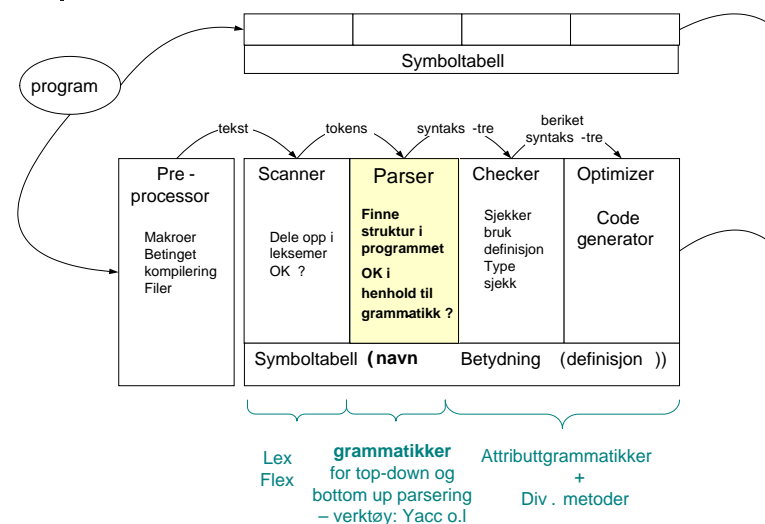


## Kontekstfrie grammatikker og syntaksanalyse (parsing)

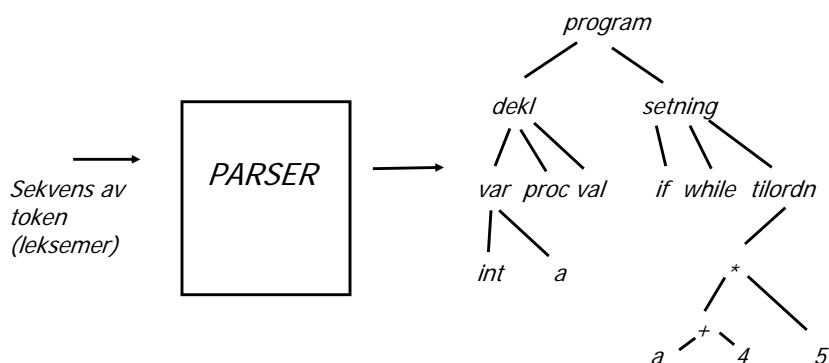
INF5110 - kap.3 i Louden + hjelpenotat (se hjemmesida)

Arne Maus  
Ifi, UiO – v2006

## Hvor er vi nå - kap. 3 (+4,5) ?



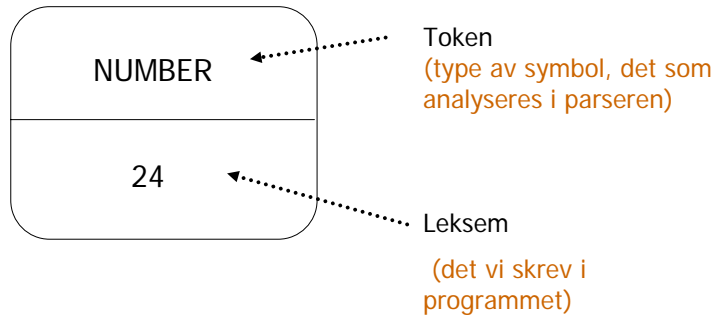
## Forenklet skisse av hva en parser gjør



## Oversikt – kap 3 (grunnleggende om grammatikker)

- Hva er en grammatikk
- Vi kommer til å lære flere (mange) typer grammatikk
  - hvorfor ?
- Kontekstfrie grammatikker
- Parserings-trær og abstrakte syntaks-trær
- Tvetydige grammatikker
- Utvidet notasjon: EBNF og syntaksdiagrammer
- Eksempel
  - Tiny

## Hva får vi fra parseren



Av og til vil vi gi eksempler med token og av og til med leksem

## Hva er en grammatikk

- En grammatikk bruker et alfabet
  - Som oftest det vi kan taste inn fra tastaturet
- En grammatikk består av noen **symboler** = sammensetninger av tegn fra alfabetet:
  - Terminal-symboler** - slike vi bruker når vi programmerer  
Dvs. for parseren er det slike tokens/leksemer vi får fra skanneren:
    - int, TALL, if, VARIABEL, .....
  - Ikke-terminal-symboler** - begreper vi bruker i grammatikken):
    - WHILE-setning, TILORDNING, KLASSEDEKL, UTTRYKK...
    - Startsymbolet : Lovlige setninger er utledet fra dette.
  - Meta-symboler**- slike hjelpesymboler/tegn vi bruker for å sette opp reglene:
- En grammatikk spesifiserer via regler lovlige sammensetninger av terminal- og ikke-terminal-symboler:  

```
<exp> ::= <exp> <op> <exp> | (<exp>) | NUMBER  
<op> ::= + | - | *
```

## Grammatikk - eksempel

- Grammatikken:  

```
<exp> ::= <exp> <op> <exp> | (<exp>) | NUMBER  
<op> ::= + | - | *
```
- Metasymboler: ::=, <, >, |
- Ikke-terminaler: exp, op
- Terminaler : NUMBER, (,), \*, +, -

::= 'leses som:' kan bestå av  
| 'betyr' eller

## Rollen til en grammatikk, semantikk og syntaks

- En grammatikk definerer *visse sider* av et språk – **syntaksen** (skrivereglene)
  - Gjelder naturlige språk (tysk, ..., norsk)
  - Gjelder alle kunstige språk som programmeringsspråk
  - Ikke alle de setningene vi kan lage i grammatikken, gir mening
- Semantikken** (meningen) med de ulike delene spesifiseres separat – se kap.6.
- En grammatikk definerer regler for sammensetting av de ulike elementene (terminalsymbolene) vi har i et språk.
  - Eks: En regel for en viss type setning på norsk er at den består av (spesifisert på BNF):  
<Setning> → <Subjekt><Predikat> | <Setning><Prep><Objekt>  
<Subjekt> → Per | Huset  
<Predikat> → løper | ler | står | slår | bygger  
<Objekt> → Kari | Huset | Per
  - En grammatikk kan lage (utlede) mange, ofte uendelig mange setninger
    - Noen følger grammatikken, men følger **ikke** semantikken og er da ikke setninger i språket
      - Eks: `int i = true;`
    - Noen følger både syntaksen og semantikken og er da korrekte setninger i språket, og vil da også kompilere
      - Men det er ikke sikkert at de løser problemet vårt (er riktige sett ut fra hva programmet skal gjøre)

## Hvorfor ikke bare én (stor) grammatikk?

- Kunne vi ikke laget en (stor) grammatikk som sa 'alt' om språket
  - F. eks det som tas av skanneren –hvordan tall og variable er definert
  - Som sa at det skal være samme type på hver side av en tilordning
  - Som ordnet alle presedensregler mellom operatører ( multiplikasjon før addisjon,... osv)
- Grammatikker spesifiserer ikke 'alt' ved språket, fordi:
  - En slik grammatikk vil bli 'uhåndterlig stor'
  - Vanskelig/"nær umulig" å formulere visse aspekter ved et språk
  - mye raskere å ta:
    - enkle ting i skanneren,
    - setningsformen i parseren
    - mer kompliserte krav i semantikk - sjekkeren
      - jfr. samlebåndsproduksjon av biler (bilen lages i flere steg)
- Må ofte jobbe med hvordan vi formulerer en grammatikk for at den skal gi en god/riktig parser
  - flere måter å formulere grammatikken for et språk

## Kontekstfrie grammatikker, BNF notasjon med variasjoner

- Bokas notasjon
 
$$\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{number}$$

$$\text{op} \rightarrow + \mid - \mid *$$
- En tradisjonell måte (Algol 60 rapporten)
 
$$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle \langle \text{op} \rangle \langle \text{exp} \rangle \mid ( \langle \text{exp} \rangle ) \mid \text{NUMBER}$$

$$\langle \text{op} \rangle ::= + \mid - \mid *$$
- Litt utvidet BNF
 
$$\text{exp} \rightarrow \text{exp} ( "+" \mid "-" \mid "*" ) \text{exp} \mid ( "(" \text{exp} ")" ) \mid \text{number}$$

## Flere måter å skrive den samme grammatikken

- Regnes som den mest basale

$$\begin{aligned} \text{exp} &\rightarrow \text{exp op exp} \\ \text{exp} &\rightarrow ( \text{exp} ) \\ \text{exp} &\rightarrow \text{number} \\ \text{op} &\rightarrow + \\ \text{op} &\rightarrow - \\ \text{op} &\rightarrow * \end{aligned}$$

- Kortest mulig

$$\begin{aligned} E &\rightarrow E O E \mid ( E ) \mid n \\ O &\rightarrow + \mid - \mid * \end{aligned}$$

## Avledning (venstrevledning)

av: **(number - number) \* number**

$$\begin{aligned} \text{exp} &\rightarrow \text{exp op exp} \\ \text{exp} &\rightarrow ( \text{exp} ) \\ \text{exp} &\rightarrow \text{number} \\ \text{op} &\rightarrow + \\ \text{op} &\rightarrow - \\ \text{op} &\rightarrow * \end{aligned}$$

Startsymbol

Produksjon brukt

- |   |  |
|---|--|
| (1) $\text{exp} \Rightarrow \text{exp op exp}$                    | $[\text{exp} \rightarrow \text{exp op exp}]$ |
| (2) $\Rightarrow (\text{exp}) \text{ op exp}$                     | $[\text{exp} \rightarrow ( \text{exp} )]$    |
| (3) $\Rightarrow (\text{exp op exp}) \text{ op exp}$              | $[\text{exp} \rightarrow \text{exp op exp}]$ |
| (4) $\Rightarrow (\text{number op exp}) \text{ op exp}$           | $[\text{exp} \rightarrow \text{number}]$     |
| (5) $\Rightarrow (\text{number} - \text{exp}) \text{ op exp}$     | $[\text{op} \rightarrow -]$                  |
| (6) $\Rightarrow (\text{number} - \text{number}) \text{ op exp}$  | $[\text{exp} \rightarrow \text{number}]$     |
| (7) $\Rightarrow (\text{number} - \text{number}) * \text{exp}$    | $[\text{op} \rightarrow *]$                  |
| (8) $\Rightarrow (\text{number} - \text{number}) * \text{number}$ | $[\text{exp} \rightarrow \text{number}]$     |

Vestrevledning = avleder med terminaler fra venstre  
Ferdig når det bare er terminal-symboler

$$L(G) = \{ s \mid \text{exp} \Rightarrow^* s \}$$

↑ grammatikk      ↑ streng med bare terminal-symboler

# Avledning (høyreavledning)

av: (number - number) \* number

```
exp → exp op exp
exp → ( exp )
exp → number
op → +
op → -
op → *
```

## Startsymbol

- (1)  $exp \Rightarrow exp\ op\ exp$
- (2)  $\Rightarrow exp\ op\ \mathbf{number}$
- (3)  $\Rightarrow exp\ * \mathbf{number}$
- (4)  $\Rightarrow ( exp ) * \mathbf{number}$
- (5)  $\Rightarrow ( exp\ op\ exp ) * \mathbf{number}$
- (6)  $\Rightarrow ( exp\ op\ \mathbf{number} ) * \mathbf{number}$
- (7)  $\Rightarrow ( exp - \mathbf{number} ) * \mathbf{number}$
- (8)  $\Rightarrow (\mathbf{number} - \mathbf{number}) * \mathbf{number}$

## Produksjon brukt

- [ $exp \rightarrow exp\ op\ exp$ ]
- [ $exp \rightarrow \mathbf{number}$ ]
- [ $op \rightarrow *$ ]
- [ $exp \rightarrow ( exp )$ ]
- [ $exp \rightarrow exp\ op\ exp$ ]
- [ $exp \rightarrow \mathbf{number}$ ]
- [ $op \rightarrow -$ ]
- [ $exp \rightarrow \mathbf{number}$ ]

Høyreavledning = avleder med terminaler fra høyre  
Ferdig når det bare er terminalsymboler

Det finnes også mange andre rekkefølger å avlede en setning fra en grammatikk

# Opplagte krav til en grammatikk

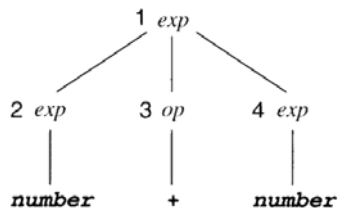
Alle ikke-terminaler:

- Må kunne inngå i en streng avledet fra startsymbolet
- Må kunne avledes videre til noe som bare inneholder terminal-symboler
- Eks:
  - A → Bx
  - B → Ay
  - C → z
- Kan aldri avlede fra A til bare terminalsymboler
- C kan ikke inngå i noen streng avledet fra A
- Altså en håpløs grammatikk

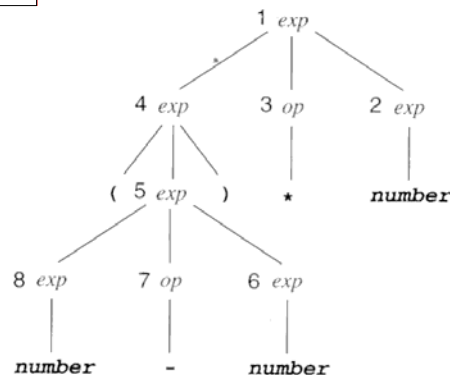
# Parserings-tre (konkret syntaks-tre)

- (1)  $exp \Rightarrow exp\ op\ exp$
- (2)  $\Rightarrow \mathbf{number}\ op\ exp$
- (3)  $\Rightarrow \mathbf{number} + exp$
- (4)  $\Rightarrow \mathbf{number} + \mathbf{number}$

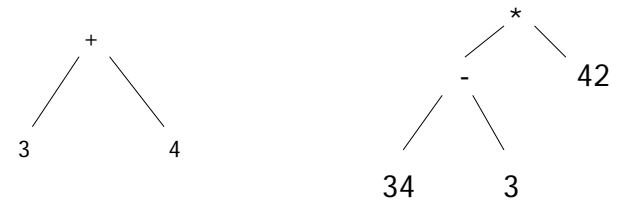
- Representasjon som er uavhengig av avlednings-rekkefølgen



- Tallene angir høyre/avledning
- Ser vi bort fra tallene, gir alle avlednings/rekkefølger det samme treet



# (Abstrakt) syntakstre – forlengs polsk notasjon – "fra tre til fil-format" – det vi egentlig trenger videre



Prefix-form av treet for (34-3)\*42: ("veltet til venstre" – prefiks traversering – først noden, så venstre sub-tre, så høyre sub-tre):  
OpExp(Times, OpExp(Minus, Const(34), Const(3)), Const(42))

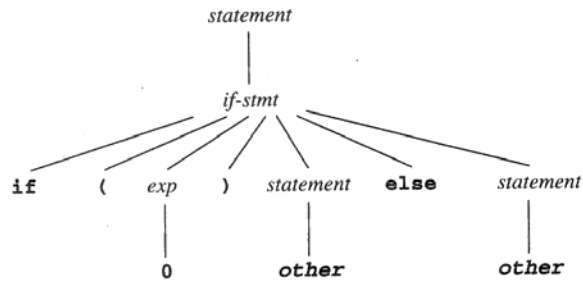
## Flere parserings-trær

G1:

$$\begin{aligned} \text{statement} &\rightarrow \text{if-stmt} \mid \text{other} \\ \text{if-stmt} &\rightarrow \text{if ( exp ) statement} \\ &\quad \mid \text{if ( exp ) statement else statement} \\ \text{exp} &\rightarrow 0 \mid 1 \end{aligned}$$

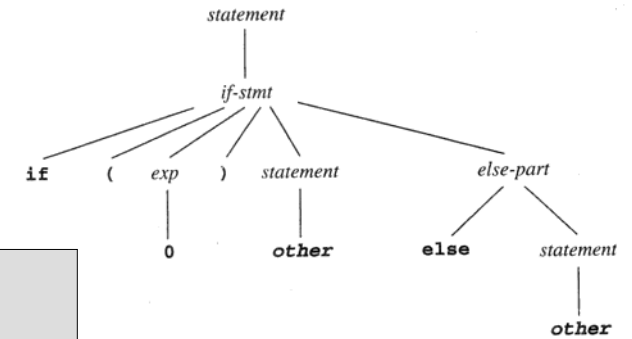
Setning:

**if (0) other else other**

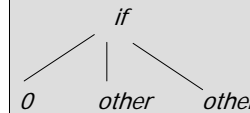


## En annen grammatikk G2 for if-setninger

G2:

$$\begin{aligned} \text{statement} &\rightarrow \text{if-stmt} \mid \text{other} \\ \text{if-stmt} &\rightarrow \text{if ( exp ) statement else-part} \\ \text{else-part} &\rightarrow \text{else statement} \mid \epsilon \\ \text{exp} &\rightarrow 0 \mid 1 \end{aligned}$$


Felles abstrakt syntaks-tre for G1 og G2:



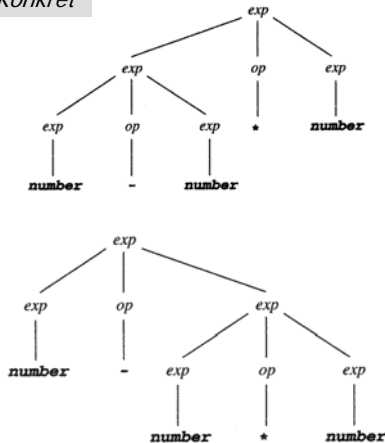
## Tvetydige grammatikker - analyse av setningen: $n - n * n$

$$\begin{aligned} \text{exp} &\rightarrow \text{exp op exp} \mid ( \text{exp} ) \mid \text{number} \\ \text{op} &\rightarrow + \mid - \mid * \end{aligned}$$

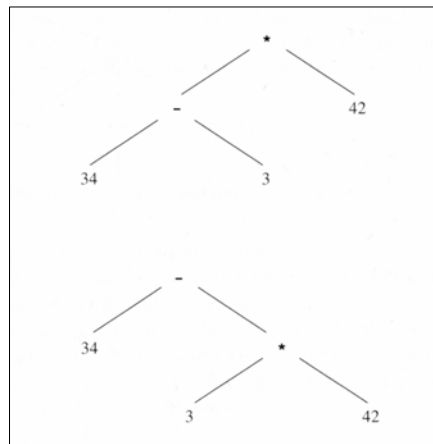
G er flertydig hvis det finnes en setning i  $L(G)$  som kan gis flere parserings-trær

Dette eksempelet er essensiell tvetydighet, angir ulike beregninger.

Konkret



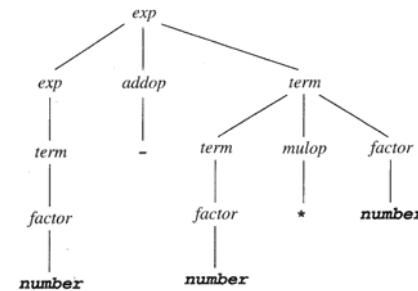
Abstrakt



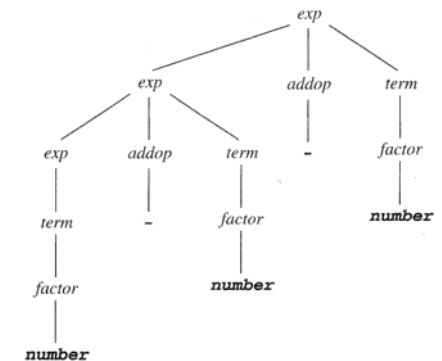
## Entydige grammatikker for samme språket, to trær for to ulike uttrykk.

$$\begin{aligned} \text{exp} &\rightarrow \text{exp addop term} \mid \text{term} \\ \text{addop} &\rightarrow + \mid - \\ \text{term} &\rightarrow \text{term mulop factor} \mid \text{factor} \\ \text{mulop} &\rightarrow * \\ \text{factor} &\rightarrow ( \text{exp} ) \mid \text{number} \end{aligned}$$

**34-3\*42**



**34-3-42**



## To eksempler på viktigheten av assosiativitet (og presedens)

a) Melding til lederen for en eksekusjons-peletong:

- Stopp ikke skyt ham !
  - Stopp ikke - skyt ham
  - Stopp - ikke skyt ham

b) uttrykket : 3 - 4 - 5

- 3 - (4-5) = 3-(-1) = 4      høyre
- (3-4) -5 = -6                      venstre

## Presedens og assosiativitet uttrykt i grammatikken

- Presedens i grammatikken for operatører
  - Noen operasjoner gjøres før andre (\* før +)
  - **Ordnes med** kaskading av definisjoner av operatorgrupper (addop, multop, expop,..) – se forrige grammatikk
- Assosiativitet i grammatikken for operatører:
  - Venstre (assositivitet) : Operatører med samme presedens utføres fra venstre mot høyre
  - Høyre (assositivitet) : Operatører med samme presedens utføres fra høyre mot venstre.
  - Ingen (assositivitet) : Tillater ingen rekkefølge av slike operasjoner. i samme .
  - **Ordnes med** at uttrykk med slike operatører bare kan utvides på den ene siden :
    - $exp \rightarrow exp \text{ addop } term \mid term$
    - ( denne utvider på venstre side og blir venstre assisiativ)
    - Ingen assosiativitet ordnes med at vi krever parenteser med bare en operator i en parentes.

## Vanlig å:

- Angi språket ved flertydige grammatikk som

$$exp \rightarrow exp \text{ op } exp \mid ( exp ) \mid \mathbf{number}$$

$$op \rightarrow + \mid - \mid *$$

- Oppgi regler for presedens og assosiativitet for hver operasjon, slik at alle setninger får ett entydig syntakstre:
  - +, lav, venstre ass.
  - \*, høy, venstre-ass
  - ↑, høyerst, høyre ass.
- Dette er helt greit for binære infix-operatører, men fungerer "vanligvis" også greit for unære postfix eller prefiks operatører

- 3+5 \* 3 \* 2 + 4 ↑2 \* 3

## Presedens og assosiativitet i Java

### Operator Precedence

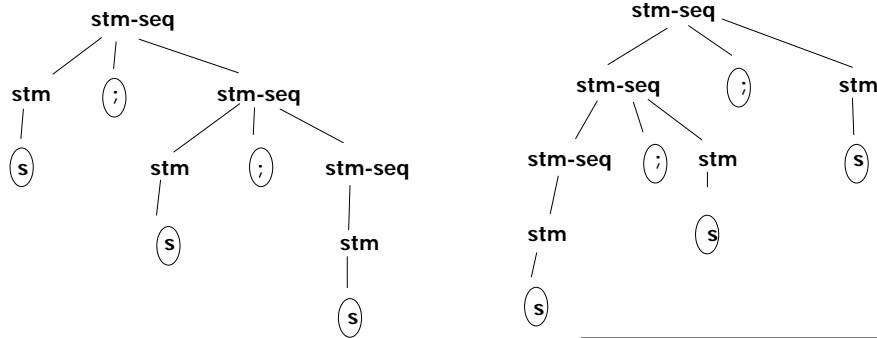
Venstre assosiativ

Java performs operations assuming the following ordering (or *precedence*) rules if parentheses are not used to determine the order of evaluation (operators on the same line are evaluated in left-to-right order subject to the conditional evaluation rule for `&&` and `||`). The operations are listed below from highest to lowest precedence (we use `(exp)` to denote an atomic or parenthesized expression):

postfix ops	<code>[] . ((exp)) (exp) ++ (exp) --</code>
prefix ops	<code>++(exp) --(exp) -(exp) ~(exp) !(exp)</code>
creation/cast	<code>new ((type))(exp)</code>
mult./div.	<code>* / %</code>
add./subt.	<code>+ -</code>
shift	<code>&lt;&lt; &gt;&gt; &gt;&gt;&gt;</code>
comparison	<code>&lt; &lt;= &gt; &gt;= instanceof</code>
equality	<code>== !=</code>
bitwise-and	<code>&amp;</code>
bitwise-xor	<code>^</code>
bitwise-or	<code> </code>
and	<code>&amp;&amp;</code>
or	<code>  </code>
conditional	<code>(bool.exp)? {true.val}: {false.val}</code>
assignment	<code>=</code>
op assignment	<code>+= -= *= /= %=</code>
bitwise assign.	<code>&gt;&gt;= &lt;&lt;= &gt;&gt;&gt;=</code>
boolean assign.	<code>&amp;= ^=  =</code>

Ikke-essensiell flertydighet

stm-seq  $\rightarrow$  stm-seq; stm | stm *venstre-ass*  
 stm-seq  $\rightarrow$  stm; stm-seq | stm *høyre-ass*  
 stm  $\rightarrow$  s



Kan like gjerne representeres som:

seq eller: seq

## "Dangelig else" - problemet

- Problem: Hvilken **if**-setning skal vi koble **else** til?

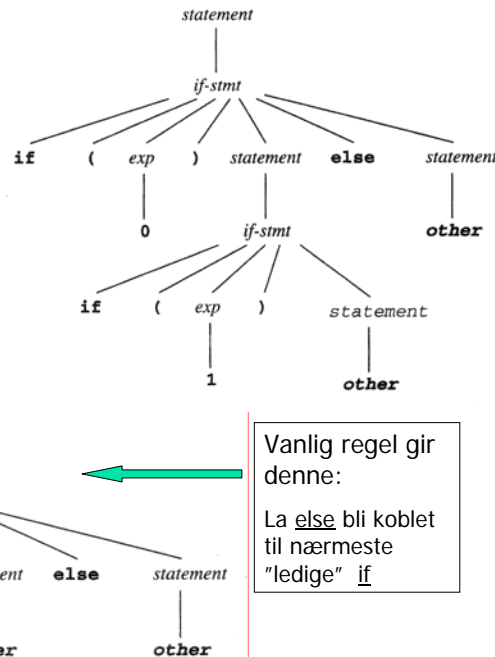
**if (0) if (1) other else other**

- Grammatikken (tvetydig):

statement  $\rightarrow$  if-stmt | **other**  
 if-stmt  $\rightarrow$  **if** ( exp ) statement  
 | **if** ( exp ) statement **else** statement  
 exp  $\rightarrow$  0 | 1

## To muligheter:

statement  $\rightarrow$  if-stmt | **other**  
 if-stmt  $\rightarrow$  **if** ( exp ) statement  
 | **if** ( exp ) statement **else** statement  
 exp  $\rightarrow$  0 | 1

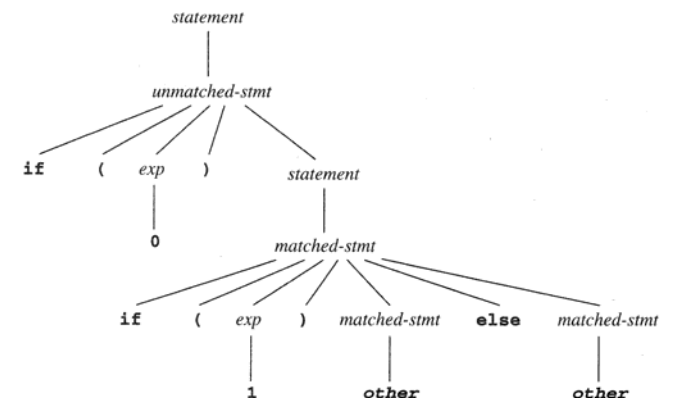


Vanlig regel gir denne:  
 La else bli koblet til nærmeste "ledige" if

## Eks: Entydig grammatikk for if-setning. Gir "vanlig" løsning

**if (0) if (1) other else other**

statement  $\rightarrow$  matched-stmt | unmatched-stmt  
 matched-stmt  $\rightarrow$  **if** ( exp ) matched-stmt **else** matched-stmt | **other**  
 unmatched-stmt  $\rightarrow$  **if** ( exp ) statement  
 | **if** ( exp ) matched-stmt **else** unmatched-stmt  
 exp  $\rightarrow$  0 | 1



Idé:  
**matched-stmt**  
 - kan *ikke* kobles med etterfølgende else  
**unmatched-stmt**  
 - kan kobles med etterfølgende else  
 Er det opplagt at denne kan generere alle "lovlige" setninger (ut fra den kortere flertydige grammatikken på forrige to foiler)?

# Utvidet BNF (EBNF)

Idé: Man kan generelt bruke "regulære uttrykk" på høyresiden i produksjoner

Vanlig:  $\alpha^*$  skrives:  $\{\alpha\}$   $\alpha$  er en streng av terminaler og ikke-terminaler  
 $\alpha?$  skrives:  $[\alpha]$

Eksempel:

$exp \rightarrow exp \left( \begin{matrix} "+" \\ "-" \\ "*" \end{matrix} \right) exp \mid \left( \begin{matrix} "(" \\ "exp" \end{matrix} \right) \mid number$

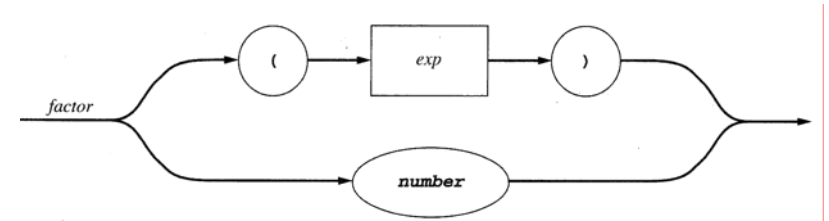
Meta-symbol                      ikke-meta

$A \rightarrow A \alpha \mid \beta$  kan skrives:  $A \rightarrow \beta \{\alpha\}$   
 $A \rightarrow \alpha A \mid \beta$  kan skrives:  $A \rightarrow \{\alpha\} \beta$   
 $stm\text{-seq} \rightarrow stm \{ ; stm \}$  eller  $\rightarrow \{stm; \} stm$   
 $if\text{-setn} \rightarrow if (expr) stm [ else stmt ]$

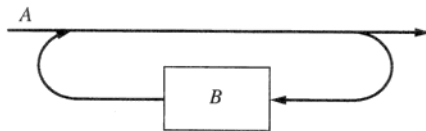
Merk: For en del metoder etc. må man gå ut fra basal BNF.

# Syntaks-diagrammer

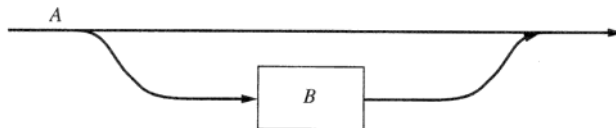
$factor \rightarrow ( exp ) \mid number$



$A \rightarrow \{ B \}$



$A \rightarrow [ B ]$

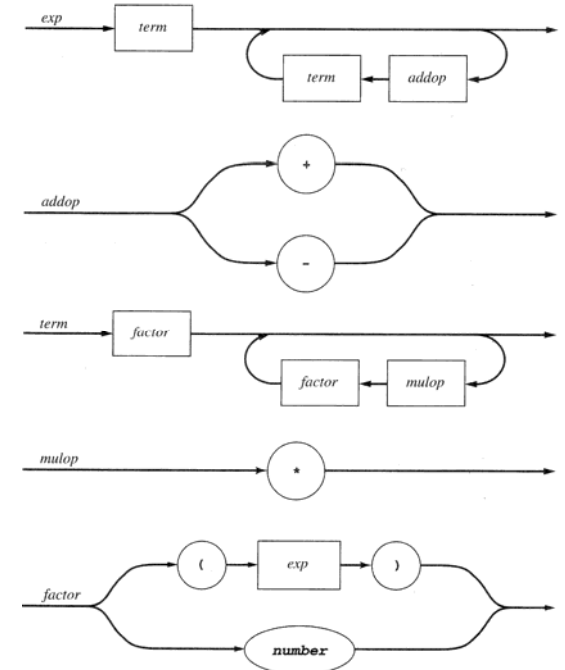


BNF-grammatikk

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

EBNF-grammatikk for samme "språket", brukes til å lage syntaksdiagrammene

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





# Chomsky-hierarkiet

*a* er vilkårlig terminalsymb.  
*β, α, γ* er vilkårlig samling av terminal- og ikke-terminalsymb.  
*A, B* er ikke-terminaler

- Type 0 – språk
  - Urestrikkerte prod.:  $\alpha \rightarrow \beta, \quad \alpha \neq \epsilon$  ( $\alpha$  er ikke-tom)
- Type 1 – språk
  - Kontekst-sensitive produksjoner:  $\beta A \gamma \rightarrow \beta \alpha \gamma$
- Type 2 – språk
  - Kontkstfrie prod.:  $A \rightarrow \alpha$
- Type 3 språk
  - Regulære språk:
    - Regulære uttrykk
    - NFA
    - DFA
  - Produksjoner bare på formen:
    - $A \rightarrow Ba$  og  $A \rightarrow a$
    - eller
    - $A \rightarrow aB$  og  $A \rightarrow a$

# BNF-grammatikk for TINY

```

program → stmt-sequence
stmt-sequence → stmt-sequence ; statement | statement
statement → if-stmt | repeat-stmt | assign-stmt | read-stmt | write-stmt
if-stmt → if exp then stmt-sequence end
           | if exp then stmt-sequence else stmt-sequence end
repeat-stmt → repeat stmt-sequence until exp
assign-stmt → identifier := exp
read-stmt → read identifier
write-stmt → write exp
exp → simple-exp comparison-op simple-exp | simple-exp
comparison-op → < | =
simple-exp → simple-exp addop term | term
addop → + | -
term → term mulop factor | factor
mulop → * | /
factor → ( exp ) | number | identifier
    
```

# Nodestruktur i C for TINY

```

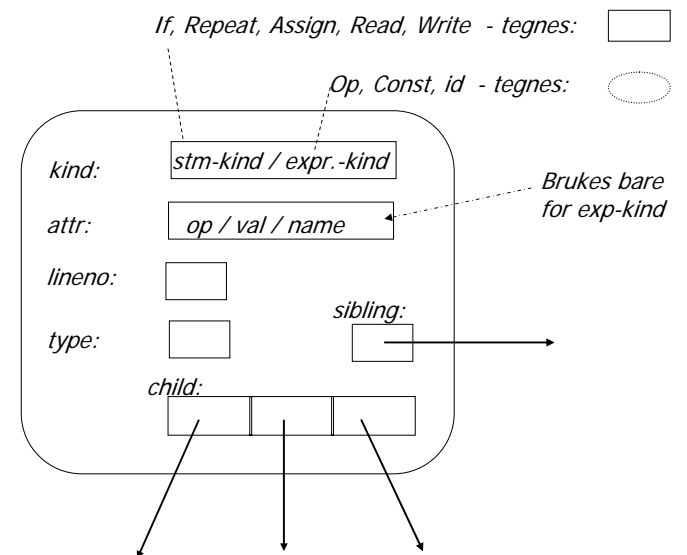
typedef enum {StmtK,ExpK} NodeKind;
typedef enum {IfK,RepeatK,AssignK,ReadK,WriteK}
           StmtKind;
typedef enum {OpK,ConstK,IdK} ExpKind;

/* ExpType is used for type checking */
typedef enum {Void,Integer,Boolean} ExpType;

#define MAXCHILDREN 3

typedef struct treeNode
{ struct treeNode * child[MAXCHILDREN];
  struct treeNode * sibling;
  int lineno;
  NodeKind nodekind;
  union { StmtKind stmt; ExpKind exp; } kind;
  union { TokenType op;
         int val;
         char * name; } attr;
  ExpType type; /* for type checking of exps */
} TreeNode;
    
```

# Nodestruktur i C for TINY



Denne nodestrukturen passer enda bedre med et OO-språk med klasser /subklasser som implementasjons-språk.

# Syntaks-tre for et program i Tiny

```
{ Sample program
  in TINY language-
  computes factorial
}
read x; { input an integer }
if 0 < x then { don't compute if x <= 0 }
  fact := 1;
  repeat
    fact := fact * x;
    x := x - 1
  until x = 0;
  write fact { output factorial of x }
end
```

```
read x; { input an integer }
if 0 < x then { don't compute if x <= 0 }
  fact := 1;
  repeat
    fact := fact * x;
    x := x - 1
  until x = 0;
  write fact { output factorial of x }
end
```

