

# Grammars en Parsing

Principes

Grammar flow analyse

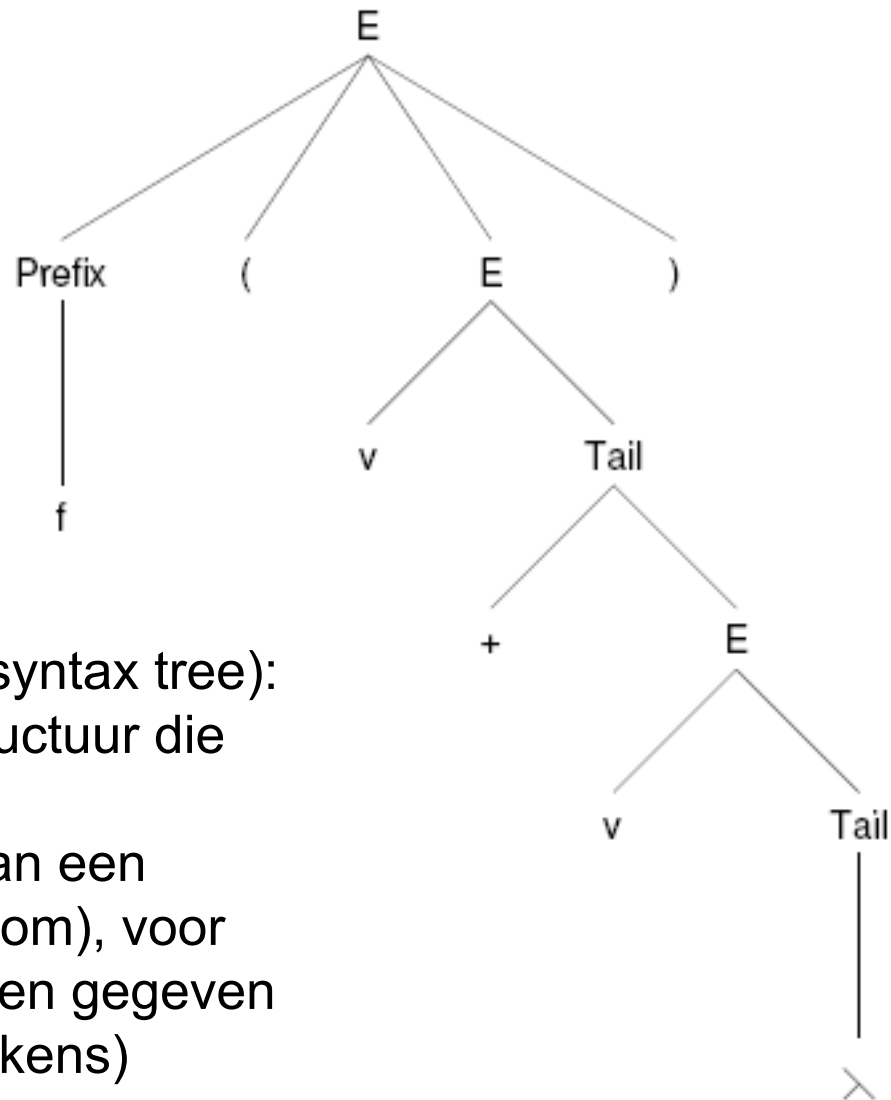
## Context-free grammars (eventueel in BNF)

- Definiëren de structuur van geldige programma's
- Drukken de gewenste betekenis uit.

1	$E$	$\rightarrow$	Prefix	(	$E$	)
2				$v$	Tail	
3	Prefix	$\rightarrow$	$f$			
4				$\lambda$		
5	Tail	$\rightarrow$	$+$	$E$		
6				$\lambda$		

Figure 4.1: A simple expression grammar.

---



Parse tree (of abstract syntax tree):  
 de belangrijkste hulpstructuur die  
 de compiler gebruikt  
 Parsing = constructie van een  
 parse tree (afleidingsboom), voor  
 een gegeven CFG en een gegeven  
 invoerwoord (=rij van tokens)

Figure 4.2: The parse tree for  $f ( v + v ) .$

Dubbelzinnigheid is ongewenst: de boom moet de juiste betekenis uitdrukken  
Niet elke CFG is geschikt voor elke parsing methode.

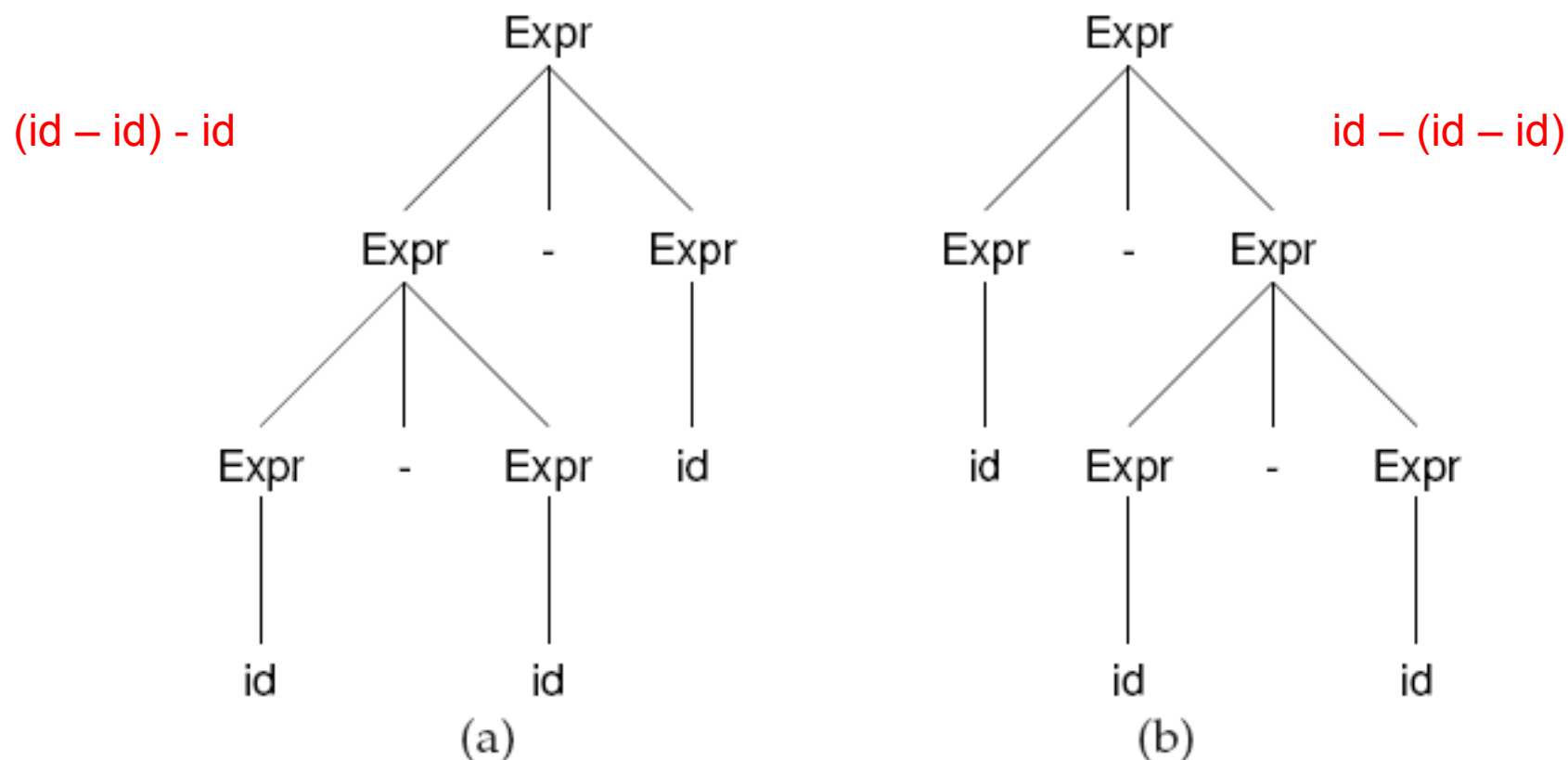


Figure 4.3: Two parse trees for id - id - id.

Voor de leesbaarheid gebruikt men vaak de Backus-Naur form, die kan gemakkelijk omgezet worden naar “echte” CFG –vorm zonder extra constructies

```
foreach  $p \in Prods$  of the form " $A \rightarrow \alpha [ X_1 \dots X_n ] \beta$ " do  
   $N \leftarrow \text{NEWNONTERM}()$   
   $p \leftarrow "A \rightarrow \alpha N \beta"$  [ ]: optioneel  
   $Prods \leftarrow Prods \cup \{ "N \rightarrow X_1 \dots X_n" \}$   
   $Prods \leftarrow Prods \cup \{ "N \rightarrow \lambda" \}$   
foreach  $p \in Prods$  of the form " $B \rightarrow \gamma \{ X_1 \dots X_m \} \delta$ " do  
   $M \leftarrow \text{NEWNONTERM}()$   
   $p \leftarrow "B \rightarrow \gamma M \delta"$  { }: iteratie  
   $Prods \leftarrow Prods \cup \{ "M \rightarrow X_1 \dots X_n M" \}$   
   $Prods \leftarrow Prods \cup \{ "M \rightarrow \lambda" \}$ 
```

Figure 4.4: Algorithm to transform a BNF grammar into standard form.

Parsing (constructie van de parse tree):  
ofwel **top-down** ofwel **bottom-up**

Voorbeeld:

Program  $\rightarrow$  begin Stmts end \$

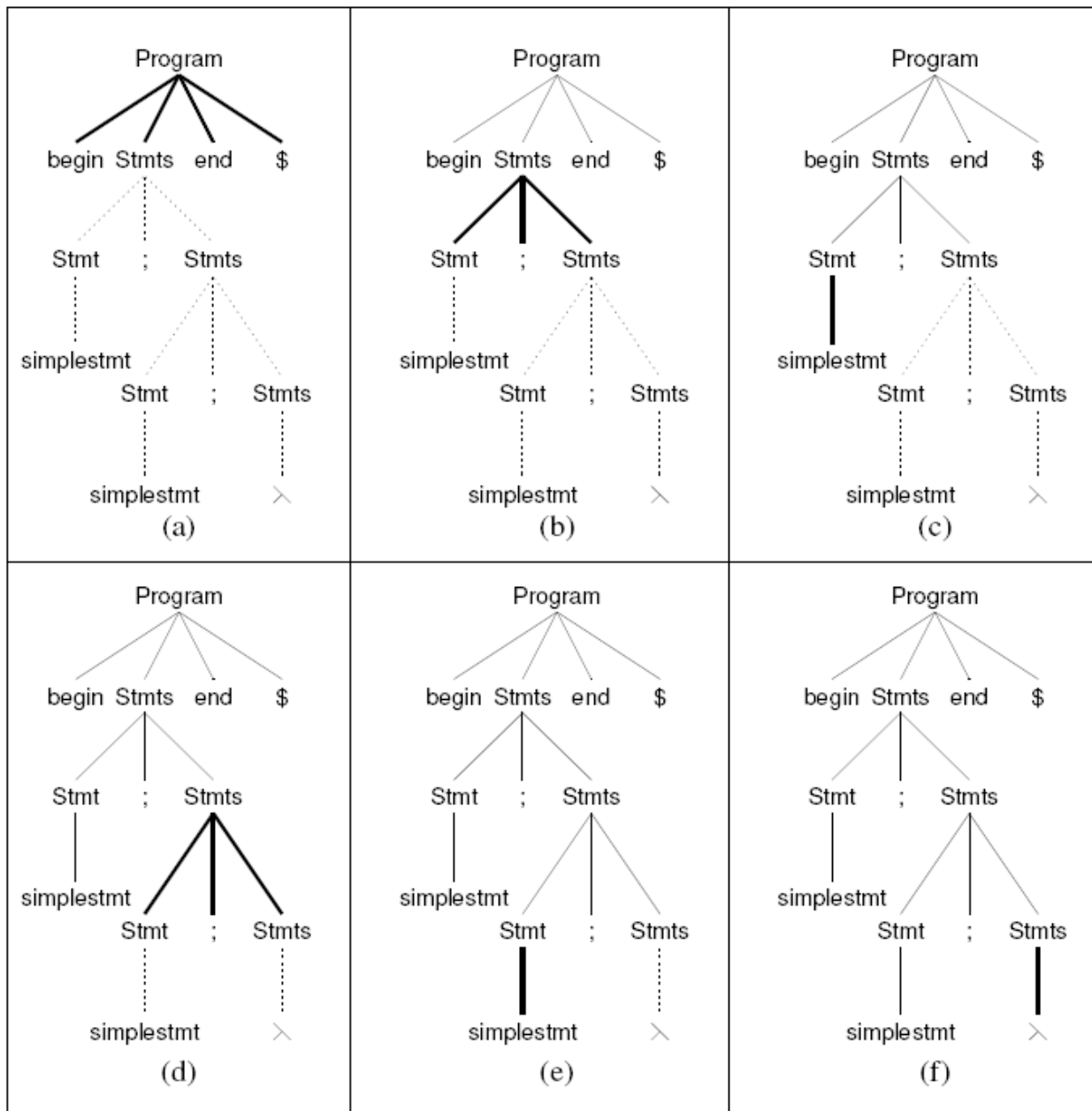
Stmts  $\rightarrow$  Stmt ; Stmts

|  $\lambda$

Stmt  $\rightarrow$  simplestmt

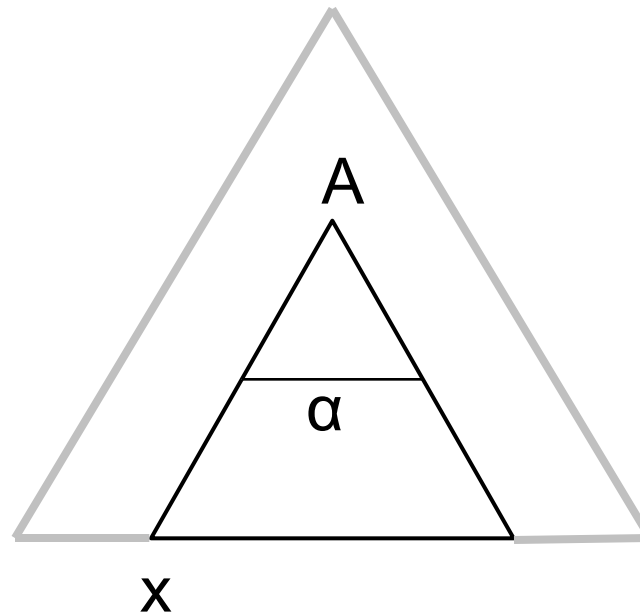
| begin Stmts end

Top-down parse van “begin simplestmt ; simplestmt ; end \$”



- De “recursive descent” methode (zie voorbeeld uit inleiding) was top-down
- Belangrijkste probleem: hoe kiezen tussen de alternatieven voor een nonterminaal  $A$ ? Test of het volgende invoersymbool  $x$  kan voorkomen als eerste symbool van een woord afgeleid van het gekozen alternatief  $\alpha$ .

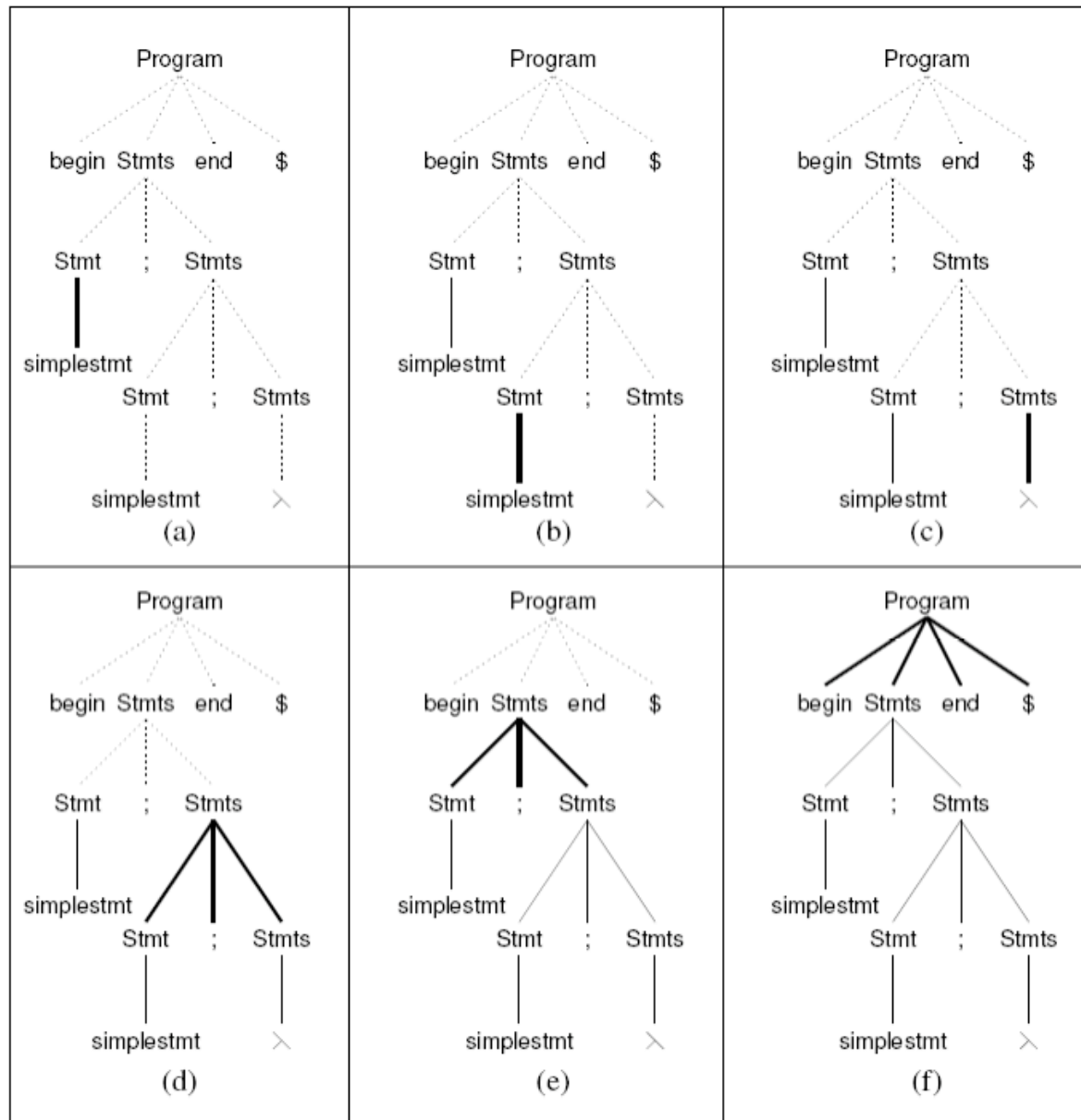
$A \rightarrow \alpha \mid \beta \mid \dots$



In het voorbeeld, stap(c): de input **simplestmt** zorgt ervoor dat het 1e alternatief gekozen wordt voor **Stmt**



Bottom-up parse van “begin simplestmt ; simplestmt ; end \$”



- Belangrijkste probleem: herkennen van de rechterkant van een productie, in een sentential form (= woord met zowel terminale als niet-terminale symbolen), om dan een “reductie” stap te kunnen uitvoeren.
- Je hebt meer informatie wanneer je de keuzes moet maken: de bottom-up methode is krachtiger (werkt voor meer grammatica's) dan de top-down methode.

# Algoritmen voor Grammar Analyse

Bij het gebruik van CFGs en parsing algoritmen gebruiken we een aantal eigenschappen van CFGs en hun symbolen die op voorhand kunnen berekend worden (m.a.w. ze hangen enkel af van de CFG, niet van het te parsen woord)

- Kan er van een gegeven symbool wel een terminaal woord afgeleid worden (= productiveness)
- Kan een gegeven symbool wel voorkomen in een woord afgeleid van het beginsymbool? (= reachability)

- Kan het lege woord afgeleid worden van een gegeven nonterminaal symbool?
- Welke terminale symbolen kunnen voorkomen als eerste symbool van een woord afgeleid van een gegeven woord  $\alpha$  ? Notatie: **First( $\alpha$ )**
- Welke symbolen kunnen volgen op een gegeven nonterminal A ? Notatie: **Follow(A)**

Productiveness: kan van een nonterminal A een  
terminaal woord (inclusief het lege) afgeleid worden?

Initialiseer W (worklist met nonterminals): initieel de nonterminals die LHS  
zijn van een productie zonder nonterminals in de RHS  
count[r]: initieel het aantal nonterminals in de RHS van de productie r  
Productive: initieel leeg.

```
while nonempty(W) {  
  X := head(W); W := tail(W);  
  if (X not in Productive) {  
    Productive := Productive  $\cup$  {X};  
    forall (productions r having X in their RHS) {  
      count[r]-- ;  
      if (count[r] = 0) W := (LHS(r)) :: W ;  
    }  
  }  
}
```

Reachability: kan een nonterminal voorkomen in een woord afgeleid van het beginsymbool?

**Reduceren** van een CFG: laat eerst de niet-productieve symbolen weg, samen met de producties waarin ze voorkomen, en doe dan hetzelfde met onbereikbare symbolen (Niet andersom!)

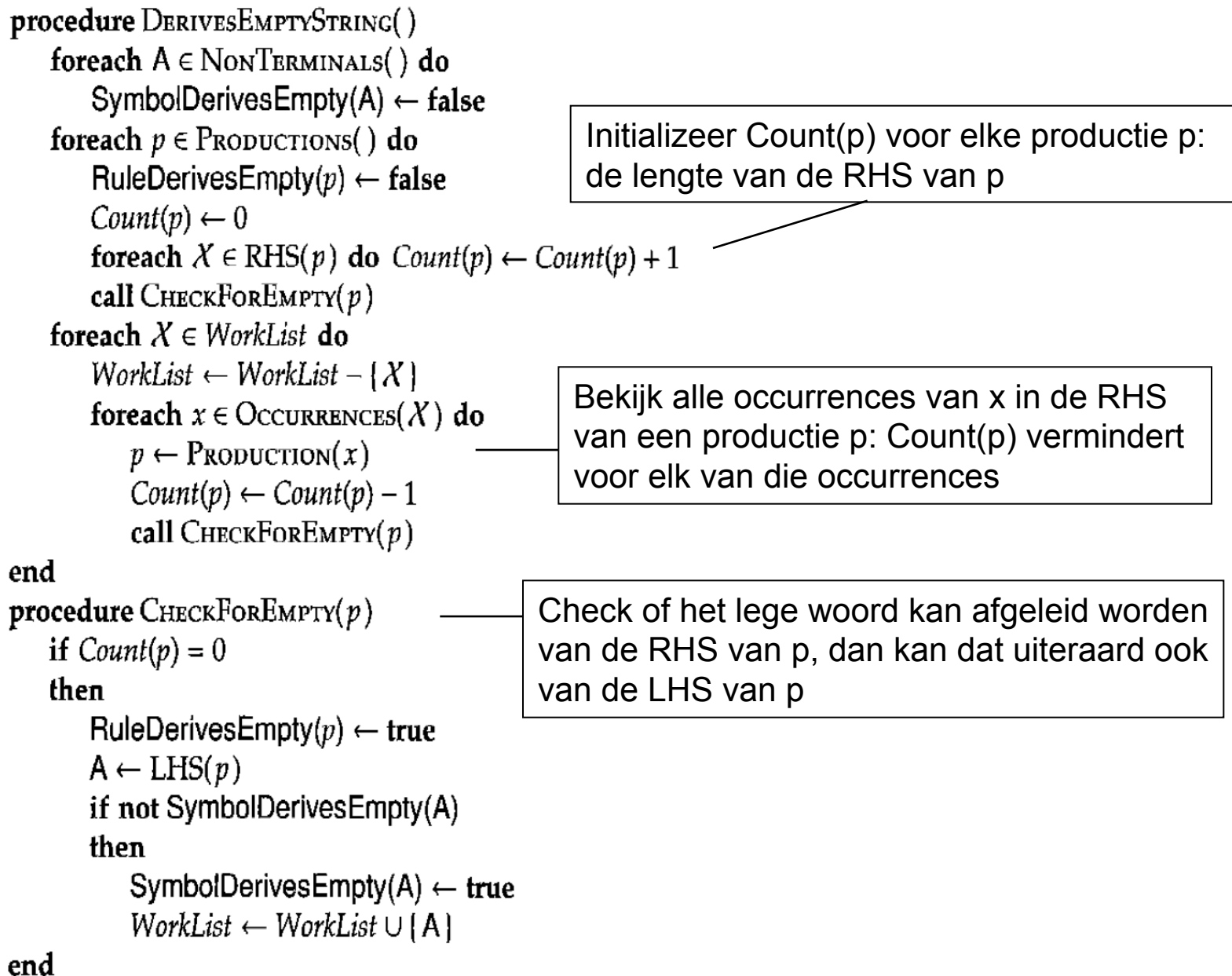


Figure 4.7: Algorithm for determining nonterminals and productions that can derive  $\lambda$ .

$\text{First}(\alpha)$  is de set van terminale symbolen die als eerste symbool voorkomen in een woord  $u$  dat van  $\alpha$  kan afgeleid worden.  $\alpha$  is een woord dat kan bestaan uit zowel terminale als niet-terminale symbolen.

$\text{First}(\alpha)$  kan als volgt berekend worden.



```

function FIRST( $\alpha$ ) returns Set
    foreach  $A \in \text{NonTerminals}()$  do  $\text{VisitedFirst}(A) \leftarrow \text{false}$            (9)
     $\text{ans} \leftarrow \text{INTERNALFIRST}(\alpha)$ 
    return ( $\text{ans}$ )
end
function INTERNALFIRST( $X\beta$ ) returns Set
    if  $X\beta = \perp$                                                          (10)
    then return ( $\emptyset$ )
    if  $X \in \Sigma$                                                          (11)
    then return ( $\{X\}$ )
     $\star$   $X$  is a nonterminal.  $\star$  / (12)
     $\text{ans} \leftarrow \emptyset$ 
    if not  $\text{VisitedFirst}(X)$ 
    then
         $\text{VisitedFirst}(X) \leftarrow \text{true}$                                      (13)
        foreach  $\text{rhs} \in \text{ProductionsFor}(X)$  do
             $\text{ans} \leftarrow \text{ans} \cup \text{INTERNALFIRST}(\text{rhs})$                  (14)
        if  $\text{SymbolDerivesEmpty}(X)$                                          (15)
        then  $\text{ans} \leftarrow \text{ans} \cup \text{INTERNALFIRST}(\beta)$ 
        return ( $\text{ans}$ )                                                     (16)
    end

```

Figure 4.8: Algorithm for computing  $\text{First}(\alpha)$ .

Voorbeeld-grammatica:

1	$E$	$\rightarrow$	Prefix	(	$E$	)
2			$v$	Tail		
3	Prefix	$\rightarrow$	$f$			
4			$\lambda$			
5	Tail	$\rightarrow$	$+$	$E$		
6			$\lambda$			

Figure 4.1: A simple expression grammar.

---

Level	First $\mathcal{X}$	$\beta$	$ans$	Marker	Done? (★=Yes)	Comment
First(Tail)						
0	Tail	$\perp$	{ }	(12)		
1	+	E	{ + }	(11)	★	Tail $\rightarrow$ +E
1	$\perp$	$\perp$	{ }	(10)	★	Tail $\rightarrow \lambda$
0			{ + }	(14)		After all rules for Tail
1	$\perp$	$\perp$	{ }	(10)	★	Since $\beta = \perp$
0			{ + }	(15)	★	Final answer
First(Prefix)						
0	Prefix	$\perp$	{ }	(12)		
1	f	$\perp$	{ f }	(11)	★	Prefix $\rightarrow$ f
1	$\perp$	$\perp$	{ }	(10)	★	Prefix $\rightarrow \lambda$
0			{ f }	(14)		After all rules for Prefix
1	$\perp$	$\perp$	{ }	(10)	★	Since $\beta = \perp$
0			{ f }	(15)	★	Final answer
First(E)						
0	E	$\perp$	{ }	(12)		
1	Prefix	( E )	{ }	(12)		E $\rightarrow$ Prefix ( E )
1			{ f }	(16)		Computation shown above
2	(	E)	{ ( }	(11)	★	Since Prefix $\Rightarrow^* \lambda$
1			{ f, ( }	(15)	★	Results due to E $\rightarrow$ Prefix ( E )
1	v	Tail	{ v }	(11)	★	E $\rightarrow$ v Tail
1	$\perp$	$\perp$	{ }	(10)		Since $\beta = \perp$
0			{ f, (, v }	(15)	★	Final answer

1  $S \rightarrow A \ B \ c$   
 2  $A \rightarrow a$   
 3  $\quad \mid \lambda$   
 4  $B \rightarrow b$   
 5  $\quad \mid \lambda$

Level	First $\mathcal{X}$	$\beta$	<i>ans</i>	Marker	Done? (★=Yes)	Comment
First (B)						
0	B	$\perp$	{ }	(12)		
1	b	$\perp$	{b}	(11)	★	$B \rightarrow b$
1	$\perp$	$\perp$	{ }	(10)	★	$B \rightarrow \lambda$
0			{b}	(15)	★	Final answer
First (A)						
0	A	$\perp$	{ }	(12)		
1	a	$\perp$	{a}	(11)	★	$A \rightarrow a$
1	$\perp$	$\perp$	{ }	(10)	★	$A \rightarrow \lambda$
0			{a}	(15)	★	Final answer
First (S)						
0	S	$\perp$	{ }	(12)		
1	A	B c	{a}	(16)		Computation shown above
2	B	c	{b}	(16)		Because $A \Rightarrow^* \lambda$ ; computation shown above
3	c	$\perp$	{c}	(11)	★	Because $B \Rightarrow^* \lambda$
2			{b,c}	(15)	★	
1			{a,b,c}	(15)	★	
0			{a,b,c}	(15)	★	

**Follow(A)** is, voor een nonterminal A, de set van terminale symbolen die op A kunnen volgen in een willekeurige sentential form, m.a.w. in een willekeurig woord dat van het startsymbool kan afgeleid worden.

**Tail(a)** is, voor een occurrence a in de RHS van een productie, het deel van die RHS dat volgt op a.

```

function FOLLOW(A) returns Set
    foreach A ∈ NonTERMINALS() do
        VisitedFollow(A) ← false
        ans ← INTERNALFOLLOW(A)
        return (ans)
    end
function INTERNALFOLLOW(A) returns Set
    ans ← ∅
    if not VisitedFollow(A)
        then
            VisitedFollow(A) ← true
            foreach a ∈ OCCURRENCES(A) do
                ans ← ans ∪ FIRST(TAIL(a))
                if ALLDERIVEEMPTY(TAIL(a))
                    then
                        targ ← LHS(PRODUCTION(a))
                        ans ← ans ∪ INTERNALFOLLOW(targ)
            return (ans)
        end
function ALLDERIVEEMPTY(γ) returns Boolean
    foreach X ∈ γ do
        if not SymbolDerivesEmpty(X) or X ∈ Σ
            then return (false)
    return (true)
end

```

Figure 4.11: Algorithm for computing Follow(*A*).

Level	Rule	Marker	Result	Comment
Follow (B)				
0				
0	$S \rightarrow A \underline{B} c$	(21)	{ c }	
0		(24)	{ c }	Returns
Follow (A)				
0				
0	$S \rightarrow \underline{A} B c$	(21)	{ b,c }	
0		(24)	{ b,c }	Returns
Follow (S)				
0				
0		(24)	{ }	Returns

Figure 4.12: Follow sets for the grammar in Figure 4.10. Note that  $\text{Follow}(S) = \{ \}$  because S does not appear on the RHS of any production.

Level	Rule	Marker	Result	Comment
		F	(Prefix)	
0		F	(Prefix)	
0	$E \rightarrow \underline{\text{Prefix}} ( E )$	(21)	{ ( }	
		F	(E)	
0		F	(E)	
0	$E \rightarrow \text{Prefix} ( \underline{E} )$	(21)	{ ) }	
0	$\text{Tail} \rightarrow + \underline{E}$	(23)	{ }	
1		F	(Tail)	
1	$E \rightarrow v \underline{\text{Tail}}$	(23)	{ }	
2		F	(E)	
		(18)	{ }	Recursion avoided
1		(24)	{ }	Returns
0		(24)	{ ) }	Returns
		F	(Tail)	
0		F	(Tail)	
0	$E \rightarrow v \underline{\text{Tail}}$	(23)	{ }	
1		F	(E)	
1	$E \rightarrow \text{Prefix} ( \underline{E} )$	(21)	{ ) }	
1	$\text{Tail} \rightarrow + \underline{E}$	(23)	{ }	
2		F	(Tail)	
		(18)	{ }	Recursion avoided
1		(24)	{ ) }	Returns
0		(24)	{ ) }	Returns

Figure 4.13: Follow sets for the nonterminals of Figure 4.1.



# Top-Down Parsing

**Predict sets** voor Top-down parsing: voor de keuze van een productie  $p$  onderscheiden we 2 gevallen: van de RHS kan het lege woord afgeleid worden, of niet.

```
function Predict( $p : A \rightarrow X_1 \dots X_m$ ) : Set  
   $ans \leftarrow \text{First}(X_1 \dots X_m)$   
  if RuleDerivesEmpty( $p$ )  
  then  
     $ans \leftarrow ans \cup \text{Follow}(A)$   
  return ( $ans$ )  
end
```

Figure 5.1: Computation of Predict sets.

---

## Voorbeeld

```
1 S → A C $
2 C → c
3   | λ
4 A → a B C d
5   | B Q
6 B → b B
7   | λ
8 Q → q
9   | λ
```

Figure 5.2: A CFGs.

---

Rule Number	N	$X_1 \dots X_m$	$\text{First}(X_1 \dots X_m)$	Derives Empty?	Follow(N)	Answer
1	S	A C \$	a,b,q,c,\$	No		a,b,q,c,\$
2	C	c	c	No		c
3		$\lambda$		Yes	d,\$	d,\$
4	A	a B C d	a	No		a
5		B Q	b,q	Yes	c,\$	b,q,c,\$
6	B	b B	b	No		b
7		$\lambda$		Yes	q,c,d,\$	q,c,d,\$
8	Q	q	q	No		q
9		$\lambda$		Yes	c,\$	c,\$

Figure 5.3: Predict calculation for the grammar of Figure 5.2.

```

function IsLL1(G) returns Boolean
  foreach A ∈ N do
    PredictSet ← ∅
    foreach p ∈ ProductionsFor(A) do
      if Predict(p) ∩ PredictSet ≠ ∅
      then return (false)
      PredictSet ← PredictSet ∪ Predict(p)
    return (true)
  end

```

Predict sets mogen niet overlappen

Figure 5.4: Algorithm to determine if a grammar *G* is LL(1).

---

```

procedure MATCH(ts, token)
  if ts.PEEK() = token
  then call ts.ADVANCE()
  else call ERROR(Expected token)
end

```

Figure 5.5: Utility for matching tokens in an input stream.

---

## Gebruik van de Predict sets

```
procedure A(ts)
  switch (...)
    case ts.PEEK() ∈ Predict( $p_1$ )
      /* Code for  $p_1$                                */
    case ts.PEEK() ∈ Predict( $p_i$ )
      /* Code for  $p_2$                                */
    /* .                                              */
    /* .                                              */
    /* .                                              */
    case ts.PEEK() ∈ Predict( $p_n$ )
      /* Code for  $p_n$                                */
    case default
      /* Syntax error                                */
  end
```

Figure 5.6: A typical recursive-descent procedure. Successful LL(1) analysis ensures that only one of the case predicates is true.

```

procedure S()
  switch (...)
    case ts.PEEK() ∈ { a, b, q, c, $ }
      call A()
      call C()
      call MATCH($)
    end
  procedure C()
    switch (...)
      case ts.PEEK() ∈ { c }
        call MATCH(c)
      case ts.PEEK() ∈ { d, $ }
        return ()
      end
    procedure A()
      switch (...)
        case ts.PEEK() ∈ { a }
          call MATCH(a)
          call B()
          call C()
          call MATCH(d)
        case ts.PEEK() ∈ { b, q, c, $ }
          call B()
          call Q()
        end
      procedure B()
        switch (...)
          case ts.PEEK() ∈ { b }
            call MATCH(b)
            call B()
          case ts.PEEK() ∈ { q, c, d, $ }
            return ()
          end
        procedure Q()
          switch (...)
            case ts.PEEK() ∈ { q }
              call MATCH(q)
            case ts.PEEK() ∈ { c, $ }
              return ()
            end
          end
        end
      end
    end
  end

```

Gebruik in een recursive descent parser

## Tabel-gestuurde LL(1) parser

```
procedure LLPARSER(ts)  
  call PUSH(S)  
  accepted  $\leftarrow$  false  
  while not accepted do  
    if TOS( )  $\in \Sigma$   
    then  
      call MATCH(ts, TOS( ))  
      if TOS( ) = $  
      then accepted  $\leftarrow$  true  
      call POP( )  
    else  
      p  $\leftarrow$  LLtable[TOS( ), ts.PEEK( )]  
      if p = 0  
      then  
        call ERROR(Syntax error—no production applicable)  
      else call APPLY(p)  
  end  
procedure APPLY(p :  $A \rightarrow X_1 \dots X_m$ )  
  call POP( )  
  for i = m downto 1 do  
    call PUSH(Xi)  
end
```

Figure 5.8: Generic LL(1) parser.



```

procedure FILLTABLE(LLtable)
  foreach  $A \in N$  do
    foreach  $a \in \Sigma$  do LLtable[A][a]  $\leftarrow 0$ 
  foreach  $A \in N$  do
    foreach  $p \in \text{ProductionsFor}(A)$  do
      foreach  $a \in \text{Predict}(p)$  do LLtable[A][a]  $\leftarrow p$ 
end

```

Figure 5.9: Construction of an LL(1) parse table.

---

Nonterminal	Lookahead					
	a	b	c	d	q	\$
S	1	1	1		1	1
C			2	3		3
A	4	5	5		5	5
B		6	7	7	7	7
Q			9		8	9

Figure 5.10: LL(1) table. The blank entries should trigger error actions in the parser.

---

Parse Stack	Action	Remaining Input
S		abbd c\$
\$CA	Apply 1: $S \rightarrow AC$	abbd c\$
\$CdCBa	Apply 4: $A \rightarrow aBCd$	abbd c\$
\$CdCB	Match	bbd c\$
\$CdCBb	Apply 6: $B \rightarrow bB$	bbd c\$
\$CdCB	Match	bdc\$
\$CdCBb	Apply 6: $B \rightarrow bB$	bdc\$
\$CdCB	Match	dc\$
\$CdC	Apply 7: $B \rightarrow \lambda$	dc\$
\$Cd	Apply 3: $C \rightarrow \lambda$	dc\$
\$C	Match	c\$
\$c	Apply 2: $C \rightarrow c$	c\$
\$	Match	\$
	Accept	

Figure 5.11: Trace of an LL(1) parse. The stack is shown in the left column, with top-of-stack as the rightmost character. The input string is shown in the right column, processed from left to right.

```

1 Stmt    → if Expr then StmtList endif
2          | if Expr then StmtList else StmtList endif
3 StmtList → StmtList ; Stmt
4          | Stmt
5 Expr     → var + Expr
6          | var

```

Figure 5.12: A grammar with common prefixes.

---

```

procedure FACTOR( )
  foreach  $A \in N$  do
     $\alpha \leftarrow \text{LongestCommonPrefix}(\text{ProductionsFor}(A))$ 
    while  $|\alpha| > 0$  do
       $V \leftarrow \text{new NonTerminal}()$ 
       $\text{Productions} \leftarrow \text{Productions} \cup \{ A \rightarrow \alpha V \}$ 
      foreach  $p \in \text{ProductionsFor}(A) \mid \text{RHS}(p) = \alpha \beta_p$  do
         $\text{Productions} \leftarrow \text{Productions} - \{ p \}$ 
         $\text{Productions} \leftarrow \text{Productions} \cup \{ V \rightarrow \beta_p \}$ 
       $\alpha \leftarrow \text{LongestCommonPrefix}(\text{ProductionsFor}(A))$ 
    end
  end

```

(13)

Figure 5.13: Factoring common prefixes.

---

```

1 Stmt    → if Expr then StmtList V1
2 V1     → endif
3         | else StmtList endif
4 StmtList → StmtList ; Stmt
5         | Stmt
6 Expr     → var V2
7 V2     → + Expr
8         | λ

```

Figure 5.14: Factored version of the grammar in Figure 5.12.

---

```

procedure ELIMINATELEFTRECURSION()
  foreach  $A \in N$  do
    if  $\exists r \in \text{ProductionsFor}(A) \mid \text{RHS}(r) = A\alpha$ 
    then
       $X \leftarrow \text{new NonTerminal}()$ 
       $Y \leftarrow \text{new NonTerminal}()$ 
      foreach  $p \in \text{ProductionsFor}(A)$  do
        if  $p = r$ 
        then  $\text{Productions} \leftarrow \text{Productions} \cup \{A \rightarrow X Y\}$  and remove  $r$ 
        else  $\text{Productions} \leftarrow \text{Productions} \cup \{X \rightarrow \text{RHS}(p)\}$ 
       $\text{Productions} \leftarrow \text{Productions} \cup \{Y \rightarrow \alpha Y, Y \rightarrow \lambda\}$ 
    end
  end

```

Figure 5.15: Eliminating left recursion.

1	Stmt	$\rightarrow$ if Expr then StmtList $V_1$
2	$V_1$	$\rightarrow$ endif
3		else StmtList endif
4	StmtList	$\rightarrow$ X Y
5	X	$\rightarrow$ Stmt
6	Y	$\rightarrow$ ; Stmt Y
7		$\lambda$
8	Expr	$\rightarrow$ var $V_2$
9	$V_2$	$\rightarrow$ + Expr
10		$\lambda$

Figure 5.16: LL(1) version of the grammar in Figure 5.14.

---

1  $S \rightarrow \text{Stmt } \$$   
2  $\text{Stmt} \rightarrow \text{if expr then Stmt else Stmt}$   
3       |  $\text{if expr then Stmt}$   
4       |  $\text{other}$

Figure 5.17: Grammar for if-then-else.

---