

Recursive descent parsing

Each language L over alphabet A has an associated **recognition problem**: Given a finite sequence in A^* , determine whether it is in L .

Many, but not all, context free languages can be recognized using a simple technique called **recursive descent parsing**.

Definition t **is a prefix of** s if and only if there is a u such that $s = tu$.

The idea is this:

- Start with a suitable CFG $(A, N, P, n_{\text{start}})$ for L
- For each nonterminal n in N create a procedure n
- Roughly speaking, the job of procedure n is to try to remove from the input a suitable prefix described by nonterminal n .
- If there is no suitable prefix, the procedure may indicate failure by setting a flag f to false.
- We use variable s to represent the remaining input sequence.
- We'll mark the end of input with a sentinel symbol $\$$ not in $A \cup N$.

Example: (tree is the start nonterminal)
$$\begin{aligned} \text{tree} &\rightarrow [\text{moreTree} \\ \text{tree} &\rightarrow \mathbf{id} \\ \text{moreTree} &\rightarrow] \\ \text{moreTree} &\rightarrow \text{tree moreTree} \end{aligned}$$

Variables:

- f is set to false if an error is encountered
- s is the remaining input. Ends with a \$.
- We assume that $t \in A^*$; so there is no \$ in t .

The main code. Is t in the language?

$$f := \text{true} \quad s := t^{\$} \quad \text{tree}() \quad f := f \wedge (s(0) = \$)$$

Where the procedures are

proc tree() // Try to remove a prefix described by tree.

if $\neg f$ then return end if

*if $s(0) = [$ then *consume()* *moreTree()**

else expect(id) end if

end tree

proc moreTree()

// Try to remove a prefix described by moreTree.

if $\neg f$ then return end if

*if $s(0) =]$ then *consume()**

*else *tree()* *moreTree()* end if*

end moreTree

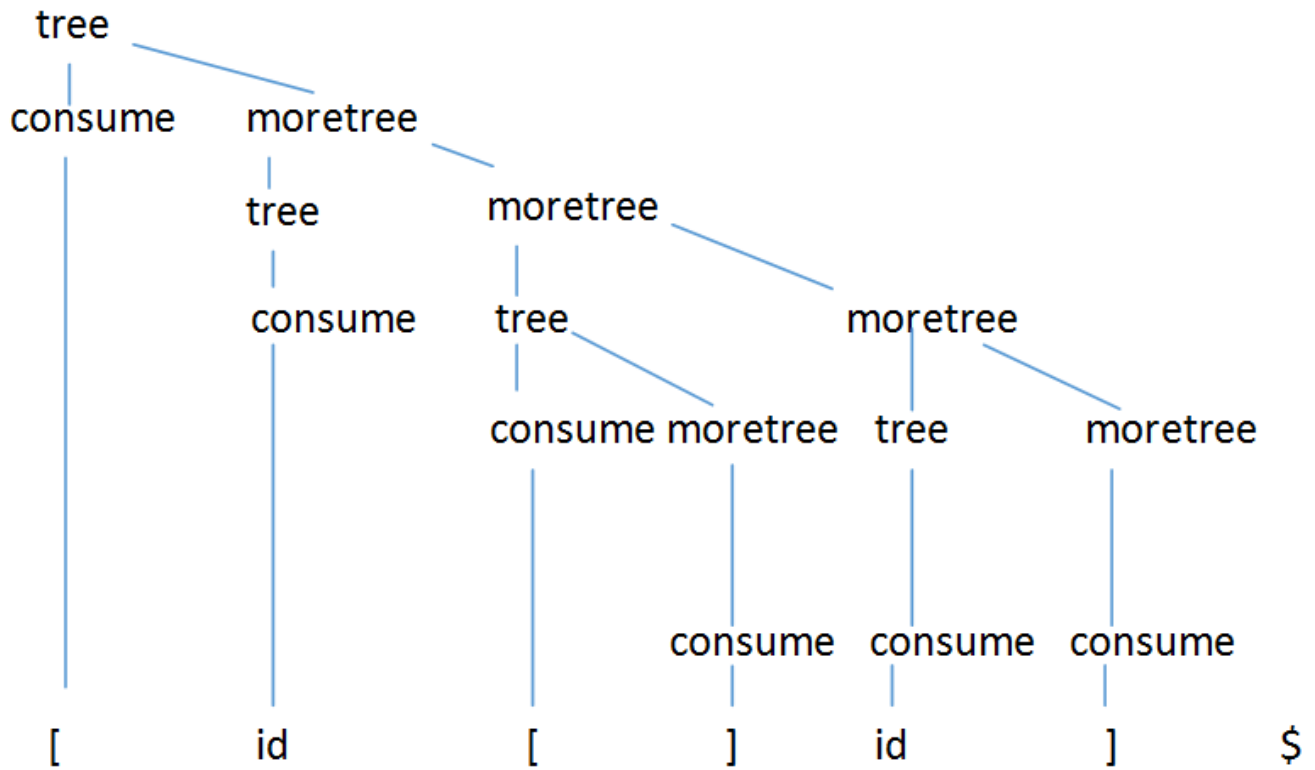
proc consume() $s := s[1, ..s.length]$ end consume

proc expect(a)

*if $s(0) = a$ then *consume()* else $f := \text{false}$ end if*

end expect

Here is an example call tree showing how this work in a successful recognition. Note how the call tree mimics the parse tree.



Specification of nonterminal procedures

The specification for procedures representing nonterminals

procedure $n()$ // Try to remove a prefix described by n

precondition: s is nonempty and ends with a $\$$

changes s, f

postcondition:

There are two possible outcomes

- Error: f is false and s still ends with a $\$$.
- Success: f is true and a prefix of s_0 , described by n , has been removed. I.e., $\exists u. s_0 = us$ and $n \xrightarrow{*} u$.

Choosing an outcome:

- If f_0 is false, Error is the only possible outcome.
- If f_0 is true but no prefix of s_0 is described by n , Error is the only possible outcome.
- If f_0 is true and $\exists t, u, v \in A^*. n_{\text{start}}\$ \xrightarrow{*} tnv\$ \xrightarrow{*} tuv\$$ and $uv\$ = s_0$, then Success is the only possible outcome (and the prefix u removed should meet these conditions).
- Otherwise it doesn't matter which outcome is chosen.

Now assume the initial value of $s \in A^*$. We can tell if s is in L as follows

$$f := \text{true}; \quad s := s^{\tau} \hat{\ } [\$]; \quad n_{\text{start}}(); \quad f := f \wedge (s(0) = \$)$$

Some handy procedures

procedure *expect*($a : A$)

// Try to remove a from the start of the s .

precondition: s ends with contains a $\$$

changes s , f

postcondition:

if f_0 and $[a]$ is a prefix of s_0)

then f and $s_0 = [a] \hat{\ } s$

else $\neg f$ and s ends with contains a $\$$

if $s(0) = a$ then *consume*()

else $f := \text{false}$ end if

end *expect*

procedure *consume*()

// Remove the first item from s

precondition: $s.\text{length} > 0$ and $s(0) \in A$

changes s

postcondition: $s = s_0[1, ..s_0.\text{length}]$

$s := s[1, ..s.\text{length}]$

end *consume*

Writing procedures that meet the specification

If a nonterminal n has productions

$$(n \rightarrow \alpha), (n \rightarrow \beta), (n \rightarrow \gamma) \in P,$$

we write a subroutine like this:

```

procedure  $n()$ 
  // For specification see slide 4
  if  $\neg f$  then return end if
  if ? then  $\llbracket \alpha \rrbracket$ 
  else if ? then  $\llbracket \beta \rrbracket$ 
  else if ? then  $\llbracket \gamma \rrbracket$ 
  else  $f := \text{false}$  end if
end  $n$ 

```

where, for $a \in A$, $m \in N$, $\alpha, \beta \in (A \cup N)^*$

$$\llbracket a \rrbracket = \text{“expect}(a)\text{”}$$

$$\llbracket m \rrbracket = \text{“}m()\text{”}$$

$$\llbracket \epsilon \rrbracket = \epsilon$$

$$\llbracket \alpha\beta \rrbracket = \llbracket \alpha \rrbracket \wedge \llbracket \beta \rrbracket$$

- Usually the boolean expressions are based on the first few items of s .
- The last case $f := \text{false}$ might be unreachable; in this case it is omitted.
- Note that $\{\neg f\} \llbracket \alpha \rrbracket \{\neg f\}$ is correct

Parsing our programming language

var $s : A^*$.

var $f : \mathbb{B}$.

procedure *main*()

 read the input into t , combining characters into symbols
 and throwing out comments and spaces

$f := \text{true}$

$s := t^{\$}$

block()

$f := f \wedge (s(0) = \$)$

 { $f = (t \text{ is in the programming language})$ }

 if f then print “yep” else print “nope” end if

end *main*

Nonterminal block

$$\text{block} \rightarrow \epsilon$$

$$\text{block} \rightarrow \text{command block}$$

procedure *block()* // Version 0

// Try to remove a prefix described by *block* .

// See the contract on slide 4

if $\neg f$ then return end if

if $s(0) \in \text{FirstComm}$ then

command() *block()*

end if

end *block*

where *FirstComm* is $\{\text{if, while}\} \cup \mathcal{I}$.

Why this works:

- When $\text{block} \rightarrow \text{command block}$ is appropriate, $s(0)$ is in $\{\text{if, while}\} \cup \mathcal{I}$;
 - * you can see this by looking at all the productions for *command*.
- When $\text{block} \rightarrow \epsilon$ is appropriate, $s(0) \in \{\$, \text{end, else}\}$;
 - * you can see this by looking at all the places *block* is used in the grammar;
 - * thus $s(0)$ is not in $\{\text{if, while}\} \cup \mathcal{I}$.
- Thus it is never right to pick the $\text{block} \rightarrow \epsilon$ production when $s(0)$ is in *FirstComm*

Note that we can apply tail recursion removal, if we want.

```

procedure block() // Version 1
  // Try to remove a prefix described by block .
  // See the contract on slide 4
  while  $f \wedge s(0) \in \{\mathbf{if}, \mathbf{while}\} \cup \mathcal{I}$  do
    command()
  end while
end block

```

Also acceptable would be

```

procedure block() // Version 2
  // Try to remove a prefix described by block .
  // See the contract on slide 4
  while  $f \wedge s(0) \in \{\mathbf{if}, \mathbf{while}\} \cup \mathcal{I}$  do
    command()
  end while
  if  $s(0) \notin \{\$, \mathbf{end}, \mathbf{else}\}$  then  $f := \mathbf{false}$  end if
end block

```

We can either detect the error here (Version 2) or leave the error to be detected later (Versions 0 and 1).

The command nonterminal

command $\rightarrow i := \text{exp}$ for all $i \in \mathcal{I}$
 command \rightarrow **if** exp **then** block **else** block **end if**
 command \rightarrow **while** exp **do** block **end while**

procedure *command*()

// Try to remove the a prefix described by *command* .

// See the contract for n a few slides back.

if $\neg f$ then return end if

if $s(0) = \mathbf{if}$ then

consume() *exp()* *expect(then)* *block()* *expect(else)*

block() *expect(end)* *expect(if)*

elseif $s(0) = \mathbf{while}$ then

consume() *exp()* *expect(do)* *block()* *expect(end)*

expect(while)

else if $s(0) \in \mathcal{I}$ then

consume() *expect(:=)* *exp()*

else

$f := \text{false}$

end if

end *command*

Parsing expressions

Recall that the rules for expressions are

$$\text{exp} \rightarrow \text{comparand}$$

$$\text{exp} \rightarrow \text{comparand} < \text{comparand}$$

Rewrite these rules to postpone the decision about which production to use until it matters

$$\text{exp} \rightarrow \text{comparand exp0}$$

$$\text{exp0} \rightarrow \epsilon$$

$$\text{exp0} \rightarrow < \text{comparand}$$

Write the procedures

```
procedure exp()
```

```
  // Try to remove a prefix described by exp.
```

```
  if  $\neg f$  then return end if
```

```
  comparand()    exp0()
```

```
end exp
```

```
procedure exp0()
```

```
  // Try to remove a prefix described by exp0
```

```
  if  $s(0) = <$  then consume() comparand() end if
```

```
end exp0
```

In-line the call to *exp0* to get

```
procedure exp()
```

```
  // Try to remove a prefix described by exp.
```

```
  if  $\neg f$  then return end if
```

```
  comparand()
```

```
  if  $s(0) = <$  then consume() comparand() end if
```

```
end exp
```

comparand \rightarrow term

comparand \rightarrow term + comparand

comparand \rightarrow term - comparand

rewrite as

comparand \rightarrow term comparand0

comparand0 \rightarrow + term comparand0

comparand0 \rightarrow - term comparand0

comparand0 \rightarrow ϵ

Write the procedures

procedure *comparand()*

 // Try to remove a prefix described by *comparand*.

 if $\neg f$ then return end if

term() *comparand0()*

end *comparand*

procedure *comparand0()*

 // Try to remove a prefix described by *comparand0*.

 if $\neg f$ then return end if

 if $s(0) \in \{+, -\}$ then *consume()* *term()* *comparand0()*

 end if

end *comparand*

After tail recursion removal and inlining, we have

procedure *comparand()*

 // Try to remove a prefix described by *comparand*.

 if $\neg f$ then return end if

term()

 while $f \wedge s(0) \in \{+, -\}$ do *consume()* *term()* end while

end *comparand*

Term is similar to comparand

term \rightarrow factor

term \rightarrow factor * term

term \rightarrow factor / term

procedure *term*()

// Try to remove a prefix described by *term*.

if $\neg f$ then return end if

factor()

while $f \wedge s(0) \in \{*, /\}$ do *consume*() *factor*() end while

end *term*

factor $\rightarrow n$ for all¹ $n \in \mathcal{N}$

factor $\rightarrow i$ for all $i \in \mathcal{I}$

factor $\rightarrow (\text{exp})$

procedure *factor*()

// Try to remove a prefix described by *factor*.

if $\neg f$ then return end if

if $s(0) \in \mathcal{N}$ then *consume*()

elseif $s(0) \in \mathcal{I}$ then *consume*()

elseif $s(0) = ($ then *consume*() *exp*() *expect*())

else $f := \text{false}$

end if

end *factor*

Exercise: find a variant expression that shows that we have no infinite loops or infinite recursion.

¹ Recall that \mathcal{N} is a finite subset of \mathbb{N} .

Generating machine code for expressions

Suppose we want to compile code for a stack machine

- The job of the code generated by procedures *factor*, *term*, *comparand*, and *exp* is to push a value.
- We'll ignore type checking and existence of variables
- We need the following instruction sequences
 - * *push*(*n*) pushes a number *n* on to the stack
 - * *fetch*(*i*) pushes the value of variable *i* onto the stack
 - * *mul* pops two values off the stack, multiplies them and pushes the result. *div* is similar to *mul*

procedure *factor*()

 if $\neg f$ then return end if

 if $s(0) \in \mathcal{N}$ then $m := m \hat{\ } \text{push}(s(0)) \text{consume}()$

 elseif $s(0) \in \mathcal{I}$ then $m := m \hat{\ } \text{fetch}(s(0)) \text{consume}()$

 elseif $s(0) = (\text{ then } \text{consume}() \text{exp}() \text{expect}())$

 else $f := \text{false}$ end if

end *factor*

term, *comparand*, and *exp* are similar to each other

procedure *term*()

 if $\neg f$ then return end if

factor()

 while $f \wedge s(0) \in \{*, /\}$ do

 val $op := s(0)$ *consume*() *factor*()

 if $op = *$ then $m := m \hat{\ } \text{mul}$ else $m := m \hat{\ } \text{div}$ end if

 end while

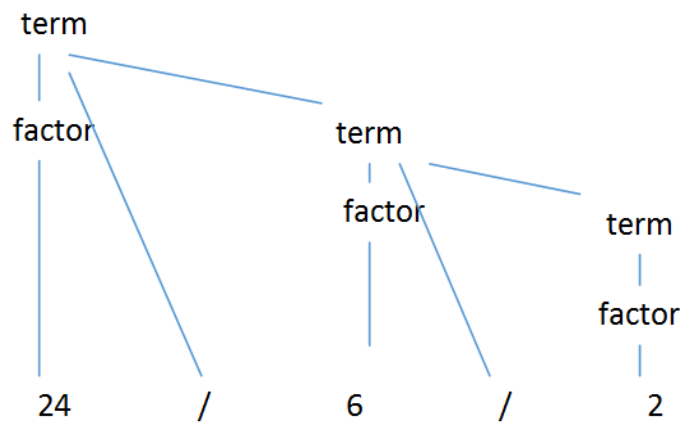
end *term*

What about associativity?

We want $-$ and $/$ to be left associative. E.g., $24/6/2$ should generate the same code as $(24/6)/2$.

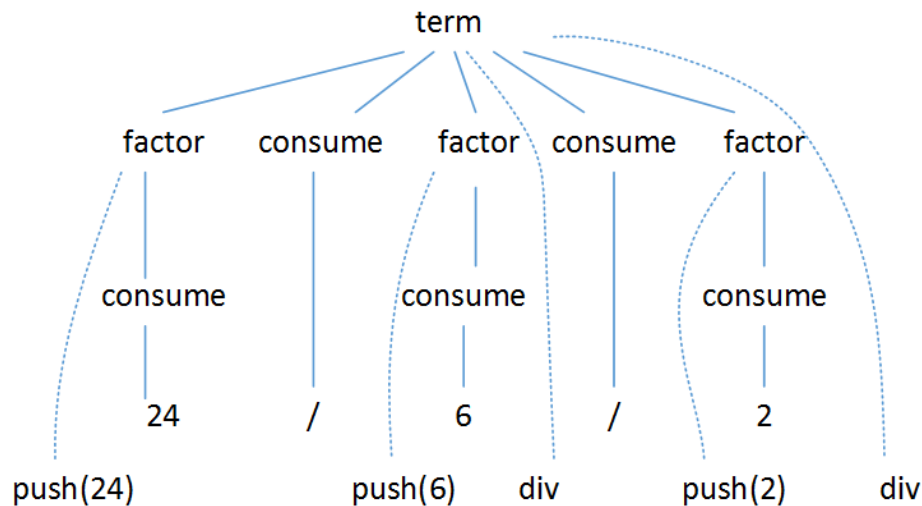
Our original grammar gets associativity “wrong” for $/$ and $-$.

Consider the parse tree for $term \xRightarrow{*} 24/6/2$.



This seems to associate the $/$ s the wrong way.

However, if you trace the actions of the compiler, you will see that the code generated for $24/6/2$ is correct because the operation is emitted at the right time.



If we look at a version without tail-call optimization, the choice is clearer.

```

procedure term()
  if  $\neg f$  then return end if
  factor()
  term0()
end term
procedure term0()
  if  $\neg f$  then return end if
  if  $s(0) \in \{*, /\}$  then
    val op :=  $s(0)$     consume()
    factor()
    // (a) emit instruction here for left associativity
    term0()
    // (b) emit instruction here for right associativity
  end if
end term0

```


Precedence

We need that $a + b * c + d * e$ generates the same code as $a + (b * c) + (d * e)$. Because of the way the grammar treats expressions, it does.

Generating code for assignment commands

Instruction

- $\text{store}(i)$ pops a value off the stack and stores it in the location for identifier i .

procedure *command*()

...

elseif $s(0) \in \mathcal{I}$ then

 val $i := s(0)$ *consume*()

expect(**:=**)

exp()

$m := m \hat{\ } \text{store}(i)$

else

...

Generating code for while commands

Instructions:

- $\text{branch}(a)$ branches to instruction a
- $\text{condBranch}(d)$ pops the stack and branches to d if the former top was false.
- I'll assume that the length of $\text{condBranch}(d)$ does not depend on d .

If the expression compiles to a sequence x and the block compiles to a sequence y , the while-loop compiles to a sequence

$$\begin{aligned}
 a & : x \\
 b & : \text{condBranch}(d) \\
 c & : y \\
 & \quad \text{branch}(a) \\
 d & :
 \end{aligned}$$

procedure *command*()

...

elseif $s(0) = \mathbf{while}$ then *consume*()

 val $a := m.\text{length}$ *exp*() *expect*(**do**)

 val $b := m.\text{length}$ $m := m \hat{\ } \text{condBranch}(0)$

 val $c := m.\text{length}$ *block*()

$m := m \hat{\ } \text{branch}(a)$

 val $d := m.\text{length}$ $m[b, ..c] := \text{condBranch}(d)$

expect(**end**) *expect*(**while**)

elseif

...

The rest of the compiler

I'll leave the rest of the compiler as an exercise:

- If commands,
- expression
- comparand
- block

Going further: Think about how you could

- Add variable declarations
- Add simple types and type checking
- Add procedures and procedure calls
- Add arrays
- Add classes and objects

When can we use recursive descent?

When can we use recursive descent parsing?

When it is possible to choose between the productions for a nonterminal based on

- Information already seen
- The next few symbols of input

In particular there is a set of grammars for which RDP is particularly easy. These grammars allow the choice to be made by looking only at the next item of input.

Such a grammar is called “LL(1)”.

LL(1)

Recall: If a nonterminal n has productions

$$(n \rightarrow \alpha), (n \rightarrow \beta), (n \rightarrow \gamma) \in P,$$

we write a subroutine like this:

```

procedure  $n()$ 
// Try to remove a prefix described by  $n$  .
  if  $\neg f$  then return end if
  if ? then  $[[\alpha]]$  else if ? then  $[[\beta]]$  else if ? then  $[[\gamma]]$ 
  else  $f := \text{false}$  end if
end

```

Often the guard only needs to look at the next input symbol.

Associate with each production $n \rightarrow \alpha$ with a “selector set” $\text{sel}(n \rightarrow \alpha) \subseteq A \cup \{\$\}$

```

procedure  $n()$ 
// Try to remove a prefix described by  $n$  .
  if  $\neg f$  then return end if
  if  $s(0) \in \text{sel}(n \rightarrow \alpha)$  then  $[[\alpha]]$ 
  else if  $s(0) \in \text{sel}(n \rightarrow \beta)$  then  $[[\beta]]$ 
  else if  $s(0) \in \text{sel}(n \rightarrow \gamma)$  then  $[[\gamma]]$ 
  else  $f := \text{false}$  end if
end  $n$ 

```

If for all distinct productions $n \rightarrow \alpha, n \rightarrow \beta$, $\text{sel}(n \rightarrow \alpha) \cap \text{sel}(n \rightarrow \beta) = \emptyset$, then the grammar is called $LL(1)$, and we can write a recursive descent parser for it.

Computing selector sets:

- **First symbols:** If $\alpha \xRightarrow{*} at$ with $a \in A$ then $a \in \text{sel}(n \rightarrow \alpha)$
- **Following symbols:** $a \in \text{sel}(n \rightarrow \alpha)$ if $\alpha \xRightarrow{*} \epsilon$ and $a \in A \cup \{\$\}$ can follow n in a derivation from $n_{\text{start}}\$\text{---}$ — i.e. if there is a derivation

$$n_{\text{start}}\$ \xRightarrow{*} tna u \Longrightarrow t\alpha a u \xRightarrow{*} tau$$
 with $t \in A^*$, $u \in (A \cup \{\$\})^*$.
- Nothing else is in $\text{sel}(n \rightarrow \alpha)$

Example: The start symbol is B

$$B \rightarrow CB \quad B \rightarrow \epsilon$$

$$C \rightarrow \text{id} := E \quad C \rightarrow \text{if } E \text{ then } B D \text{ end if}$$

$$D \rightarrow \text{else } B \quad D \rightarrow \epsilon$$

$$E \rightarrow \text{id}$$

The selector set of $B \rightarrow CB$ is the first symbols of CB which are $\{\text{id}, \text{if}\}$.

The selector set of $B \rightarrow \epsilon$ is the symbols that can follow B which are $\{\text{else}, \text{end}, \$\}$.

Exercise. Show that the following grammar, with start symbol B , is not LL(1)

$$B \rightarrow CB \quad B \rightarrow \epsilon$$

$$C \rightarrow \text{id} := E \quad C \rightarrow \{B\} \quad C \rightarrow \text{if } E \text{ then } C D$$

$$D \rightarrow \text{else } C \quad D \rightarrow \epsilon$$

$$E \rightarrow \text{id}$$

If a grammar is not LL(1), we can still often use recursive descent, e.g., by looking more symbols ahead.

Here are a few tricks of the trade to make a grammar LL(1), or at least more suitable for RDP.

- Factor: Example: Replace

$$\text{command} \rightarrow \mathbf{id} := \text{exp}$$

$$\text{command} \rightarrow \mathbf{id} (\text{args})$$

with

$$\text{command} \rightarrow \mathbf{id} \text{ more}$$

$$\text{more} \rightarrow := \text{exp}$$

$$\text{more} \rightarrow (\text{args})$$

More generally, replace productions

$$n \rightarrow \alpha a \beta$$

$$n \rightarrow \alpha b \gamma,$$

where ~~$a, b \in A$~~ and $\alpha, \beta, \gamma \in (A \cup N)^*$, with

$$n \rightarrow \alpha p$$

$$p \rightarrow \alpha \beta$$

$$p \rightarrow b \gamma,$$

where p is a fresh nonterminal.

- Remove left recursion: Example: Replace

$$\text{type} \rightarrow \text{type} []$$

$$\text{type} \rightarrow \mathbf{int}$$

with

$$\text{type} \rightarrow \mathbf{int} \text{type0}$$

$$\text{type0} \rightarrow [] \text{type0}$$

$$\text{type0} \rightarrow \epsilon$$

More generally, replace

$$n \rightarrow n\alpha$$

$$n \rightarrow \beta,$$

where $\alpha, \beta \in (A \cup N)^*$, with

$$n \rightarrow \beta p$$

$$p \rightarrow \alpha p$$

$$p \rightarrow \epsilon,$$

where p is a fresh nonterminal.

Most formats can be parsed by recursive descent, one way or another.

Tools

While writing recursive descent parsers is straightforward for simple grammars, it can be error prone and tedious as grammars evolve and get larger.

Luckily there are a large number of tools that convert grammars to parsers. Examples:

- JavaCC.
 - * Allows grammars in which the RHS of each production is a regular expression.
 - * Produces recursive descent parsers written in Java or C++.
 - * Calculates the guard expressions automatically for most grammars
 - * Allows the programmer to intervene in cases the automatic rules don't handle
 - * Allows the programmer to annotate the grammar with bits of Java (or C++) code that are interpolated into the parser.
- ANTLR 4
 - * Similar to JavaCC
 - * Automatic treatment of left recursion and operator precedence
- Yacc/Bison
 - * Produces bottom-up parsers
 - * Handles a large class of grammars automatically
 - * No need to factor or remove left recursion