

Kap.4 – del I Top Down Parsing INF5110 – v2006

Arne Maus
Ifi, UiO

Innhold

- Motivering
 - Boka gir først parsingen uten First/Follow-mengder og så innfører dem.
 - Vi tar "teorien" først
- First og Follow-mengder
- Fjerning av venstre-rekursjon
- Venstre-faktorisering
- Rekursive nedstigning
- LL(1) grammatikker

First og Follow-mengder, hvorfor

- Når vi leser token og skal analysere en ikke-terminal, ønsker vi:
 - finne ut hvilken produksjon (høyreside-alternativ) som kommer nå - ser på:
 - First-mengdene
Det vi leste må da være starten på
 - I kap 4: Sjekke at det som etterfølger det vi nå har lest og tolket som en ikke-terminal, er lovlig etterfølger (feil-sjekking)
I kap. 5 (bottom up) finne riktig reduksjon av det vi hittil har lest på input .
 - Follow-mengdene

FirstMengder

Def $\left\{ \begin{array}{l} \text{First}(A) = \{ a \mid \text{finnes avledning } A \Rightarrow^* a \alpha \} \\ \text{Dessuten: Om } A \text{ "er utnullbar", så er } \epsilon \in \text{First}(A) \end{array} \right.$

Per def er da: $\text{First}(a) = \{ a \}$
Def "A er utnullbar" $\Leftrightarrow A \Rightarrow^* \epsilon$

Let X be a grammar symbol (a terminal or nonterminal) or ϵ . Then the set **First**(X) consisting of terminals, and possibly ϵ , is defined as follows.

1. If X is a terminal or ϵ , then $\text{First}(X) = \{X\}$.
2. If X is a nonterminal, then for each production choice $X \rightarrow X_1 X_2 \dots X_n$, $\text{First}(X)$ contains $\text{First}(X_i) - \{\epsilon\}$. If also for some $i < n$, all the sets $\text{First}(X_1), \dots, \text{First}(X_i)$ contain ϵ , then $\text{First}(X)$ contains $\text{First}(X_{i+1}) - \{\epsilon\}$. If all the sets $\text{First}(X_1), \dots, \text{First}(X_n)$ contain ϵ , then $\text{First}(X)$ also contains ϵ .

Now define **First**(α) for any string $\alpha = X_1 X_2 \dots X_n$ (a string of terminals and non-terminals), as follows. $\text{First}(\alpha)$ contains $\text{First}(X_i) - \{\epsilon\}$. For each $i = 2, \dots, n$, if $\text{First}(X_k)$ contains ϵ for all $k = 1, \dots, i - 1$, then $\text{First}(\alpha)$ contains $\text{First}(X_i) - \{\epsilon\}$. Finally, if for all $i = 1, \dots, n$, $\text{First}(X_i)$ contains ϵ , then $\text{First}(\alpha)$ contains ϵ .

```
for all nonterminals A do First(A) := {}  
while there are changes to any First(A) do  
  for each production choice  $A \rightarrow X_1 X_2 \dots X_n$  do  
    k := 1; Continue := true;  
    while Continue = true and k <= n do  
      add First(Xk) - {ε} to First(A);  
      if ε is not in First(Xk) then Continue := false;  
      k := k + 1;  
    if Continue = true then add ε to First(A);
```

Algoritme:

- Gjør steg 1.
- Gjenta steg 2 til alle Firstmengdene har stabilisert seg.

Eks. 4.9 Beregning av First-mengde

- (1) $exp \rightarrow exp \text{ addop } term$
- (2) $exp \rightarrow term$
- (3) $addop \rightarrow +$
- (4) $addop \rightarrow -$
- (5) $term \rightarrow term \text{ mulop } factor$
- (6) $term \rightarrow factor$
- (7) $mulop \rightarrow *$
- (8) $factor \rightarrow (exp)$
- (9) $factor \rightarrow \mathbf{number}$

Grammar rule	Pass 1	Pass 2	Pass 3
$exp \rightarrow exp \text{ addop } term$			
$exp \rightarrow term$			First(exp) = { (, number }
$addop \rightarrow +$	First($addop$) = { + }		
$addop \rightarrow -$	First($addop$) = { +, - }		
$term \rightarrow term \text{ mulop } factor$			
$term \rightarrow factor$		First($term$) = { (, number }	
$mulop \rightarrow *$	First($mulop$) = { * }		
$factor \rightarrow (exp)$	First($factor$) = { (}		
$factor \rightarrow \mathbf{number}$	First($factor$) = { (, number }		

Kan bare fylle ut en tabell

	First
exp	
addop	
term	
mulop	
factor	

Beregning av Followmengder

$\beta \gamma$ virkelige strenger

Def Follow(A) = { a | finnes avledning $S \$ \Rightarrow^* \beta A a \gamma$ }

Dvs. Det finnes en utledning fra startsymbolet S hvor A kommer rett før a (a er en terminal)

Given a nonterminal A, the set Follow(A), consisting of terminals, and possibly \$, is defined as follows.

1. If A is the start symbol, then \$ is in Follow(A).
2. If there is a production $B \rightarrow \alpha A \gamma$, then First(γ) - { ϵ } is in Follow(A).
3. If there is a production $B \rightarrow \alpha A \gamma$ such that ϵ is in First(γ), then Follow(A) contains Follow(B).

~~Follow(start-symbol) := { \$ } ;
for all nonterminals A \neq start-symbol do Follow(A) := { } ;
while there are changes to any Follow sets do
for each production $A \rightarrow X_1 X_2 \dots X_n$ do
for each X_i that is a nonterminal do
add First($X_{i+1} X_{i+2} \dots X_n$) - { ϵ } to Follow(X_i)
(* Note: if $i=n$, then $X_{i+1} X_{i+2} \dots X_n = \epsilon$ *)
if ϵ is in First($X_{i+1} X_{i+2} \dots X_n$) then
add Follow(A) to Follow(X_i)~~

Algoritme:

- Gjør steg 1.
- Gjenta steg 2 og 3 til alle Follow-mengdene har stabilisert seg.

Beregning av Followmengder 4.12

Follow(exp) = { \$, +, -,) }
Follow($addop$) = { (, **number** }
Follow($term$) = { \$, +, -, *,) }
Follow($mulop$) = { (, **number** }
Follow($factor$) = { \$, +, -, *,) }

- (1) $exp \rightarrow exp \text{ addop } term$
- (2) $exp \rightarrow term$
- (3) $addop \rightarrow +$
- (4) $addop \rightarrow -$
- (5) $term \rightarrow term \text{ mulop } factor$
- (6) $term \rightarrow factor$
- (7) $mulop \rightarrow *$
- (8) $factor \rightarrow (exp)$
- (9) $factor \rightarrow \mathbf{number}$

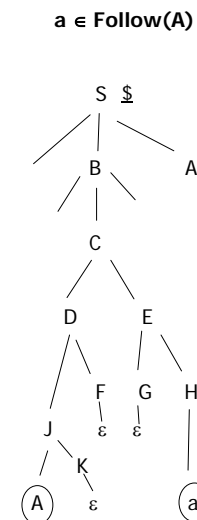
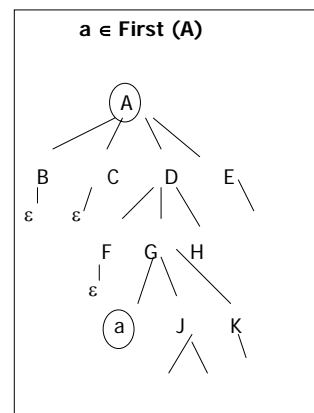
First(exp) = { (, **number** }
First($term$) = { (, **number** }
First($factor$) = { (, **number** }
First($addop$) = { +, - }
First($mulop$) = { * }

Grammar rule	Pass 1	Pass 2
$exp \rightarrow exp \text{ addop } term$	Follow(exp) = { \$, +, - }	Follow($term$) = { \$, +, -, *,) }
$exp \rightarrow term$		
$term \rightarrow term \text{ mulop } factor$	Follow($term$) = { \$, +, -, * }	Follow($factor$) = { \$, +, -, *,) }
$term \rightarrow factor$		
$factor \rightarrow (exp)$	Follow(exp) = { \$, +, -,) }	

Kan bare fylle ut en tabell

	Follow
exp	
addop	
term	
mulop	
factor	

Generelt om situasjonen for neste symbol i First- eller Follow-mengden

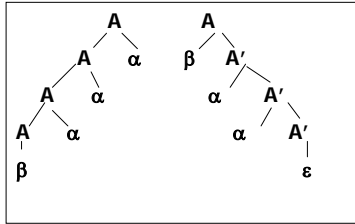


Fjerning av venstre-rekursjon - case 1

Kan skrives på EBNF

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

$$A \rightarrow \beta \{ \alpha \}$$



$$A \rightarrow \beta A'$$

$$A' \rightarrow \alpha A' \mid \epsilon$$

Eksempel:

$$exp \rightarrow exp \text{ addop } term \mid term$$

This is of the form $A \rightarrow A\alpha \mid \beta$, with $A = exp$, $\alpha = \text{addop } term$, and $\beta = term$.
Rewriting this rule to remove left recursion, we obtain

$$exp \rightarrow term \exp'$$

$$\exp' \rightarrow \text{addop } term \exp' \mid \epsilon$$

Fjerning av venstre-rekursjon

Case 2:

$$A \rightarrow A\alpha_1 \mid A\alpha_2 \mid \dots \mid A\alpha_n \mid \beta_1 \mid \beta_2 \mid \dots \mid \beta_m$$



$$A \rightarrow \beta_1 A' \mid \beta_2 A' \mid \dots \mid \beta_m A'$$

$$A' \rightarrow \alpha_1 A' \mid \alpha_2 A' \mid \dots \mid \alpha_n A' \mid \epsilon$$

Case 3 (indirekte venstre-rekursjon)

$$A \rightarrow Ba \mid Aa \mid c$$

$$B \rightarrow Bb \mid Ab \mid d$$

- Her er det feil i boka i algoritmen s.159
- Ikke med som pensum

$$A \rightarrow BaA' \mid cA'$$

$$A' \rightarrow aA' \mid \epsilon$$

$$B \rightarrow Bb \mid Ab \mid d$$

$$A \rightarrow BaA' \mid cA'$$

$$A' \rightarrow aA' \mid \epsilon$$

$$B \rightarrow Bb \mid BaA'b \mid cA'b \mid d$$

$$A \rightarrow BaA' \mid cA'$$

$$A' \rightarrow aA' \mid \epsilon$$

$$B \rightarrow cA'bB' \mid dB'$$

$$B' \rightarrow bB' \mid aA'bB' \mid \epsilon$$

Eks side 160 – Fjerning av venstre-rekursjon

Den tradisjonelle entydige grammatikken

$$exp \rightarrow exp \text{ addop } term \mid term$$

$$addop \rightarrow + \mid -$$

$$term \rightarrow term \text{ mulop } factor \mid factor$$

$$mulop \rightarrow *$$

$$factor \rightarrow (exp) \mid \mathbf{number}$$

Med fjernet venstre-rekursjon

$$exp \rightarrow term \exp'$$

$$\exp' \rightarrow \text{addop } term \exp' \mid \epsilon$$

$$addop \rightarrow + \mid -$$

$$term \rightarrow factor \text{ term}'$$

$$\text{term}' \rightarrow \text{mulop } factor \text{ term}' \mid \epsilon$$

$$mulop \rightarrow *$$

$$factor \rightarrow (exp) \mid \mathbf{number}$$

Venstre-faktorisering

Problem Ønsker ikke alternativer som starter likt:

$A \rightarrow abc B \mid ab C \mid DE$

Løsning: $A \rightarrow ab A' \mid DE$

$A' \rightarrow c B \mid C$

<p>while there are changes to the grammar do</p> <p>for each nonterminal A do</p> <p> let α be a prefix of maximal length that is shared by two or more production choices for A</p> <p>if $\alpha \neq \epsilon$ then</p> <p> let $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ be all the production choices for A and suppose that $\alpha_1, \dots, \alpha_k$ share α, so that $A \rightarrow \alpha \beta_1 \mid \dots \mid \alpha \beta_k \mid \alpha_{k+1} \mid \dots \mid \alpha_n$, the β_j's share no common prefix, and the $\alpha_{k+1}, \dots, \alpha_n$ do not share α</p> <p> replace the rule $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$ by the rules</p> <p> $A \rightarrow \alpha A' \mid \alpha_{k+1} \mid \dots \mid \alpha_n$</p> <p> $A' \rightarrow \beta_1 \mid \dots \mid \beta_k$</p>	<p>} Gjenta så lenge det er muligheter</p> <p>} Velg lengst mulig prefix</p> <p>} Gjør som over. Pass på å få med <u>alle</u> alternativer med dette prefixet</p>
--	---

eks1:

$stmt\text{-}sequence \rightarrow stmt ; stmt\text{-}sequence \mid stmt$
 $stmt \rightarrow s$

$stmt\text{-}sequence \rightarrow stmt stmt\text{-}seq'$
 $stmt\text{-}seq' \rightarrow ; stmt\text{-}sequence \mid \epsilon$

Litt om stoffet i kap. 4

- First og Follow-mengder
 - Boka tar det sist, vi først
- Begrepet LL(1) – parsering
 - Boka reserverer dette begrepet for den metoden (kap 4.2 s. 152) der man bruker en eksplisitt stakk (isteden for rekursive metoder, som i "recursive decent" – framgangsmåten)
 - Det vanlige er å bruke betegnelsen LL(1) –parsering også når rec.decent metoden brukes slavisk på en ren BNF-grammatikk
 - Kravet til en LL(1) – grammatikk kommer like tydelig fram ut fra begge disse metodene. Vi skal tenke like mye rec.decent som eksplisitt stakk. (boka gjør bare det siste).
 - Ofte brukes betegnelsen LL(1) – parsering også om rec.decent metoden brukt ut fra syntaksdiagrammer eller EBNF, men da er det uklart hva LL(1) – kravet til en grammatikk er.

"Rekursiv nedstigning " – parsering

$exp \rightarrow exp\ addop\ term \mid term$
 $addop \rightarrow + \mid -$
 $term \rightarrow term\ mulop\ factor \mid factor$
 $mulop \rightarrow *$
 $factor \rightarrow (exp) \mid \mathbf{number}$

- Skriv en funksjon /prosedyre/metode for hver ikke-terminal
- La denne lese/sjekk høyreside-alternativene

token : if ← Global variabel

"Typisk"
 rec.dec.
 prosedyre
 for det
 enkle'
 tilfellet.

```

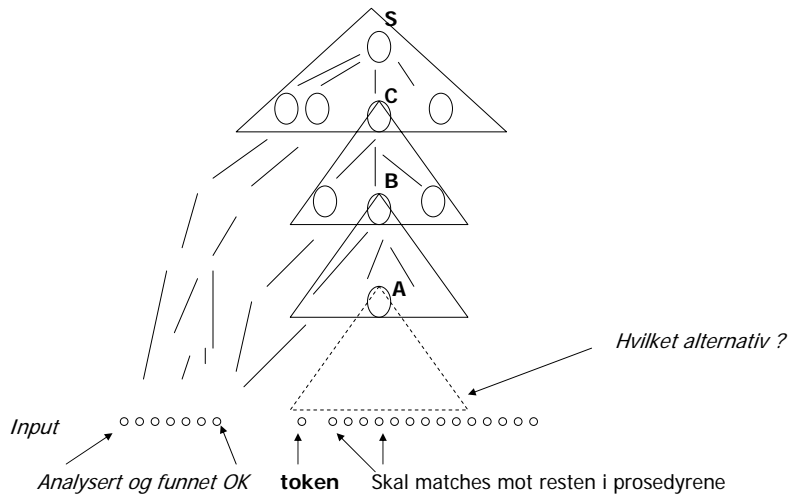
procedure factor ;
begin
  case token of
    (: match( ) ;
      exp ;
      match( ) ) ;
    number :
      match(number) ;
    else error ;
  end case ;
end factor ;
  
```

Sjekker at angitt terminal kommer, og "leser til neste". Brukes ofte bare for å lese (sjekken må slå til).

```

procedure match ( expectedToken ) ;
begin
  if token = expectedToken then
    getToken ;
  else
    error ;
  end if ;
end match ;
  
```

Situasjonen under rekursiv parsering



Litt mer kompliserte tilfeller kan løses med EBNF

Opprinnelig

$if-stmt \rightarrow \mathbf{if} (exp) statement [\mathbf{else} statement]$

Skrives ut som:

$if-stmt \rightarrow \mathbf{if} (exp) statement$
 $\quad \quad \quad | \mathbf{if} (exp) statement \mathbf{else} statement$

R-D-prosedyre:

```

procedure ifStmt ;
begin
  match (if) ;
  match ( ( ) ;
  exp ;
  match ( ) ;
  statement ;
  if token = else then
    match (else) ;
    statement ;
  end if ;
end ifStmt ;
    
```

N.B.: Kunne også bruke venstre-faktorisering. Da ville dette bli en egen prosedyre "elsePart":

$ifStmt \rightarrow \mathbf{if} (exp) stmt \mathbf{elsePart}$
 $\mathbf{elsePart} \rightarrow \epsilon \mid \mathbf{else} stmt$

Venstre-rekursjon gir problemer

Gir uendelig mange rekursive kall

N.B. Kan også fjernevenstre-rekursjon på trad måte, se senere foil.

$exp \rightarrow exp \text{ addop } term \mid term$

Bruker EBNF:

$exp \rightarrow term \{ \text{ addop } term \}$

$term \rightarrow term \text{ multop } factor \mid factor$

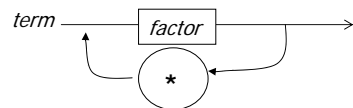
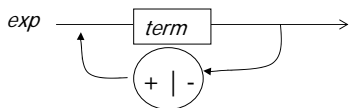
$term \rightarrow factor \{ \text{ multop } factor \}$

```

procedure exp ;
begin
  term ;
  while token = + or token = - do
    match (token) ;
    term ;
  end while ;
end exp ;
    
```

```

procedure term ;
begin
  factor ;
  while token = * do
    match (token) ;
    factor ;
  end while ;
end term ;
    
```



Hvordan "lage noe" under rec.-decent parsing?

- Mål: Ønsker å bygge abstrakt syntaks-tre
- Men foreløpig:
 - beregner verdien av et uttrykk (med venstre-assosiativitet)

```

function exp : integer ;
var temp : integer ;
begin
  temp := term ;
  while token = + or token = - do
    case token of
      + : match (+) ;
          temp := temp + term ;
      - : match (-) ;
          temp := temp - term ;
    end case ;
  end while ;
  return temp ;
end exp ;
    
```

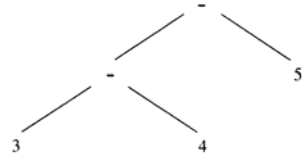
Kall

- Kan lett bygges ut til full "kalkulator"

3 + 4 + 5

Bygging av syntaks-treet

3-4-5:



```
function exp : syntaxTree ;
var temp, newtemp : syntaxTree ;
begin
temp := term ;
while token = + or token = - do
case token of
+ : match (+) ;
newtemp := makeOpNode(+);
leftChild(newtemp) := temp ;
rightChild(newtemp) := term ;
temp := newtemp ;
- : match (-) ;
newtemp := makeOpNode(-);
leftChild(newtemp) := temp ;
rightChild(newtemp) := term ;
temp := newtemp ;
end case ;
end while ;
return temp ;
end exp ;
```

Alternativt:
newtemp.leftChild

Kall

Kall

Merk: Dersom det bare er en "term", så lages ingen ny node. Vi leverer den vi har fått

Skriue prosedyre for:

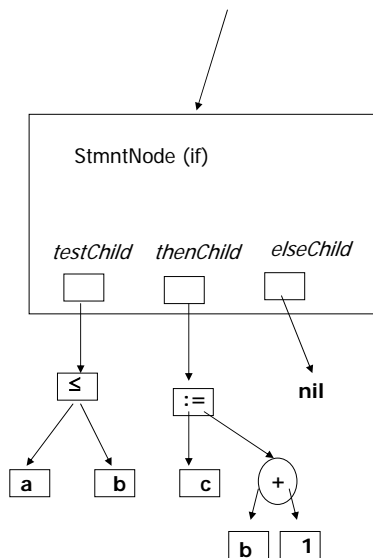
$factor \rightarrow (exp) | \underline{number}$

```
proc factor : syntaxTree;
var fact: syntaxTree;
begin
case token of
( : match ( ( ) ;
fact = exp ;
match ( ) ;
number :
match ( number ) ;
fact =makeNumberNode( number ) ;
else error(.....) ;
end case
return fact;
end factor;
```

Node for generering for if-setning

if-stmt -> if (exp) stmt [else tstmt]

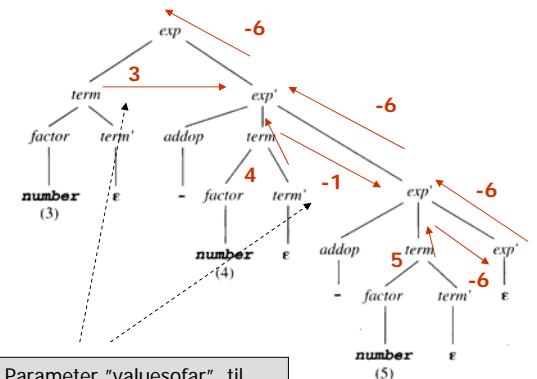
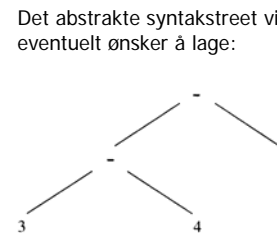
```
function ifStatement : syntaxTree ;
var temp : syntaxTree ;
begin
match (if) ;
match ( ( ) ;
temp := makeStmtNode(if) ;
testChild(temp) := exp ;
match ( ) ;
thenChild(temp) := statement ;
if token = else then
match (else) ;
elseChild(temp) := statement ;
else
elseChild(temp) := nil ;
end if ;
end ifStatement ;
```



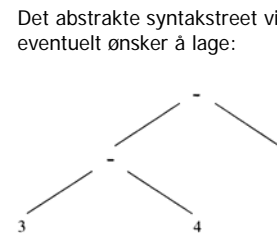
Rec.decent etter tradisjonell fjerning av venstre-rekursjon (treet er nå høyre assiativt istedenfor venstre). Det må rettes opp.

Lage tre eller beregne verdi : 3 - 4 - 5

```
exp → term exp'
exp' → addop term exp' | ε
addop → + | -
term → factor term'
term' → mulop factor term' | ε
mulop → *
factor → ( exp ) | number
```



Parameter "valuesofar" til prosedyren "exp"
(For trebygging ville den være: "rootOfTreeSoFar")



Prosedyrer for beregning av verdi (tre-bygging tilsvarende)

```

exp → term exp'
exp' → addop term exp' | ε
addop → + | -
term → factor term'
term' → mulop factor term' | ε
mulop → *
factor → ( exp ) | number
    
```

Bare analyse

```

procedure exp ;
begin
  term ;
  exp' ;
end exp ;
    
```

```

procedure exp' ;
begin
  case token of
  + : match (+) ;
    term ;
    exp' ;
  - : match (-) ;
    term ;
    exp' ;
  end case ;
end exp' ;
    
```

Med beregning (trebygging tilsvarende)

```

function exp : integer ;
var temp : integer ;
begin
  temp := term ;
  return exp'(temp) ;
end exp ;
    
```

NB.: Parameter

```

function exp' ( valsofar : integer ) : integer ;
begin
  if token = + or token = - then
  case token of
  + : match (+) ;
    valsofar := valsofar + term ;
  - : match (-) ;
    valsofar := valsofar - term ;
  end case ;
  return exp'(valsofar) ; ← ε -alternativet
  else return valsofar ;
end exp' ;
    
```

ε -alternativet
Leverer verdien uendret oppover igjen.

LL(1) – grammatikk

$$A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3 \mid \dots \mid \alpha_n$$

- LL(1) -kravet for en "ren BNF-grammatikk". Det som kreves for at en Rek. desc.-parsering skal fungere direkte fra grammatikken.
- For å avgjøre om en grammatikk er "LL(1)": Sett opp tabell $M[N,T]$ med mulige aksjoner for alle mulige situasjoner, slik:

1. Altså, $a \in \text{First}(\alpha)$

- If $A \rightarrow \alpha$ is a production choice, and there is a derivation $\alpha \Rightarrow^* a \beta$, where a is a token, then add $A \rightarrow \alpha$ to the table entry $M[A, a]$.
- If $A \rightarrow \alpha$ is a production choice, and there are derivations $\alpha \Rightarrow^* \epsilon$ and $S \Rightarrow^* \beta A a \gamma$, where S is the start symbol and a is a token (or $\$$), then add $A \rightarrow \alpha$ to the table entry $M[A, a]$.

2. Altså, dersom:
- α er utnullbar, og
 - $a \in \text{Follow}(A)$

Definisjon:
En grammatikk er LL(1) dersom $M[N,T]$ er entydig for alle situasjoner (eller angir "error")

Oppsett av LL(1) –tabell

```

statement → if-stmt | other
if-stmt → if ( exp ) statement else-part
else-part → else statement | ε
exp → 0 | 1
    
```

- Venstre-faktorisering utført
- Er ikke vestrerekursiv

	First	Follow
statement	other, if	\$, else
if-stmt	if	\$, else
else-part	else, ε	\$, else
exp	0, 1)

$M[N, T]$	if	other	else	0	1	\$
statement	statement → if-stmt	statement → other				
if-stmt	if-stmt → if (exp) statement else-part					
else-part			else-part → else statement else-part → ε			else-part → ε
exp				exp → 0	exp → 1	

Merk:

- fjernet venstre-rek.
- Utført venstre-fakt.
- er **ikke** nok til å garantere LL(1)-grammatikk.
- fordi tabellen er ikke entydig**