

Dagens tema:

Kapittel 4

Parsering ovenfra-ned (top-down)

Vi har med alle foilene til kap. 4 her, også de som ble gjennomgått mot slutten av forelesning 7. februar

Pensum er korta ned litt i forhold til det i 2013



Kap. 4: Parsering ovenfra-ned (top-down)

Dette bør leses om igjen etter kapittel 4:

- *First* og *Follow*-mengder
 - Boka tar det et stykke uti kap 4, vi tok det først (forrige gang)
- *LL(1)-parsering* og boka og pensum:
 - Det som i *boka* kalles LL(1)-parsering (4.2.1, 4.2.2 og 4.2.4) er en metode for top-down parsering med en eksplisitt stakk. Denne metoden er *ikke* med som pensum.
 - Vi konsentrerer oss i dette kapittelet om "recursive descent"-parsering, en intuitiv metode som mange sikkert har vært litt borti (INF2100, ++)
 - Ofte brukes betegnelsen LL(1)-parsering også om "rec. desc."-metoden brukt ut fra syntaksdiagrammer, EBNF eller ren BNF, men man har da gjerne et ikke-teknisk forhold til det om ting helt sikkert fungerer riktig.
 - Vi skal se på det tekniske kravet for at en ren BNF-grammatikk kan parseres rett fram med "rec. desc."-metoden.
 - **LL(1)-grammatikk i vår betydning:**
Det er en **ren BNF-grammatikk** som tilfredstiller kravet fra forrige punkt

Vi skal først og fremst se på metoden «Recursive descent» («Rekursiv parsing»??)

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

Neste token \rightarrow \leftarrow «Global» variabel

if n + n * (n + n) <= ...

“Typisk”
rec.decent-
prosedyre
for siste
produksjon
over, som
blir veldig
enkel.

```
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). Da er det
egentlig nok å
kalle “getToken”
(men her kalles
alltid “match”)

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

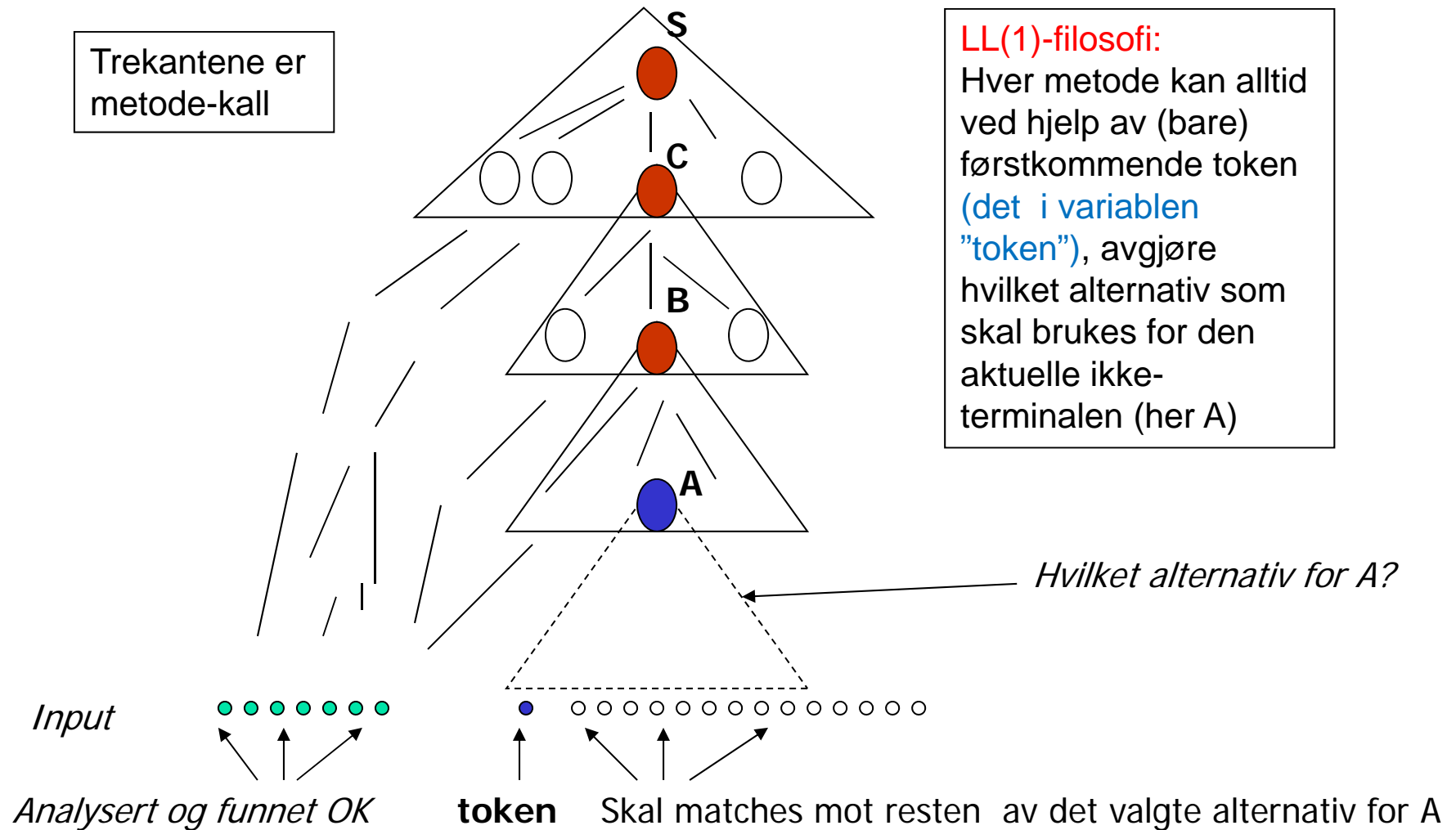
Hoved-idé:

- Skriv en funksjon/ prosedyre/metode for hver ikke-terminal
- La denne finne **riktig alternativ (helst fra bare variabelen “token”)**, og gjør parsing av den videre input ut fra det
- Og sjekk på veien at alt er OK

Situasjonen under rekursiv parsing

Kalles altså "top down"-parsing (ovenfra-ned-parsing)

Trekantene er metode-kall





I mer kompliserte tilfeller virker ofte ikke ren BNF bra, men med venstrefaktorisering eller EBNF går det her ofte greit

Opprinnelig:

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

Kan i EBNF skrives ut som:

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

R-D-prosedyre:

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

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

ifStmt \rightarrow if (exp) stmt elsePart
elsePart \rightarrow ϵ | else stmt

Venstre-rekursjon gir problemer ved ren BNF. Men går ofte greit med EBNF

Gir uendelig mange rekursive kall

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

NB: Kan også fjerne venstre-rekursjon på trad. måte, se neste foil.

Bruker EBNF:

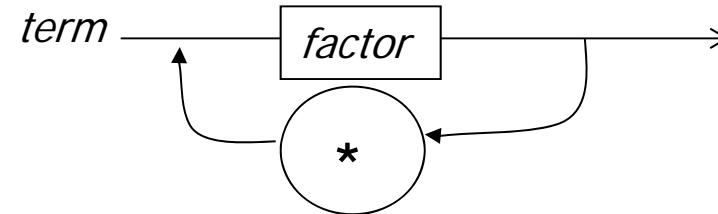
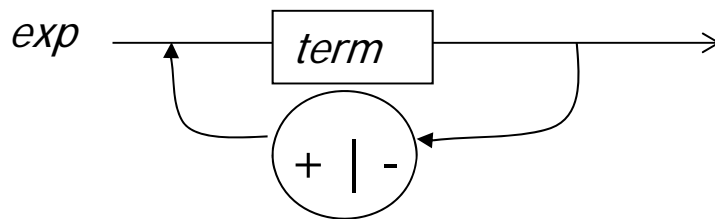
$exp \rightarrow term \{ \text{addop } term \}$

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

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

```
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 ;
```



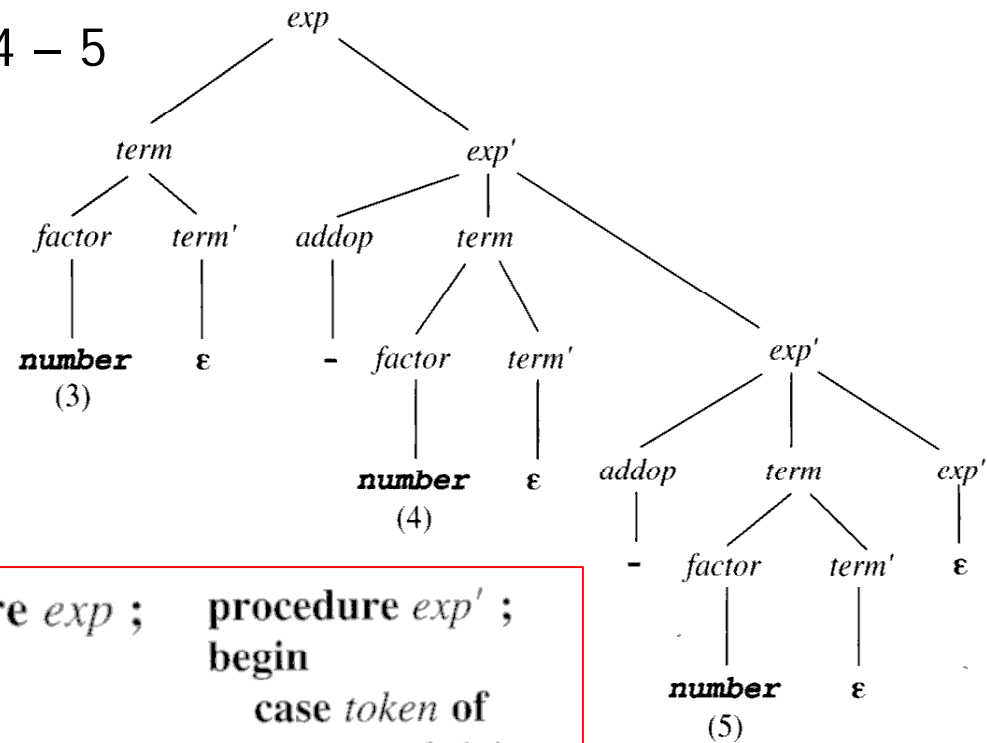
Rec.decent etter tradisjonell fjerning av venstre-rekursjon

Treet blir nå høyre-assosiativt istedenfor venstre.

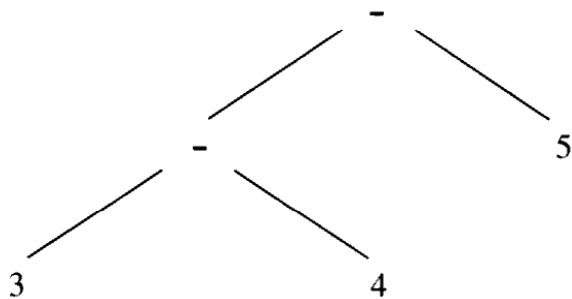
Dette må det korrigeres for, men det er ikke pensum!

Uttrykk: 3 - 4 - 5

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



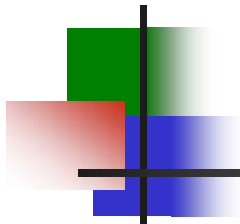
Det abstrakte syntakstreet vi egentlig ønsket å lage:



```

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

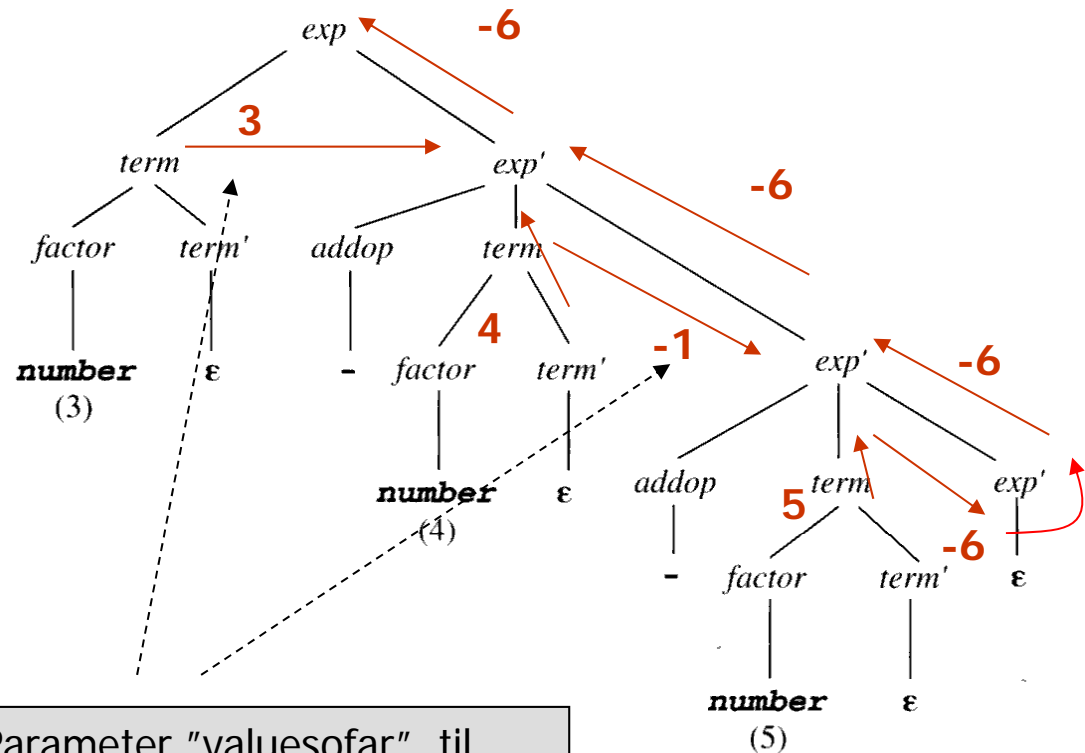
procedure exp' ;
begin
    case token of
      + : match (+) ;
        term ;
        exp' ;
      - : match (-) ;
        term ;
        exp' ;
    end case ;
end exp' ;
    
```



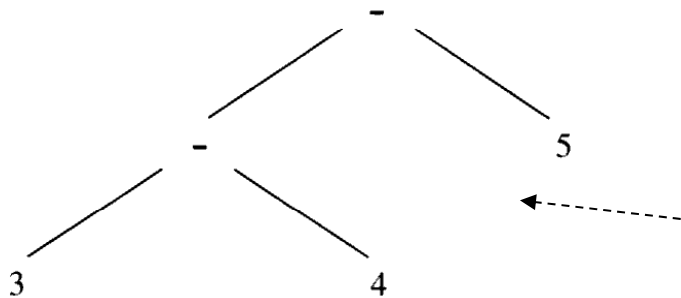
Ikke pensum: Rec.decent etter tradisjonell fjerning av venstre-rekursjon (treet er nå høyre-assosiativt istedenfor venstre).

Lage tre eller beregne verdi : 3 - 4 - 5

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



Det abstrakte syntakstreet vi ønsker å lage:



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

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

- Mål: Ønsker å bygge abstrakt syntaks-tre
- Tar utgangspunkt i foil 6 (uten nummer)
- Men foreløpig (som kan være forvirrende!):
 - 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"

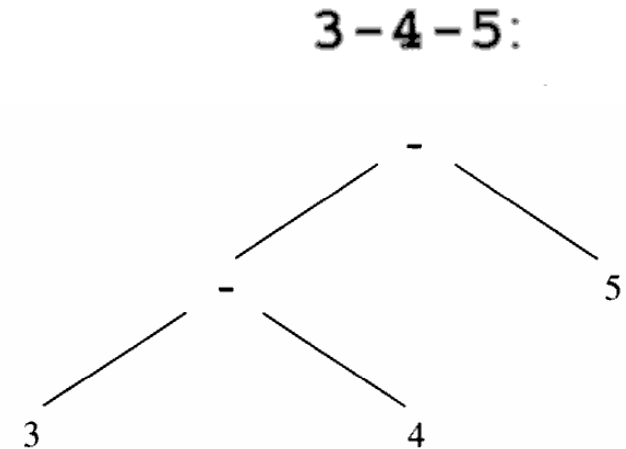
3 + 4 + 5

Bygging av abstrakt syntaks-tre

Igjen med utgangspunkt i foil 6

```
function exp : syntaxTree ;  
var temp, newtemp : syntaxTree ;  
begin  
  temp := term ; ← NB: Kall  
  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



NB: Kall

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



Prosedyre med trebygging for:

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

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

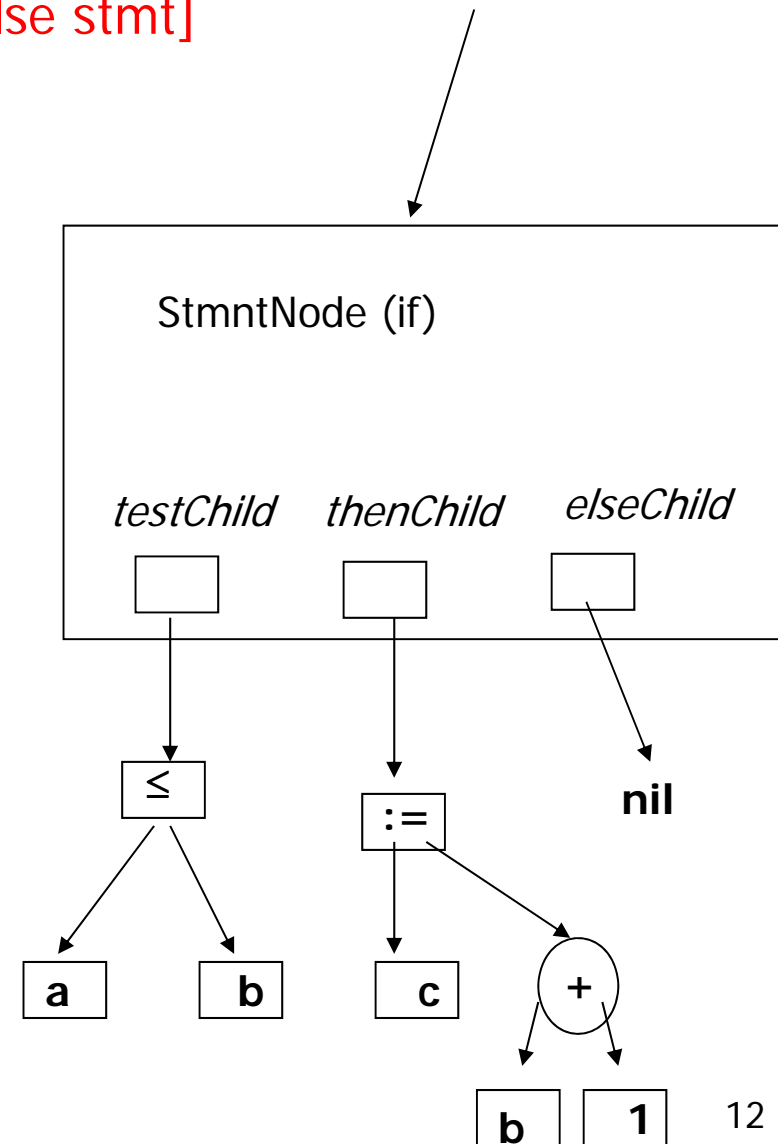
Gir "dummy-test"



Parsering av if-setning, med tre-generering

EBNF: **if-stmt** -> **if (exp) stmt [else stmt]**

```
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 ;
```



Viktig detalj om hva som skal stå i «token» ved inngang og utgang av de rekursive metodene

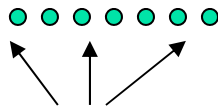
Som før:

Trekantene er metode-kall.

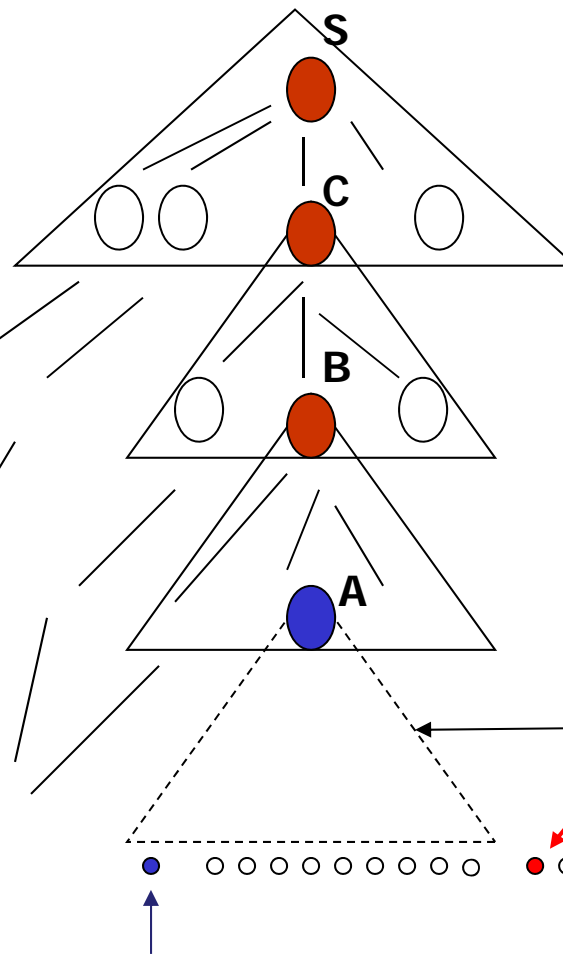
LL(1)-filosofi:

Hver metode kan alltid ved hjelp av (bare) førstkomende token, avgjøre hvilket alternativ som skal brukes for den aktuelle ikke-terminalen (her A)

Input



Analysert og funnet OK



Det valgte alternativet for A

Dette skal ligge i «token» når metoden for A kalles

Regel for hva som skal ligge i «token» ved inngang og utgang av en syntaks-metode (si for A):

- Ved *inngang* skal første symbolet i den syntaktiske konstruksjonen A ligge i «token»
- Ved *utgang* skal første symbolet *etter* konstruksjonen A ligge i «token»

Dette skal ligge i «token» når metoden for A returnerer

LL(1) – grammatikk

- LL(1) -kravet for en "ren BNF-grammatikk":

Det som kreves for at en rek. descent-parsing skal fungere direkte fra grammatikken uten omskrivninger.

- For å avgjøre om en grammatikk er "LL(1)": Sett opp tabell $M[N,T]$ med mulige aksjoner for alle mulige situasjoner, slik:

1. 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 $M[A, a]$.
2. If $A \rightarrow \alpha$ is a production choice, and there are derivations $\alpha \Rightarrow^* \varepsilon$ 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 $M[A, a]$.

Definisjon av LL(1):

En grammatikk er LL(1) dersom $M[A,a]$ er entydig for alle situasjoner (eller angir "error")

1. Altså, dersom

$a \in \text{First}(\alpha)$

så legg

$A \rightarrow \alpha$ inn i $M[A, a]$

2. Altså, dersom:

• α er utnullbar, og

• $a \in \text{Follow}(A)$

så legg $A \rightarrow \alpha$ inn i $M[A, a]$

Oppsett av LL(1) –tabell

$statement \rightarrow if-stmt \mid \mathbf{other}$
 $if-stmt \rightarrow \mathbf{if} (exp) statement \mathit{else-part}$
 $\mathit{else-part} \rightarrow \mathbf{else} statement \mid \epsilon$
 $exp \rightarrow \mathbf{0} \mid \mathbf{1}$

- 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: Selv om man:

- fjerner venstre-rek.
- Utfører venstre-fakt.
- er det generelt **ikke** nok til å *garantere* LL(1)-grammatikk.

For tabellen ble ikke entydig her

Ikke Pensum: Kan også være greit å snakke om den "utvidede startmengden", *Efirst*, til en produksjon

NB: Ikke i boka, men kan øke forståelsen!

- Den "utvidede startmengden", *Efirst*, til en *produksjon* $A \rightarrow \alpha$ er *det og bare det* som kan ligge som *førstkommende token* i input dersom $A \rightarrow \alpha$ et riktig valg på dette stadiet under parseringen.
- Her må man også tenke på tilfellet at α kan være utnullbar, og da kan også $\text{Follow}(A)$ komme som *førstkommende token*
- Mengden kan beregnes slik:
 $Efirst(A \rightarrow \alpha) = \text{First}(\alpha)$, pluss, om α er utnullbar, $\text{Follow}(A)$.
- Vi ser da at produksjonen $A \rightarrow \alpha$ skal inn i $M[A, a]$ hvis og bare hvis $a \in Efirst(A \rightarrow \alpha)$
- Dermed, får vi en alternativ definisjon av LL(1):

Anta at: $A \rightarrow \alpha_1 \mid \alpha_2 \mid \alpha_3 \mid \dots \mid \alpha_n$

Da må alle:

$Efirst(A \rightarrow \alpha_1), Efirst(A \rightarrow \alpha_2), \dots, Efirst(A \rightarrow \alpha_n)$
være parvis disjunkte. Merk at dette også innebærer at bare ett av alternativene kan være utnullbare

På forrige side vil da både

$Efirst(\text{else-part} \rightarrow \text{else statmt})$

og

$Efirst(\text{else-part} \rightarrow \epsilon)$

inneholde **else**

Dermed er grammatikken *ikke*

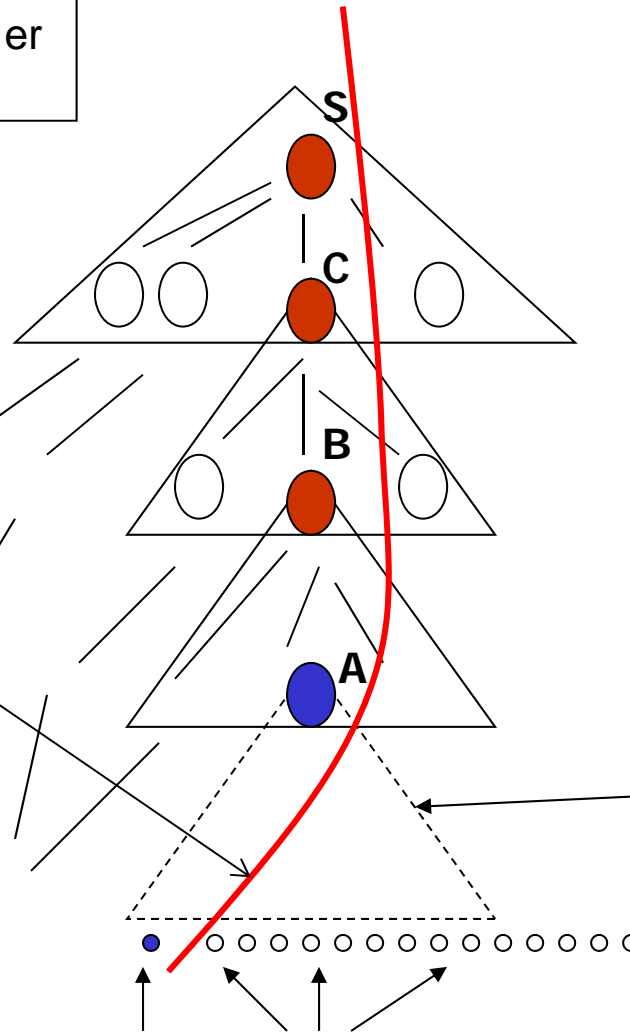
LL(1)

Tidligere foil:

Situasjonen under rekursiv parsing

Trekantene er metode-kall

Det er nå laget ferdig et abstrakt syntakstre for alt til venstre for **den røde linjen**



Merk:

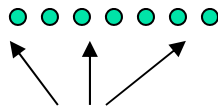
En grammatikk som er flertydig kan ikke være en LL(1)-grammatikk

Altså:

En LL(1)-grammatikk er alltid entydig!

Hvilket alternativ for A?

Input



Analysert og funnet OK

token

Skal matches mot resten av det valgte alternativ for A

LL(1) –tabell for uttrykks-grammatik

Vi har tidligere funnet First og Follow som en oppgave

Har fjernet venstre-rekursjon. Er den nå LL(1)?

Opprinnelig utgave. Ikke LL(1)!

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

$addop \rightarrow + \mid -$

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

$mulop \rightarrow *$

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

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

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

$addop \rightarrow + \mid -$

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

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

$mulop \rightarrow *$

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

Vi får da følgende First- og Follow-mengder:

$\text{First}(exp) = \{ (, \mathbf{number} \}$

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

$\text{First}(exp') = \{ +, -, \varepsilon \}$

$\text{Follow}(exp') = \{ \$,) \}$

$\text{First}(addop) = \{ +, - \}$

$\text{Follow}(addop) = \{ (, \mathbf{number} \}$

$\text{First}(term) = \{ (, \mathbf{number} \}$

$\text{Follow}(term) = \{ \$,), +, - \}$

$\text{First}(term') = \{ *, \varepsilon \}$

$\text{Follow}(term') = \{ \$,), +, - \}$

$\text{First}(mulop) = \{ * \}$

$\text{Follow}(mulop) = \{ (, \mathbf{number} \}$

$\text{First}(factor) = \{ (, \mathbf{number} \}$

$\text{Follow}(factor) = \{ \$,), +, -, * \}$

I

$M[N, T]$	(number)	+	-	*	\$
<i>exp</i>	$exp \rightarrow$ <i>term exp'</i>	$exp \rightarrow$ <i>term exp'</i>					
<i>exp'</i>			$exp' \rightarrow \epsilon$	$exp' \rightarrow$ <i>addop</i> <i>term exp'</i>	$exp' \rightarrow$ <i>addop</i> <i>term exp'</i>		$exp' \rightarrow \epsilon$
<i>addop</i>				$addop \rightarrow$ +	$addop \rightarrow$ -		
<i>term</i>	$term \rightarrow$ <i>factor</i> <i>term'</i>	$term \rightarrow$ <i>factor</i> <i>term'</i>					
<i>term'</i>			$term' \rightarrow$ ϵ	$term' \rightarrow \epsilon$	$term' \rightarrow \epsilon$	$term' \rightarrow$ <i>mulop</i> <i>factor</i> <i>term'</i>	$term' \rightarrow$ ϵ
<i>mulop</i>						$mulop \rightarrow$ *	
<i>factor</i>	$factor \rightarrow$ (<i>exp</i>)	$factor \rightarrow$ number					



Når kompilatoren oppdager feil

Vi ser på denne og den neste foilen ved slutten av kap. 5

- Minstekrav:
 - Tester løpende at programmet er OK, og gir fornuftig feilmelding ved feil (men stopper kanskje)
- Vanlig krav ved feil ("error recovery"):
 - Gir fornuftig feilmelding.
 - "Blar forbi feilen" og fortsetter kompileringen (og blar forbi så lite som mulig).
 - Vanligvis vil man slutte å lage maskinkode etter feil (Men noe "feilrettende" kompilatorer forsøker det – lite brukt)
 - Det er etter *syntaksfeil* det er vanskeligst å ta opp tråden igjen.



Feilmeldinger m.m.

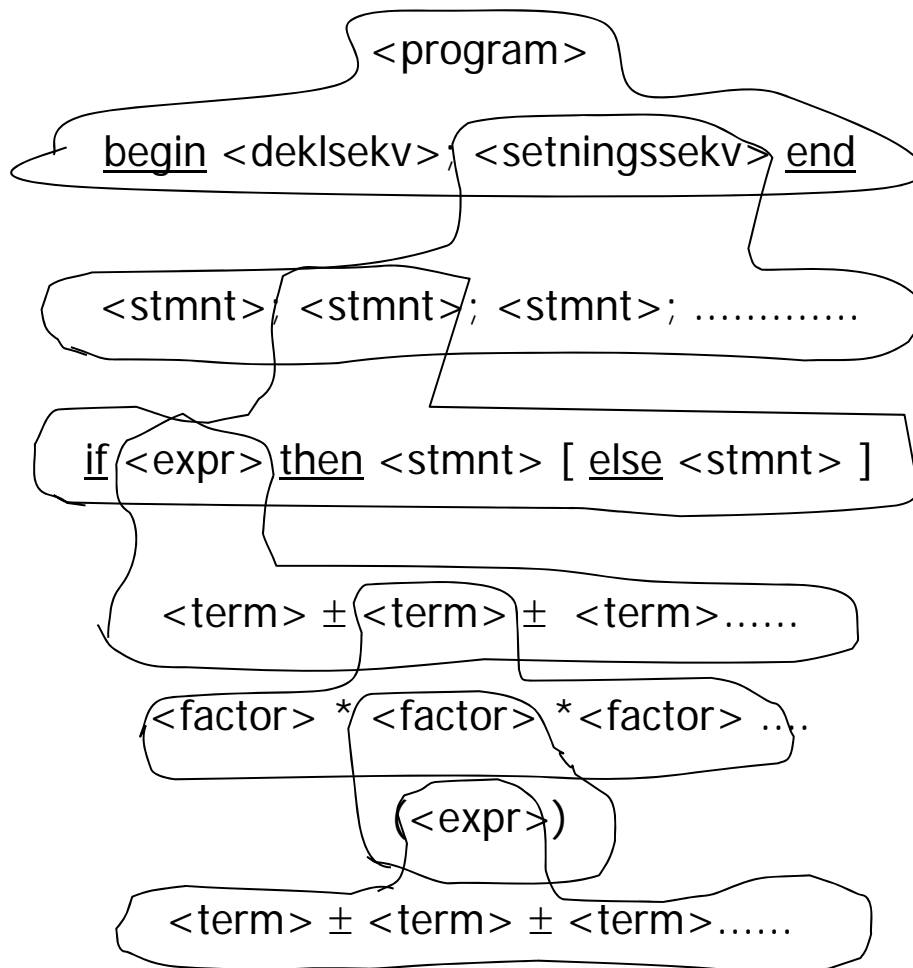
Ser på denne i forbindelse med Kap. 5

- Viktig:
 - Forsøke å unngå feilmeldinger som bare er følgefeil
 - Rapportere feil så tidlig som mulig, helst så snart det man har lest *ikke kan forlenges til et riktig program*
 - Man må passe på at man ikke blir gående i løkke *uten å lese noe fra input* (verken med eller uten rapportering feil)
- Hvilken feilmelding skal man gi?
 - Anta at man ved "factor" har valgt alternativet " (exp) ", og at man, når exp-metoden returnerer, ikke finner en ") ".
 - Hvilken feilmelding skal man da gi??
 - Man kan melde: "Sluttparentes mangler"
 - Men dette kan ofte være forvirrende, f.eks. om det som stod var: (a + b c).
 - Her vil exp-metoden stoppe etter " (a + b ", og det vil være forvirrende her med feilmeldingen "Sluttparentes mangler"
 - Man bør derfor heller gi meldingen: "Noe galt i et uttrykk, eller sluttparentes mangler"

Behandling av Syntaksfeil

"recursive decent". Ikke pensum 2014.

Metode: "Panic mode" og synkroniserings-mengde



Synch-set (stakk eller parameter):

\$

end

; First(stmnt)

navn if while for ...

then First(stmnt) else

+ - First(term)

(tall navn

* First(factor)

)

+ - (tall navn



Syntaksfeil ved "recursive descent"

Ikke pensum 2014.

Ut fra skissen på forrige foil er det greit å finne:

- hvem som skal ta opp tråden
- "hvor" denne skal fortsette eksekveringen

Vi antar at \$ bare legges på stakken av start-symbol-metoden
Unionen av alle på stakken kalles "synkroniseringsmengden", SM

Algoritme:

For hvert input-symbol framover, test om det er med i SM

I så fall:

- Let gjennom SM-stakken fra det siste vi la på, og finn den metoden som
 - sist ble kalt,
 - og som kan ta opp tråden på dette input-symbolet
- Denne metoden vet selv hvor den skal fortsette, ut fra input-symbolet

Det som *ikke er greit*, er å programmere dette *uten at den vakre strukturen* ved "rec. descent" blir helt ødelagt.

Spørsmål, som foreleseren ikke har gått til bunns i:

Kan man få til noe bra ved å bruke Javas unntaksmekanisme??

Uttrykksprosedyrer ved "error recovery". Ikke pensum 2014

Filosofien her er litt annerledes (og noe uklar?)

```
procedure exp ( synchset ) ;
begin
  checkinput ( { (, number }, synchset ) ;
  if not ( token in synchset ) then
    term ( synchset ) ;
    while token = + or token = - do
      match ( token ) ;
      term ( synchset ) ;
    end while ;
    checkinput ( synchset, { (, number } ) ;
  end if ;
end exp ;
```

Også { +, - } ?

Hovedfilosofi

"checkinput" kalles to ganger: Først for å sjekke at konstruksjonen starter riktig, etterpå for å sjekke at symbolet etter konstruksjonen er lovlig.

Bruker parameter, ikke stakk

Prosedyrene må selv ta opp tråden riktig når de får igjen kontrollen:

match(t) er som før:

- tester input mot t
- kaller eventuelt "error" (som nå returnerer!)
- kaller ikke "scanto(...)"

if token in {(,number} then ...

```
procedure factor ( synchset ) ;
begin
  checkinput ( { (, number }, synchset ) ;
  if not ( token in synchset ) then
    case token of
      ( : match ( ( ) ;
        exp ( { } ) ) ;
        match ( ) ;
      number :
        match ( number ) ;
      else error ;
    end case ;
    checkinput ( synchset, { (, number } ) ; *
  end if ;
end factor ;
```

Hvorfor ikke også "synchset"?

```
procedure scanto ( synchset ) ;
begin
  while not ( token in synchset  $\cup$  { $ } ) do
    getToken ;
  end scanto ;
```

```
procedure checkinput ( firstset, followset ) ;
begin
  if not ( token in firstset ) then
    error ;
    scanto ( firstset  $\cup$  followset ) ;
  end if ;
end;
```