterm & \rightarrow rterm rfactor \mid rfactor \\ rfactor & \rightarrow rfactor * \mid rprimary \\ rprimary & \rightarrow a \mid b \end{array}$$

It is a common practice to use formal notations to describe the syntactic rules of languages for a bunch of reasons such as precise-expressive specification, ease of maintenance and improvement, automatic generation recognizers, and similar. For most of the languages that need a parser, the syntactic specification has a property of recursion which is not supported by the type-3 languages i.e. regular expressions.

Any BNF like notation defines common constructs as noted below.

Non-terminals: These are the symbols for the atomic components of the syntax definition reported by the lexical analyzer.

Terminals: These are the symbols making reference to the derivation rules. This references allow the designer to define recursion / cycles in a controlled manner.

Rules: These are the constructs combining the symbols with simple operators such as alternation, Kleene star and derivatives, and helper elements such as Epsilon, Parentheses, and similar. In some contexts, Kleene star and its derivatives are avoided for sake of simplicity as they can be mapped to expressions with help of interim symbols, alternation and Epsilon without loss of expressivity.



Grammars

BNF for Human Reader

The syntax of the standard, in RFC #733, was originally specified in the Backus-Naur Form (BNF) meta-language. Ken L. Harrenstien, of SRI International, was responsible for re-coding the BNF into an augmented BNF that makes the representation smaller and easier to understand.

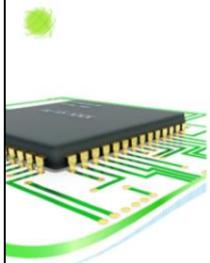
RFC 822

```
....  
dtext      = <any CHAR excluding "[",      ; => may be folded  
          "]", "\", & CR, & including  
          linear-white-space>  
  
comment    = "(" *(ctext / quoted-pair / comment) ")"  
  
ctext      = <any CHAR excluding "(",      ; => may be folded  
          ")", "\", & CR, & including  
          linear-white-space>  
  
quoted-pair = "\" CHAR                      ; may quote any char  
  
phrase     = 1*word                         ; Sequence of words  
  
word       = atom / quoted-string
```

The excerpt is from RFC 822, "Standard for the Format of ARPA Internet Text Messages; David H. Crocker,"
<https://datatracker.ietf.org/doc/html/rfc822>

A common notation to specify the intended syntax is known as Backus-Naur Form (BNF), which dates back to the 1960s. BNF and similar notations are widely applied to define the grammar of the intended language as type-2 grammars which are called as context free grammars (CFGs). The notations used in books or standard documents (RFC's for example) have a certain level of flexibility as their main concern is to describe the grammar for human readers. Automatic parser generators on the other hand must be strict in their meta-language (notation) definition, which is similar to the BNF notation at some level. Each project defines its own way of building declarations.

The excerpt is from RFC 822, "Standard for the Format of ARPA Internet Text Messages; David H. Crocker;", <https://datatracker.ietf.org/doc/html/rfc822>



Grammars

(E)BNF for a Language Processor

From a Parser Generator Project

```
*****\n    Lexical elements and rules referenced by both RFC822 and MIME.\n*****\nmsg_id      : Matches(<:LessThan) addr_spec Matches(>:GreaterThan);\naddr_spec   : local_part AtSign domain;\nlocal_part   : word *(Dot word);\nword         : Matches(+!{#0-#32,specials} :atom) |\n                Matches(" *(!{#0,"\\\",#13}|\\{#1-#255}) ":"quoted_string");\ndomain       : sub_domain *(Dot sub_domain);\nsub_domain   : domain_ref | Matches([*(!{#0,[,],\\\",#13}|\\{#1-#255}):domain_literal]);\ndomain_ref   : atom;
```

The excerpt is from a proprietary parser definition language.



Grammars

BNF Simplified

* Transformation

$u e^* w \rightarrow u e' w$
 $e' : e'e | \epsilon$ (Left recursive form)
 $e' : e e' | \epsilon$ (Right recursive form)

+ Transformation

$u e^+ w \rightarrow u e e' w$
 $e' : e'e | \epsilon$ (Left recursive form)
 $e' : e e' | \epsilon$ (Right recursive form)

? Transformation

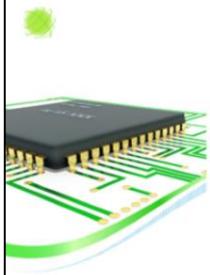
$u e? w \rightarrow u e' w$
 $e' : e | \epsilon$

Parentheses

$u (e) w \rightarrow u e' w$
 $e' : e$

Alternative Notations

[e], {e}, prefix, postfix, etc...



Derivation

Definition

Derivation is a sequence of rewriting steps that begins with the grammar's start symbol and ends with a sentence in the language.

Definition from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

If uAv is a sentential form in derivation of string and production rule A exists such that $A \rightarrow w$, then $uAv \rightarrow uwv . A \hat{\mid} V$, and $u, v, w \in S^*$, where S^* is the set of the sentential forms allowed by the R .

Basically, the choice of A in a sentential form can be made arbitrarily. But, there are two basic ways of making this choice that we are interested in. The first way is choosing the leftmost non-terminal, the second way is choosing the rightmost non-terminal. If a derivation chooses the leftmost non-terminals only, it is called leftmost derivation (LMD). Similarly, if a derivation chooses the rightmost non-terminals only, it is called rightmost derivation (RMD). A sentential form can be defined alternatively as an intermediate form that is reachable from S . So, a left sentential form is a sentential form that exists as a step in LMD; a right sentential form is a sentential form that exists as a step in RMD.



Sentential Form

Definition

Sentential form is a string of symbols that occurs as a step in a valid derivation.

A step is making use of a rule on the latest sentential form.

- Choose a variable.
- Replace with the associated production.

$$G = (V, \Sigma, R, S)$$

$$f_n = uAv \quad A \in V \quad A \rightarrow w \in R \quad f_{n+1} = uwv \quad uAv, uwv \in S^*$$

$$u, v \in (V \cup \Sigma)^*$$

Definition from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

If uAv is a sentential form in derivation of string and production rule A exists such that $A \rightarrow w$, then $uAv \rightarrow uwv$. $A \in V$, and $u, v, w \in S^*$, where S^* is the set of the sentential forms allowed by the R .

Basically, the choice of A in a sentential form can be made arbitrarily.



Derivation

$S^* \Rightarrow \alpha$ means α can be found as a sentence starting from S by zero or more derivation steps.

$\alpha \Rightarrow^* \beta$ means β is reachable from α by zero or more derivation steps.

$\alpha \Rightarrow^* \beta$ means β is reachable from α by one or more derivation steps.

$\alpha \Rightarrow^* \beta$ and $\beta \Rightarrow^* \gamma$ then $\alpha \Rightarrow^* \gamma$

G1 :

Example: Input: aaba

$S \rightarrow aA$

$S \rightarrow abC$

$A \rightarrow aA$

$A \rightarrow bA$

$A \rightarrow C$

$A \rightarrow \epsilon$

$C \rightarrow c$

G1 :

$S \rightarrow aA \mid abC$

$A \rightarrow aA \mid bA \mid C \mid \epsilon$

$C \rightarrow c$

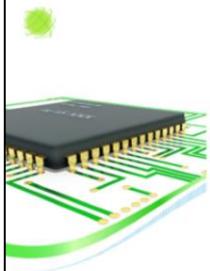
S by $S \rightarrow aA$

aA by $A \rightarrow aA$

aaA by $A \rightarrow bA$

$aabA \rightarrow aA$

$aaba$ This is a sentence. This is rightmost derivation by the way.



Derivations

LMD, RMD, Left and Right Sentential Forms

Leftmost Derivation (LMD) chooses the leftmost variable.

Rightmost Derivation (RMD) chooses the rightmost variable.

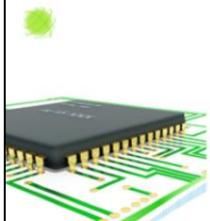
Any step in LMD is leftmost sentential form.

Any step in RMD is rightmost sentential form.

Finding sentence aaba by G1

S use $S \rightarrow aA$
aA use $A \rightarrow aA$
aaA use $A \rightarrow bA$
aabA use $A \rightarrow aA$
aabaA use $A \rightarrow \epsilon$
aaba

But, there are two basic ways of making this choice that we are interested in. The first way is choosing the leftmost non-terminal, the second way is choosing the rightmost non-terminal. If a derivation chooses the leftmost non-terminals only, it is called leftmost derivation (LMD). Similarly, if a derivation chooses the rightmost non-terminals only, it is called rightmost derivation (RMD). A sentential form can be defined alternatively as an intermediate form that is reachable from S. So, a left sentential form is a sentential form that exists as a step in LMD; a right sentential form is a sentential form that exists as a step in RMD.



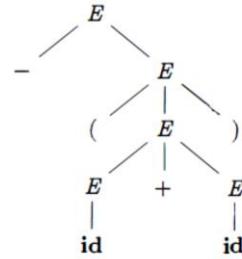
Parse Trees

Successive steps taken in a derivation can be represented with help of trees, which are known as parse trees.

G: $E \rightarrow E + E \mid E * E \mid - E \mid (E) \mid id$

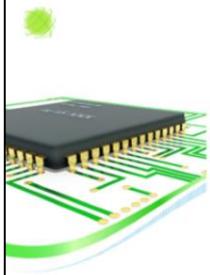
$E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow - E$
 $E \rightarrow (E)$
 $E \rightarrow id$

Input: $- (id + id)$



Example from "Aho, A.V, Ullman J.D, Sethi R., Lam M.S; Dragon Book"

Successive steps taken in a derivation can be represented with help of trees, which are known as parse trees.



Ambiguity

A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

G: $E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$

$E \rightarrow E + E$
 $E \rightarrow E * E$
 $E \rightarrow -E$
 $E \rightarrow (E)$
 $E \rightarrow id$

Input: $id + id * id$

**Create leftmost derivations
for this sentence!**

Definition and example from "Aho, A.V., Ullman J.D., Sethi R., Lam M.S; Dragon Book"

A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

A well-known ambiguity problem is the dangling-else problem!



Ambiguity Elimination

Rewriting

G:

$E \rightarrow E + E \mid E * E \mid -E \mid (E) \mid id$

$E \rightarrow E + E$

$E \rightarrow E * E$

$E \rightarrow -E$

$E \rightarrow (E)$

$E \rightarrow id$

G:

$E \rightarrow E + T \mid T$

$T \rightarrow T * F \mid F$

$F \rightarrow -E \mid (E) \mid id$

$E \rightarrow E + T$

$E \rightarrow T$

$T \rightarrow T * F$

$T \rightarrow F$

$F \rightarrow -E$

$F \rightarrow (E)$

$F \rightarrow id$

Inputs: $id + id * id$, $id + id + id$, $id * id * id$

Definition and example from "Aho, A.V, Ullman J.D, Sethi R., Lam M.S; Dragon Book"

A grammar that produces more than one parse tree for some sentence is said to be ambiguous.



Parsing

as Tree Building Method

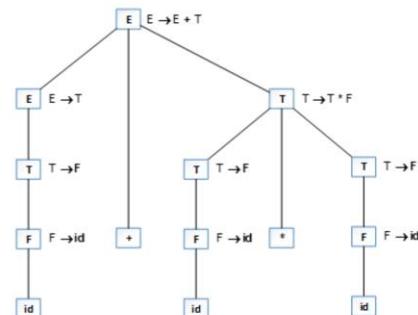
Top-Down Parsing

Begins with the root and grows the tree toward the leaves.

Bottom-Up Parsing

Begins with the leaves and grows the tree toward the root.

Input: id + id * id



Definition from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

A grammar that produces more than one parse tree for some sentence is said to be ambiguous.

Create the parse tree top-down and bottom-up fashion.



Top Down Parsing

The Left Recursion Problem

G:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow -E \mid (E) \mid id \end{aligned}$$

Input: id + id * id

Immediate Recursions

$$\begin{aligned} E &\rightarrow E + T \\ E &\rightarrow T \\ T &\rightarrow T * F \\ T &\rightarrow F \\ F &\rightarrow -E \\ F &\rightarrow (E) \\ F &\rightarrow id \end{aligned}$$

$$A \rightarrow A \alpha \mid \beta$$

Elimination

$$\begin{aligned} A &\rightarrow \beta A' \\ A' &\rightarrow \alpha A' \mid \epsilon \end{aligned}$$

Elimination rule from "Aho, A.V, Ullman J.D, Sethi R., Lam M.S; Dragon Book"

A top down parser establishes the parse tree starting from the root and tries each production of the variable by expanding it as a child.

Use leftmost derivations to create the parse tree. When you apply the rules in the order they were declared, you cannot avoid endless recursion.



Left Recursion Elimination

When Grammar has no Cycles and no ϵ -productions

General Elimination Algorithm

arrange the nonterminals in some order A_1, A_2, \dots, A_n .

```
for ( each  $i$  from 1 to  $n$  ) {
    for ( each  $j$  from 1 to  $i - 1$  ) {
        replace each production of the form  $A_i \rightarrow A_j \gamma$  by the
        productions  $A_i \rightarrow \delta_1 \gamma \mid \delta_2 \gamma \mid \dots \mid \delta_k \gamma$ , where
         $A_j \rightarrow \delta_1 \mid \delta_2 \mid \dots \mid \delta_k$  are all current  $A_j$ -productions
    }
    eliminate the immediate left recursion among the  $A_i$ -productions
}
```

Algorithm from "Aho, A.V, Ullman J.D, Sethi R., Lam M.S; Dragon Book"

The algorithm given above does not guarantee when the preconditions are not met.



Left Recursion Elimination

General Elimination Algorithm

G:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow -E \mid (E) \mid id \end{aligned}$$

$$\begin{aligned} E &\rightarrow E + T \\ E &\rightarrow T \\ T &\rightarrow T * F \\ T &\rightarrow F \\ F &\rightarrow -E \\ F &\rightarrow (E) \\ F &\rightarrow id \end{aligned}$$

Input: id + id * id

G:

$$\begin{aligned} E &\rightarrow E + T \mid T \\ T &\rightarrow T * F \mid F \\ F &\rightarrow -E \mid (E) \mid id \end{aligned}$$

$$\begin{aligned} E &\rightarrow T E' \\ E' &\rightarrow +T E' \mid \epsilon \\ T &\rightarrow F T' \\ T' &\rightarrow * F T' \mid \epsilon \\ F &\rightarrow -E \\ F &\rightarrow (E) \\ F &\rightarrow id \end{aligned}$$

When the algorithm applied to the grammar G, it is transformed to a more complicated form. But now, it is recursion free.

Use leftmost derivations again to create the parse tree.



Backtracking

The
leftmost,
top-down
parser

```
root ← node for the start symbol, S;  
focus ← root;  
push(null);  
word ← NextWord();  
  
while (true) do;  
    if (focus is a nonterminal) then begin;  
        pick next rule to expand focus ( $A \rightarrow \beta_1, \beta_2, \dots, \beta_n$ );  
        build nodes for  $\beta_1, \beta_2, \dots, \beta_n$  as children of focus;  
        push( $\beta_n, \beta_{n-1}, \dots, \beta_2$ );  
        focus ←  $\beta_1$ ;  
    end;  
    else if (word matches focus) then begin;  
        word ← NextWord();  
        focus ← pop();  
    end;  
    else if (word = eof and focus = null)  
        then accept the input and return root;  
    else backtrack;  
end;
```

Algorithm from "Cooper, K.D., Torczon, L.; Engineering A Compiler"



Backtracking

Backtrack-free Grammar / Predictive Grammar
a CFG for which the leftmost, top-down parser can always predict the correct rule with lookahead of at most one word.

Elimination Using First and Follow Sets

FIRST Set

For a grammar symbol α ,
 $\text{FIRST}(\alpha)$ is the set of terminals that can appear at the start of a sentence derived from α .

FOLLOW Set

For a nonterminal α , $\text{FOLLOW}(\alpha)$ contains the set of words that can occur immediately after α in a sentence.

Definitions from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

First Set Construction



```

for each  $\alpha \in (T \cup eof \cup \epsilon)$  do;
    FIRST( $\alpha$ )  $\leftarrow \alpha$ ;
end;
for each  $A \in NT$  do;
    FIRST( $A$ )  $\leftarrow \emptyset$ ;
end;

while (FIRST sets are still changing) do;
    for each  $p \in P$ , where  $p$  has the form  $A \rightarrow \beta$  do;
        if  $\beta$  is  $\beta_1\beta_2\dots\beta_k$ , where  $\beta_i \in T \cup NT$ , then begin;
            rhs  $\leftarrow$  FIRST( $\beta_1$ )  $- \{\epsilon\}$ ;
            i  $\leftarrow 1$ ;
            while ( $\epsilon \in$  FIRST( $\beta_i$ ) and  $i \leq k-1$ ) do;
                rhs  $\leftarrow$  rhs  $\cup$  (FIRST( $\beta_{i+1}$ )  $- \{\epsilon\}$ );
                i  $\leftarrow i + 1$ ;
            end;
        end;
        if  $i = k$  and  $\epsilon \in$  FIRST( $\beta_k$ )
            then rhs  $\leftarrow$  rhs  $\cup \{\epsilon\}$ ;
        FIRST( $A$ )  $\leftarrow$  FIRST( $A$ )  $\cup$  rhs;
    end;
end;

```

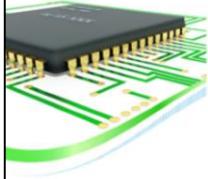
G = (V, Σ, R, S)

What are T, NT, P?

Algorithm from "Cooper, K.D., Torczon, L.; Engineering A Compiler"



First Set Example



G:

$E \rightarrow E + T \cdot T$
 $T \rightarrow T^* F | F$
 $F \rightarrow -E | (E) | id$

$E \rightarrow T E'$
 $E' \rightarrow +T E'$
 $E' \rightarrow \epsilon$

Build the first sets!

$T \rightarrow F T'$
 $T' \rightarrow * F T'$
 $T' \rightarrow \epsilon$

$F \rightarrow -E$
 $F \rightarrow (E)$
 $F \rightarrow id$

Build the first sets.



Follow Set Construction

```
for each  $A \in NT$  do;
    FOLLOW( $A$ )  $\leftarrow \emptyset$ ;
end;

FOLLOW( $S$ )  $\leftarrow \{\text{eof}\}$ ;
while (FOLLOW sets are still changing) do;
    for each  $p \in P$  of the form  $A \rightarrow \beta_1\beta_2\dots\beta_k$  do;
        TRAILER  $\leftarrow$  FOLLOW( $A$ );
        for  $i \leftarrow k$  down to 1 do;
            if  $\beta_i \in NT$  then begin;
                FOLLOW( $\beta_i$ )  $\leftarrow$  FOLLOW( $\beta_i$ )  $\cup$  TRAILER;
                if  $\epsilon \in \text{FIRST}(\beta_i)$ 
                    then TRAILER  $\leftarrow$  TRAILER  $\cup$  ( $\text{FIRST}(\beta_i) - \epsilon$ );
                else TRAILER  $\leftarrow$  FIRST( $\beta_i$ );
            end;
            else TRAILER  $\leftarrow$  FIRST( $\beta_i$ ); // is { $\beta_i$ }
        end;
    end;
end;
```

Algorithm from "Cooper, K.D., Torczon, L.; Engineering A Compiler"



Follow Set Example

G:

$E \rightarrow E + T \mid T$
 $T \rightarrow T^* F \mid F$
 $F \rightarrow -E \mid (E) \mid id$

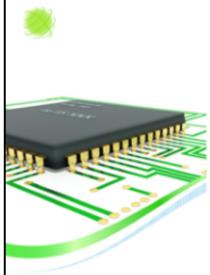
$E \rightarrow T E'$
 $E' \rightarrow +T E'$
 $E' \rightarrow \epsilon$

Build the follow sets!

$T \rightarrow F T'$
 $T' \rightarrow * F T'$
 $T' \rightarrow \epsilon$

$F \rightarrow -E$
 $F \rightarrow (E)$
 $F \rightarrow id$

Build the first sets.



Backtracking

Backtrack-free Test

Test each non-terminal with multiple RHS!

$$First^+(A \rightarrow \beta) = \begin{cases} First(\beta), & \varepsilon \notin First(\beta) \\ First(\beta) \cup Follow(A), & \text{otherwise} \end{cases}$$

$$First^+(A \rightarrow \beta_i) \cap First^+(A \rightarrow \beta_j) = \emptyset, \quad \forall 1 \leq i, j \leq n, \quad i \neq j$$

Test method from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

To test if a grammar is backtrack-free any non-terminal having more than one RHS (alternatives) must be tested. For a backtrack-free grammar, the First+ sets generated for alternatives are distinct.



Left Factoring

G:

$$\begin{aligned} E &\rightarrow E + T \cdot T \\ T &\rightarrow T^* F | F \\ F &\rightarrow -E \mid (E) \mid id \end{aligned}$$

$$\begin{aligned} E &\rightarrow T \cdot E' \\ E' &\rightarrow +T \cdot E' \\ E' &\rightarrow \epsilon \end{aligned}$$

$$\begin{aligned} T &\rightarrow F \cdot T' \\ T' &\rightarrow * F \cdot T' \\ T' &\rightarrow \epsilon \end{aligned}$$

$$\begin{aligned} F &\rightarrow -E \\ F &\rightarrow (E) \\ F &\rightarrow id \\ F &\rightarrow id(E) \\ F &\rightarrow id[E] \end{aligned}$$

When array and function call like structures added!

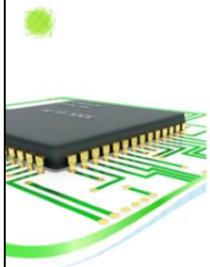
G:

$$\begin{aligned} E &\rightarrow E + T \cdot T \\ T &\rightarrow T^* F | F \\ F &\rightarrow -E \mid (E) \mid id \end{aligned}$$

$$\begin{aligned} &\vdots \\ F &\rightarrow -E \\ F &\rightarrow (E) \\ F &\rightarrow id G \end{aligned}$$

$$\begin{aligned} G &\rightarrow (E) \\ G &\rightarrow [E] \\ G &\rightarrow \epsilon \end{aligned}$$

When we add two more productions for F to create a flavor of function call and array access backtracking test for F fails! Rewriting with left factoring helps eliminate common First+ sets.



Predictive Parsers

Stages of Test & Elimination

- Ambiguity
- Left Recursion
- Backtracking
- Left Factoring

Can you revise the leftmost, top-down parser?

Now, we can focus on

- Recursive Descent Parsers
- Table Driven LL(1) Parser

When we add two more productions for F to create a flavor of function call and array access backtracking test for F fails! Rewriting with left factoring helps eliminate common First+ sets.



Recursive Descent Parsing

```
Expr( )
/* Expr → Term Expr' */
if (Term())
    then return EPrime();
else Fail();

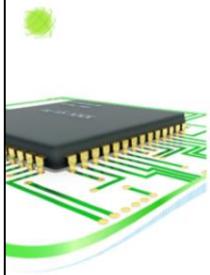
EPrime()
/* Expr' → + Term Expr' */
/* Expr' → - Term Expr' */
if (word = + or word = - )
    then begin;
        word ← NextWord();
        if (Term())
            then return EPrime();
        else Fail();
    end;
```

A procedure for each NT

- Hand coded
- Or, automatically generated based on grammar
- Ease of customization, debugging

Sample pseudo-code from "Cooper, K.D., Torczon, L.; Engineering A Compiler"

Recursive descent parsing is one way of implementing top down parsing strategy. It is easy to construct and follow. It can be built by hand or automatically. A standard debugger becomes a tool to diagnose and fix parsing related problems!



LL(1) Parsing

LL(1)

- 1 lookahead token
- Leftmost derivation
- Left to right scanning

LL(1) Parsers

- Use standard skeleton procedure
- Use parse tables
- Can be generated by parser generators



LL(1) Parsing

**The leftmost,
top-down
parser**

(non-predictive)

```
root ← node for the start symbol, S;
focus ← root;
push(null);
word ← NextWord();

while (true) do:
    if (focus is a nonterminal) then begin:
        pick next rule to expand focus ( $A \rightarrow \beta_1, \beta_2, \dots, \beta_n$ );
        build nodes for  $\beta_1, \beta_2, \dots, \beta_n$  as children of focus;
        push( $\beta_n, \beta_{n-1}, \dots, \beta_2$ );
        focus ←  $\beta_1$ ;
    end;
    else if (word matches focus) then begin:
        word ← NextWord();
        focus ← pop()
    end;
    else if (word = eof and focus = null)
        then accept the input and return root;
        else backtrack;
    end;
```

Sample pseudo-code from "Cooper, K.D., Torczon, L.; Engineering A Compiler"



LL(1) Parser

Predictive
No
backtracking

LL(1) Parsing

```
word ← NextWord();
push eof onto Stack;
push the start symbol, S, onto Stack;
focus ← top of Stack;
loop forever:
    if (focus = eof and word = eof)
        then report success and exit the loop;
    else if (focus| ∈ T or focus = eof) then begin;
        if focus matches word then begin;
            pop Stack;
            word ← NextWord();
        end;
        else report an error looking for symbol at top of stack;
    end;
    else begin; /* focus is a nonterminal */
        if Table[focus,word] is  $A \rightarrow B_1B_2\dots B_k$  then begin;
            pop Stack;
            for i ← k to 1 by -1 do;
                if ( $B_i \neq \epsilon$ )
                    then push  $B_i$  onto Stack;
            end;
            else report an error expanding focus;
        end;
        focus ← top of Stack;
    end;
end;           Sample pseudo-code from "Cooper, K.D., Torczon, L.; Engineering A Compiler"
```



LL(1) Parser

Predictive

No backtracking

0	<i>Goal</i>	\rightarrow	<i>Expr</i>		6	<i>Term'</i>	\rightarrow	x	<i>Factor Term'</i>
1	<i>Expr</i>	\rightarrow	<i>Term Expr'</i>		7			\div	<i>Factor Term'</i>
2	<i>Expr'</i>	\rightarrow	+	<i>Term Expr'</i>	8			ϵ	
3			-	<i>Term Expr'</i>	9	<i>Factor</i>	\rightarrow	()
4			ϵ		10			num	
5	<i>Term</i>	\rightarrow	<i>Factor Term'</i>		11			name	

Use First+ to create the parse table!

	eof	+	-	\times	\div	()	name	num
<i>Goal</i>	—	—	—	—	—	0	—	0	0
<i>Expr</i>	—	—	—	—	—	1	—	1	1
<i>Expr'</i>	4	2	3	—	—	—	4	—	—
<i>Term</i>	—	—	—	—	—	5	—	5	5
<i>Term'</i>	8	8	8	6	7	—	8	—	—
<i>Factor</i>	—	—	—	—	—	9	—	11	10

Grammars and table from "Cooper, K.D., Torczon, L.; Engineering A Compiler"



Bottom Up Parsing

LR Grammars

LR(n)

- n lookahead token
- Rightmost derivation
- Left to right scanning

Bottom-Up Parsing

Begins with the leaves and grows the tree toward the root.

LR Grammars

- Shift Reduce Parsing
- LR(0) Automaton
- Using Follow: SLR Parsing
- Using Lookahead: LR(1) Canonical Items
- Saving Resources: LALR

We will start with shift reduce parsing, which conveys the basic operation of the bottom up parsing.



Bottom Up Parsing

Shift Reduce Parsing

Grammar:

1. $P \rightarrow E$
2. $E \rightarrow E + T$
3. $E \rightarrow T$
4. $T \rightarrow id(E)$
5. $T \rightarrow id$

Input: id(id+id)

Stack	Input	Action
id	id (id + id) \$	shift
id ((id + id) \$	shift
id (id	id + id) \$	shift
id (id +	+ id) \$	reduce $T \rightarrow id$
id (T	+ id) \$	reduce $E \rightarrow T$
id (E	+ id) \$	shift
id (E +	id) \$	shift
id (E + id) \$	reduce $T \rightarrow id$
id (E + T) \$	reduce $E \rightarrow E + T$
id (E)	\$	shift
T	\$	reduce $E \rightarrow T$
E	\$	reduce $P \rightarrow E$
P	\$	accept

Grammar and trace from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

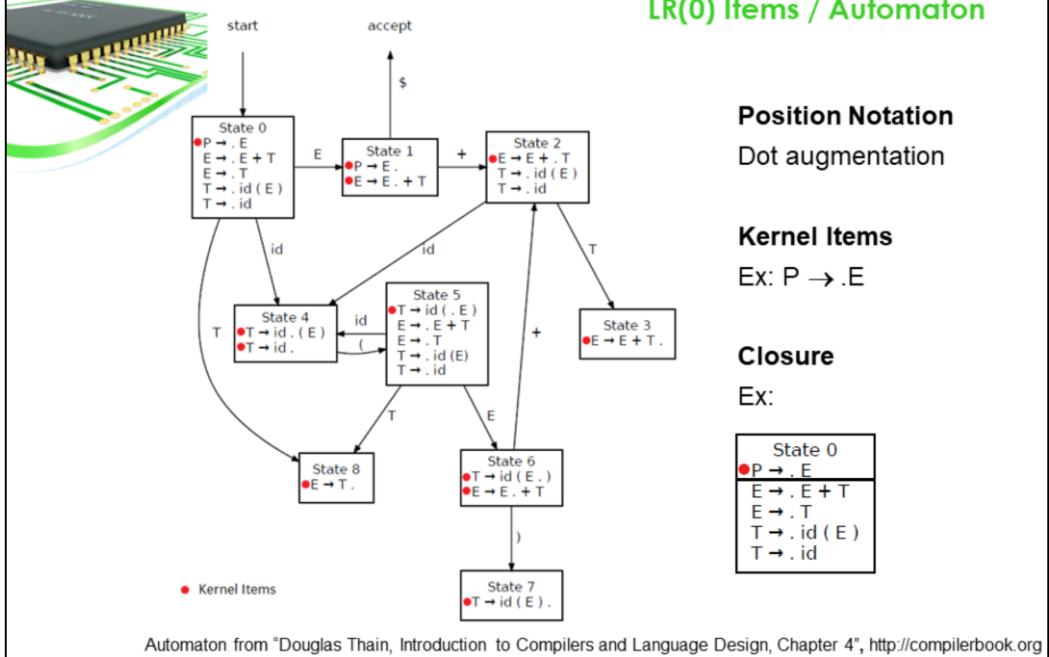
This is a sample to illustrate the operation of the shift reduce parsing, which applies a strategy to start from the input and reaching the goal symbol by applying reductions. The sequence of operations applied is result of purposeful selection to achieve successful completion.

The second action ('shift') could also be ('reduce'), which would not lead us a result. In other words, "id(" is a viable prefix for reduction while "T(", "E(" are not. It is critical to keep the left side in viable prefix forms. Another relevant definition to give at this point is the "handle", which is a substring that matches a production at certain form in right sentential form. See the Stack and Input to observe right sentential forms in reverse order. Handles are reduced (or pruned) in the process of recognition.

A strategy with backtracking would be proposed! Like we did in top-down parsing, we will seek the ways of making good choices with the clues (or the context) that can be obtained from the input stream so that predictive actions can be taken. In this regard, we will define requirements of the well-formed grammars. To start with a rule, let's note that the grammar must not be ambiguous in which case there is more than one handle for some stages of shift reduce parsing.

Bottom Up Parsing

LR(0) Items / Automaton



Position Notation

Dot augmentation

Kernel Items

Ex: $P \rightarrow .E$

Closure

Ex:

State 0
$P \rightarrow .E$
$E \rightarrow .E + T$
$E \rightarrow .T$
$T \rightarrow .id(E)$
$T \rightarrow .id$

The LR(0) items are the productions augmented with a dot that signifies the position on the rule. The substring on the left is the recognized part, the substring on the right are the remaining substring to be recognized.

1. If there is alternatives for the starting production create a root production $S' \rightarrow \Sigma$, ανδ ασσυμε Σ ασ σταρτινγ προδυχτιον.
2. To create start state, create the kernel item using the starting production
3. Choose a state to process and create closure. To create closure of a state, kernel items are the starting points. Scan items in the closure, on each item, for each non-terminal appearing on the right of dot, include the productions for that nonterminal with a dot at start. Repeat the process until the closure becomes stable .i.e. no new items added.
4. Scan items and for each item in the closure, for each w appearing on the right of dot ($A \rightarrow u.wv$) try to find a transition and destination state on the diagram. If not found create a transition and a destination state. Create a kernel item in the form of $A \rightarrow uw.v$.

5. Goto step 3 if there is at least one state that was not processed.



Bottom Up Parsing

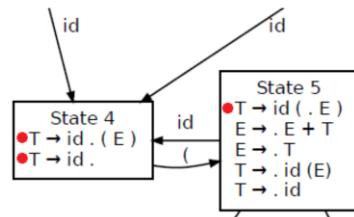
Conflicts

Shift - Reduce Conflicts

From the example grammar (State 4)

$$T \rightarrow id . E$$

$$T \rightarrow id .$$



Reduce - Reduce Conflicts

Ex: If we had followings on a state

$$S \rightarrow id (E) .$$

$$E \rightarrow id (E) .$$

Examples from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

Conflicts have been dealt with by using predictions. Our first step for predictive actions is the SLR Parsing.



Bottom Up Parsing

SLR Parsing

Prediction Using Follow Sets

State	GOTO		ACTION			
	E	T	id	()	+	\$
0	G1	G8	S4			
1				S2	R1	
2		G3	S4			
3				R2	R2	R2
4			S5	R5	R5	R5
5	G6	G8	S4			
6				S7	S2	
7				R4	R4	R4
8				R3	R3	R3

Grammar:

1. $P \rightarrow E$
2. $E \rightarrow E + T$
3. $E \rightarrow T$
4. $T \rightarrow id (E)$
5. $T \rightarrow id$

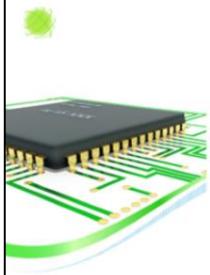
Table and grammar from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

Create the First and Follow Sets.

Create S (Shift) actions for each terminal transition edge using the destination state number.

Create G (Goto) actions for each non-terminal transition edge using the destination state number.

Create R (Reduce) actions for each terminal in the follow set. Use the production number to note R.



Bottom Up Parsing

SLR Parsing

SLR Parsing Algorithm.

Let S be a stack of LR(0) automaton states. Push S_0 onto S .
Let a be the first input token.

Loop:

 Let s be the top of the stack.
 If ACTION[s, a] is **accept**:
 Parse complete.
 Else if ACTION[s, a] is **shift** t :
 Push state t on the stack.
 Let a be the next input token.
 Else if ACTION[s, a] is **reduce** $A \rightarrow \beta$:
 Pop states corresponding to β from the stack.
 Let t be the top of stack.
 Push GOTO[t, A] onto the stack.
 Otherwise:
 Halt with a parse error.

Algorithm from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

This is the SLR parsing algorithm driven by the SLR table and the input.



Bottom Up Parsing

SLR Parsing

	Stack	Symbols	Input	Action
Grammar:				
1. $P \rightarrow E$	0 4	id	id (id + id) \$	shift 4
2. $E \rightarrow E + T$	0 4 5	id ((id + id) \$	shift 5
3. $E \rightarrow T$	0 4 5 4	id (id	+ id) \$	shift 4
4. $T \rightarrow id (E)$	0 4 5 8	id (T	+ id) \$	reduce $T \rightarrow id$
5. $T \rightarrow id$	0 4 5 6	id (E	+ id) \$	reduce $E \rightarrow T$
Input: id(id+id)	0 4 5 6 2	id (E +	id) \$	shift 2
	0 4 5 6 2 4	id (E + id) \$	shift 4
	0 4 5 6 2 3	id (E + T) \$	reduce $T \rightarrow id$
	0 4 5 6 6	id (E)	\$	reduce $E \rightarrow E + T$
	0 8	T	\$	shift 7
	0 1	E	\$	reduce $T \rightarrow id(E)$
				reduce $E \rightarrow T$
				accept

Execution trace from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

Use of SLR table eliminated the dilemma that occurred in the execution that we examined while introducing shift-reduce parsing. Note that, the symbols and input are displayed to make visible the sentential forms. The next token and the stack of the states are the essential parts to implement table driven shift reduce algorithm.

Bottom Up Parsing

LR(1) Items

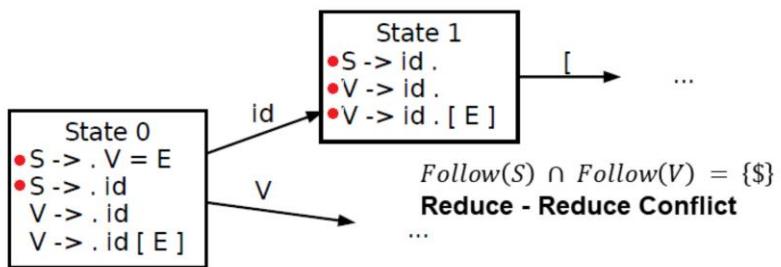


SLR fails when follow sets are not distinct on the same state

	First	Follow
S	id	\$
V	id	= \$)
E	Id	\$)

Grammar:

1. $S \rightarrow V = E$
2. $S \rightarrow id$
3. $V \rightarrow id$
4. $V \rightarrow id (E)$
5. $E \rightarrow V$



Grammar and diagram from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

The state 1 has two productions with follow sets having a common element causing reduce-reduce conflict. Use of lookahead instead of follow may help improve the power of the bottom-up parsing.



Bottom Up Parsing

LR(1) Items

Canonical Form

Like LR(0), Augmented by Lookahead

Grammar:

1. $S \rightarrow V = E$
2. $S \rightarrow id$
3. $V \rightarrow id$
4. $V \rightarrow id (E)$
5. $E \rightarrow V$

Writing Lookahead for Each Item in State

In case $A \rightarrow \alpha . B, \{L\}$ create $B \rightarrow . \gamma, \{L\}$

In case $A \rightarrow \alpha . B \beta, \{L\}$

create $B \rightarrow . \gamma, \text{First}(\beta)$ if β cannot be ϵ

create $B \rightarrow . \gamma, \text{First}(\beta) \cup \{L\}$ if β can be ϵ

Grammar from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

The LR(0) items are the productions augmented with a dot that signifies the position on the rule. The substring on the left is the recognized part, the substring on the right are the remaining substring to be recognized.

1. If there is alternatives for the starting production create a root production $S' \rightarrow \Sigma$, ανδ ασσυμε Σ ασ σταρτινγ προδυχτιον.
2. To create start state, create the kernel item using the starting production(s) with lookahead {\$}
3. Choose a state to process and create closure. To create closure of a state, kernel items are the starting points. Scan items in the closure, on each item, for each non-terminal appearing on the right of dot, include the productions for that nonterminal with a dot at start and the calculated lookahead. Repeat the process until the closure becomes stable .i.e. no new items added.
4. Merge items having LR(0) items with different lookahead by unioning lookahead sets on the finalized closure.
5. Scan items and for each item in the closure, for each w appearing on the

right of dot ($A \rightarrow u.wv, \{L\}$) try to find a transition and destination state on the diagram. If not found create a transition and a destination state. Create a kernel item in the form of $A \rightarrow uw.v \{L\}$.

6. Goto step 3 if there is at least one state that was not processed.



Bottom Up Parsing

LR(1) Items

Example Case for P (State 0)

As kernel item add (1) $P \rightarrow . E, \{\$\}$

Using (1) add (2) $E \rightarrow . E + T, \{\$\}$

Using (1) add (3) $E \rightarrow . T, \{\$\}$

Using (2) add (4) $E \rightarrow . E + T, \{+\}$

Using (3) add (5) $T \rightarrow . id(E), \{\$\}$

Using (3) add (6) $T \rightarrow . id, \{\$\}$

Using (4) add (7) $E \rightarrow . T, \{+\}$

No new item using (5), (6)

Using (7) add (8) $T \rightarrow . id(E), \{+\}$

Using (7) add (9) $T \rightarrow . id, \{+\}$

No new items using (8), (9)

Merge the items having common LR(0) items.

Closure

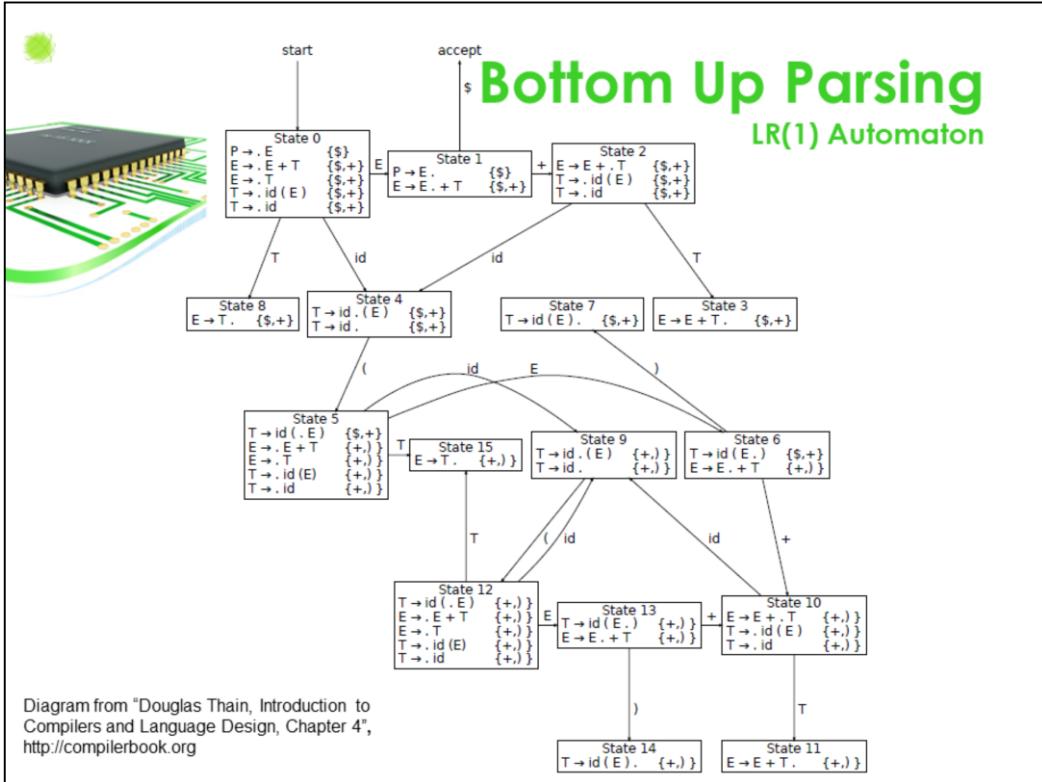
$P \rightarrow . E, \{\$\}$
$E \rightarrow . E + T, \{\$\}$
$E \rightarrow . T, \{\$\}$
$T \rightarrow . id(E), \{\$\}$
$T \rightarrow . id, \{\$\}$

Grammar:

1. $P \rightarrow E$
2. $E \rightarrow E + T$
3. $E \rightarrow T$
4. $T \rightarrow id(E)$
5. $T \rightarrow id$

Grammar and diagram from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

This is step by step evolution of closure for the state 0.



Applying the principles similar to the LR(0) case, the LR(1) automaton is built as seen. Note that there are many states (16 in this case) for even a small piece of grammar. Developing the LR(1) states and items by hand is time consuming and it is easy to make mistakes. The parser generators like Yacc, Bison take care of handling generation of the tables and error free code.



Bottom Up Parsing

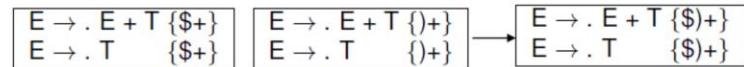
LALR Parsing

Merge States Having Common Cores

Identify the States Having Common Cores

Write a New State with LR(0) items and Union Lookaheads

Example: Not from the Previous Diagram!



Identify the states that can be merged in the previous diagram!

Reconstruct the diagram!

Note on Grammars:

$LL(1) \subset SLR \subset LALR \subset LR(1) \subset CFG$

Example and comparison from "Douglas Thain, Introduction to Compilers and Language Design, Chapter 4", <http://compilerbook.org>

LALR parsers weaker than the LR(1) parsers but stronger than the SLR parsers.



Final Words

Other Parsing Strategies

- LL(*)
- Parse Expression Grammars (PEG)
- ...more

"BOTTOM LINE: We do all this math to ensure that the source code is translated to internal representation unambiguously, predictably." (*)

(*) Remark from Cem Bozşahin, notes on Syntactic Analysis

Studying LL(1), SLR, LALR, LR(1) parsers are useful to understand the insights and the challenges in syntactic analysis. The study helps as relate the syntactic analysis with context free grammars along with the properties and limitations. Note that even syntactic analysis may transcend the boundaries of context free grammars. Even though the parsers are mostly generated with help of parser generators, manual interventions may be needed to obtain the intended parsers. The strategies of implementing parsers are far more richer than the strategies introduced in this presentation.

<https://www.antlr.org/papers/LL-star-PLDI11.pdf>

<https://bford.info/pub/lang/peg.pdf>