

INF5110, onsdag 8. februar 2012
Mer fra kap 4:
Parsering ovenfra-ned ("top down")
Stein Krogdahl Ifi, UiO



Merk: Ikke forelesning på onsdag 15. februar

Oppgaver som gjennomgås i dag:

(Svar ligger bakerst. Ikke titt!)

- Spørsmålene på de to siste foilene fra onsdag 1/2 (Bl.a. lag et subklasse-hierarki som ville passe for nodene i et abstrakt syntakstre for Tiny)
- Formuler en **entydig** uttrykks-grammatikk som også har (høyreassosiativ) eksponensiering
- Beskriv en algoritme som finner alle de ikke-terminaler som kan produsere den tomme streng (uten også å finne First-mengden). Vi kaller dem altså de "utnullbare" ikke-terminaler.



Annonse:

Søk jobb på "Package Template"-kompilatoren!

- I et prosjekt (SWAT-prosjektet) har vi designet en språkmekanisme som kalles "Package Templates", forkortet PT
- Vi har beskrevet hvordan denne mekanismen kan legges inn som en utvidelse av Java, til et språk "JPT"
- Vi er nå kommet langt med en kompilator for JPT, og mye er allerede på plass.
- Vi hadde to til å jobbe på kompilatoren, men nå har vi bare en (Steinar Kaldager)
- Vi trenger en til, bl.a. til å jobbe med uttesting.
- Lønn: ca. 170 kr/time ?
- Søknadsfrist: fredag 10/2, til steinkr@ifi.uio.no
- Bare send en mail med relevante opplysninger (også karakterer i relevante fag)

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

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

- 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

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

NB: Ikke i boka, men er pensum!

- Den "utvidede startmengden", *Efirst*, til en produksjon det som kan ligge som **førstkommende token** i input dersom $A \rightarrow \alpha$ et riktig valg på dette stadiet under parsingen.
- Her må man også tenke på tilfellet at α kan være utnullbar, og da kan også Follow(**A**) komme som **førstkommende token**
- Mengden kan beregnes slik:
 $Efirst(A \rightarrow \alpha) = First(\alpha)$, pluss, om α er utnullbar, 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 er utnullbare

På forrige side vil da både

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

og

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

*inneholde **else***

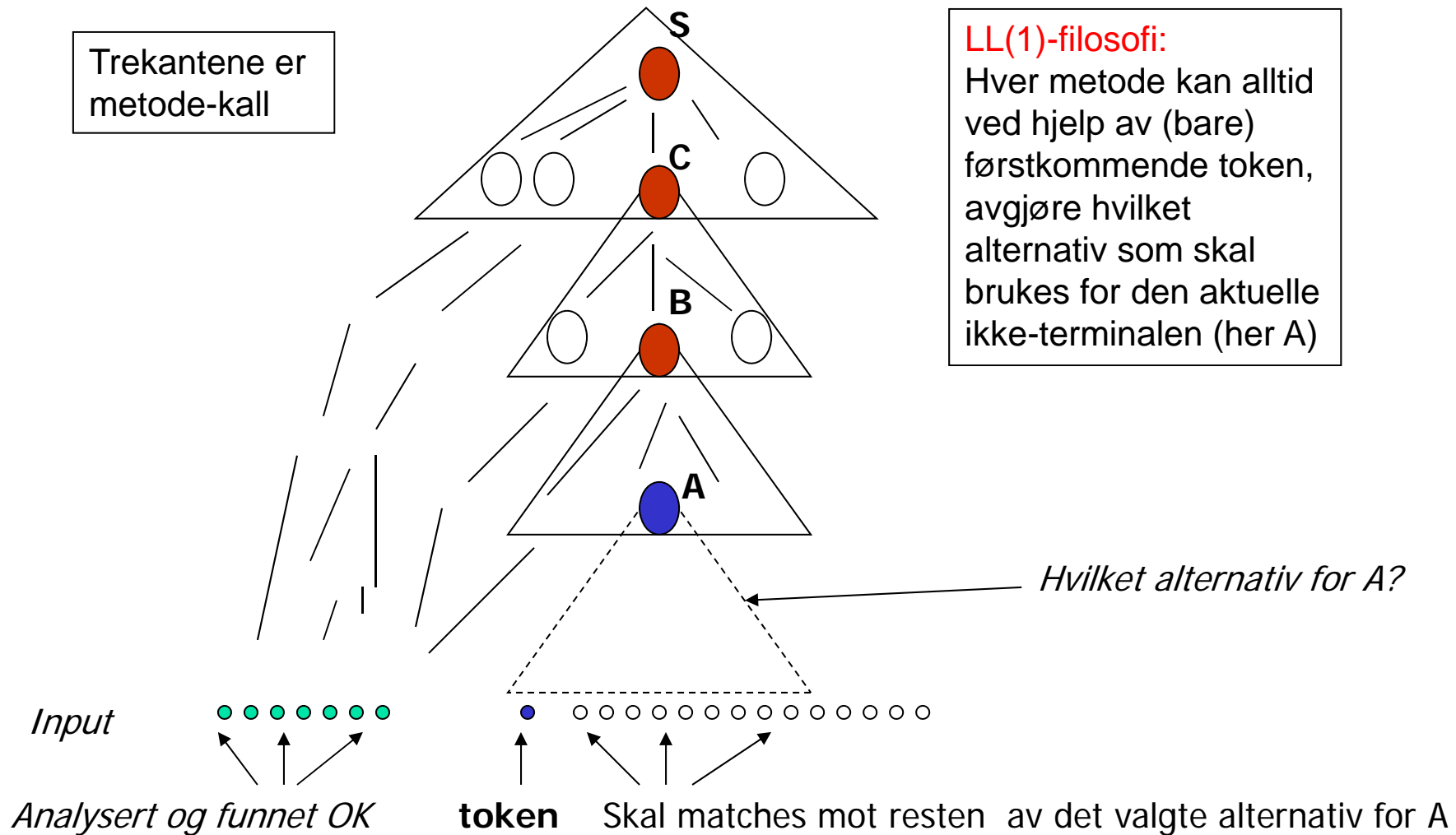
*Dermed er grammatikken *ikke**

LL(1)

Tidligere foil:

Situasjonen under rekursiv parsing

Trekantene er metode-kall



LL(1) –tabell for uttrykks-grammatikk

Har fjernet venstre-
rekursjon:

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

Vi får da følgende First- og
Follow-mengder:

$$\text{First}(\text{exp}) = \{ (, \mathbf{number} \}$$
$$\text{First}(\text{exp}') = \{ +, -, \varepsilon \}$$
$$\text{First}(\text{addop}) = \{ +, - \}$$
$$\text{First}(\text{term}) = \{ (, \mathbf{number} \}$$
$$\text{First}(\text{term}') = \{ *, \varepsilon \}$$
$$\text{First}(\text{mulop}) = \{ * \}$$
$$\text{First}(\text{factor}) = \{ (, \mathbf{number} \}$$
$$\text{Follow}(\text{exp}) = \{ \$,) \}$$
$$\text{Follow}(\text{exp}') = \{ \$,) \}$$
$$\text{Follow}(\text{addop}) = \{ (, \mathbf{number} \}$$
$$\text{Follow}(\text{term}) = \{ \$,), +, - \}$$
$$\text{Follow}(\text{term}') = \{ \$,), +, - \}$$
$$\text{Follow}(\text{mulop}) = \{ (, \mathbf{number} \}$$
$$\text{Follow}(\text{factor}) = \{ \$,), +, -, * \}$$

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					



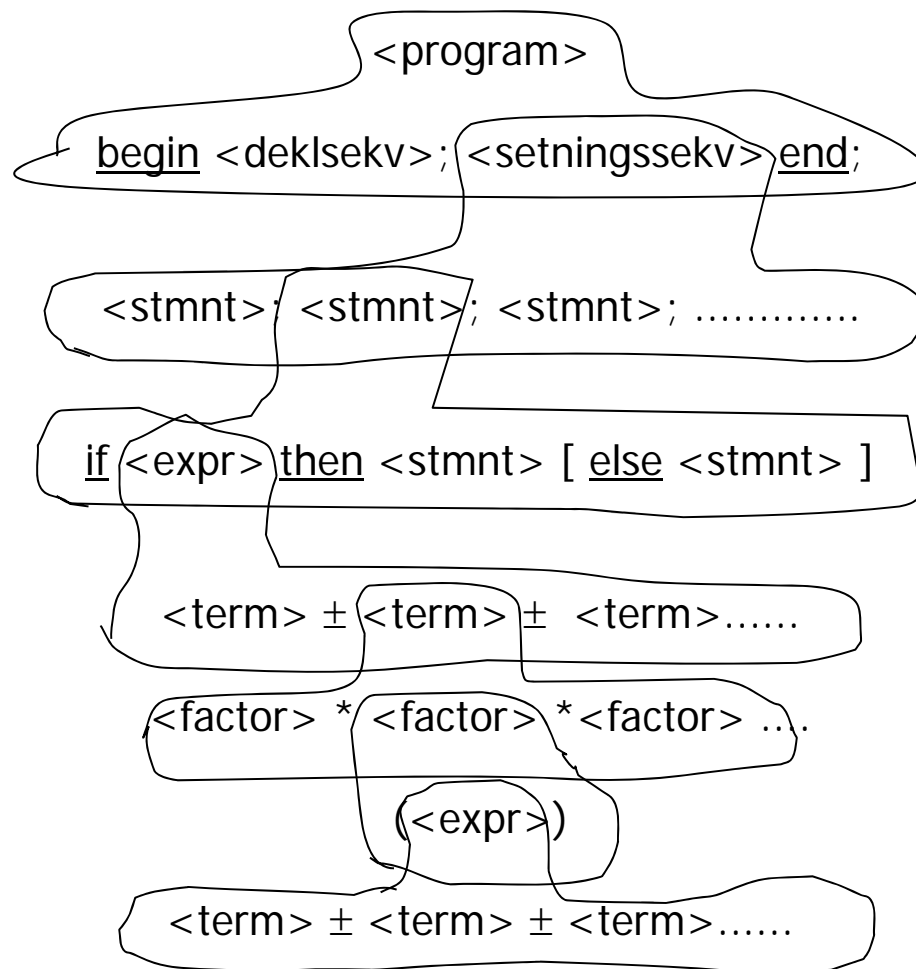
Når kompilatoren oppdager feil

- Minstekrav:
 - Tester løpende at programmet er OK, og gir fornuftig feilmelding ved feil (men stopper)
- Vanlig krav ved feil ("error recovery"):
 - Gir fornuftig feilmelding.
 - "Blar forbi feilen" og fortsetter kompileringen (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 for *syntaksfeil* det er vanskeligst å ta opp tråden etter feil.
- 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 rapportere feil *uten å lese noe fra input.*

Behandling av Syntaksfeil

ved "recursive decent" parsing.

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 "rec. descent"

Ut fra skissen 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, 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.

Uttrykksprosedyrer ved "error recovery"

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 { +, - } ?

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

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(...)

```
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;
```



Løsningsforslag på oppgaver

Noen spørsmål om Tiny-grammatikken

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**

- Er grammatikken entydig?
- Hva om vi vil tillate tomme setninger
- Hva om vi vill ha semikolon etter og ikke mellom setningene?
- Hva slags assosiativitet og presedens er det for operatorene?

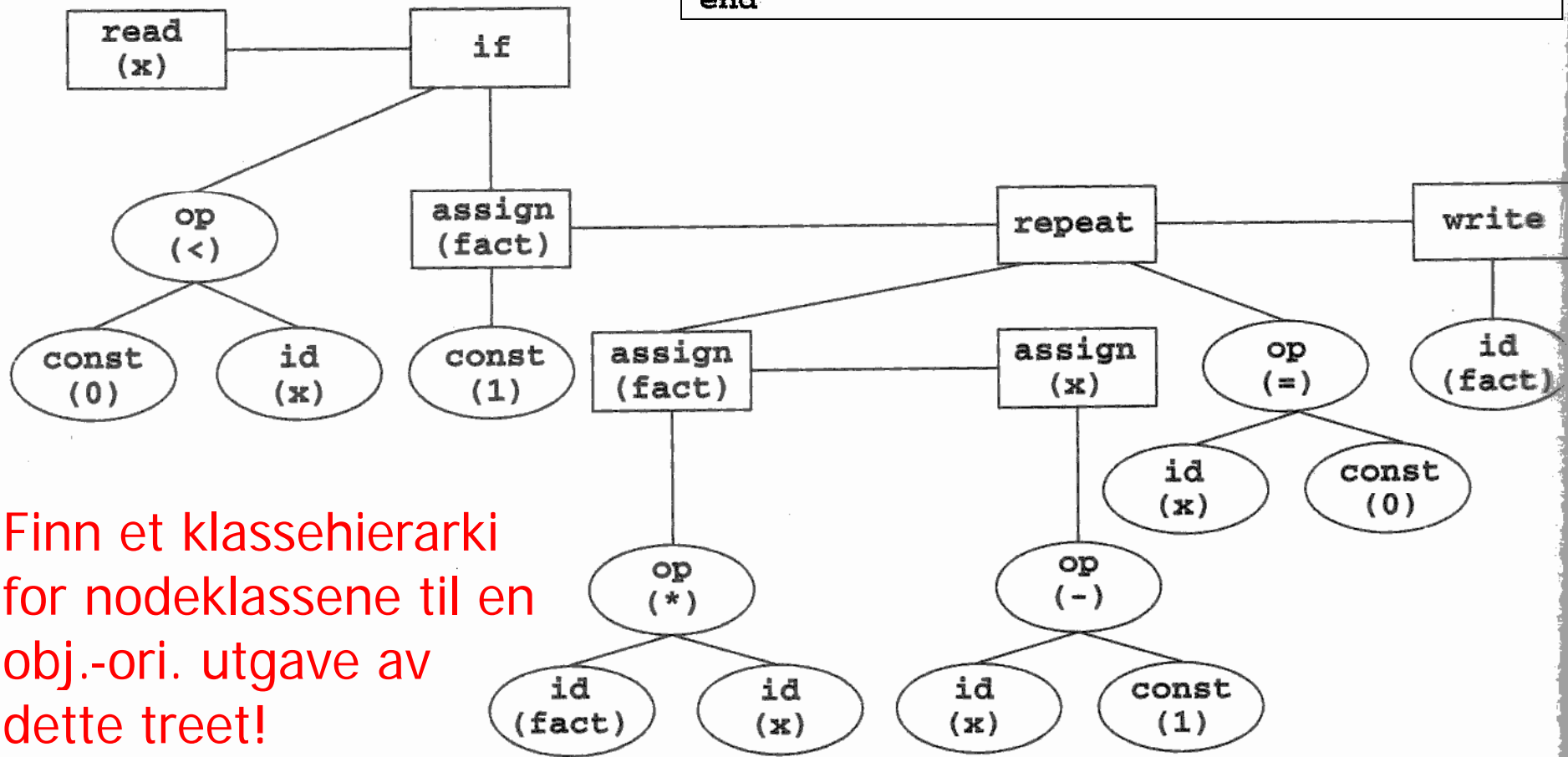


Svar på spørsmålene på forrige foil

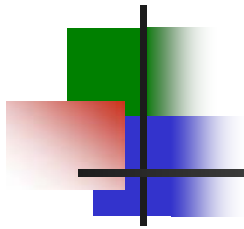
- Er grammatikken entydig?
 - Dette er generelt *uavgjørbart* for generelle BNF-grammatikker
 - Vi skal se på metode for å avgjøre det for mange praktiske grammatikker
 - Denne er ihvertfall delt opp i presedens-nivåer, og har assosiativites-angivelse, og er nok derved entydig. If-setninger med både **else** og **end** er ikke noe problem for entydigheten.
- Hva om vi vil tillate tomme setninger
 - Det er bare å sette til et tomt alternativ for *statement*
 - Den ser ut til fremdeles å være entydig
- Hva om vi vil ha semikolon etter og ikke mellom setningene?
 - Bytt ut reglen for *stmt-sequence* med:
stmt-sequence - \rightarrow *stmt-sequence statement ; | statement ;*
- Hva slags assosiativitet og presedens er det for operatorene?
 - Høyest presedens * / Venstre-assosiativ
 - Midlere presedens + - Venstre-assosiativ
 - Lavest presedens < = Ikke-assosiativitet (bare to operander)
 - Kunne altså brukt flertydig grammatikk, med disse tilleggs-reglene

```
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
```

Abstrakt syntakstre for Tiny-programmet:



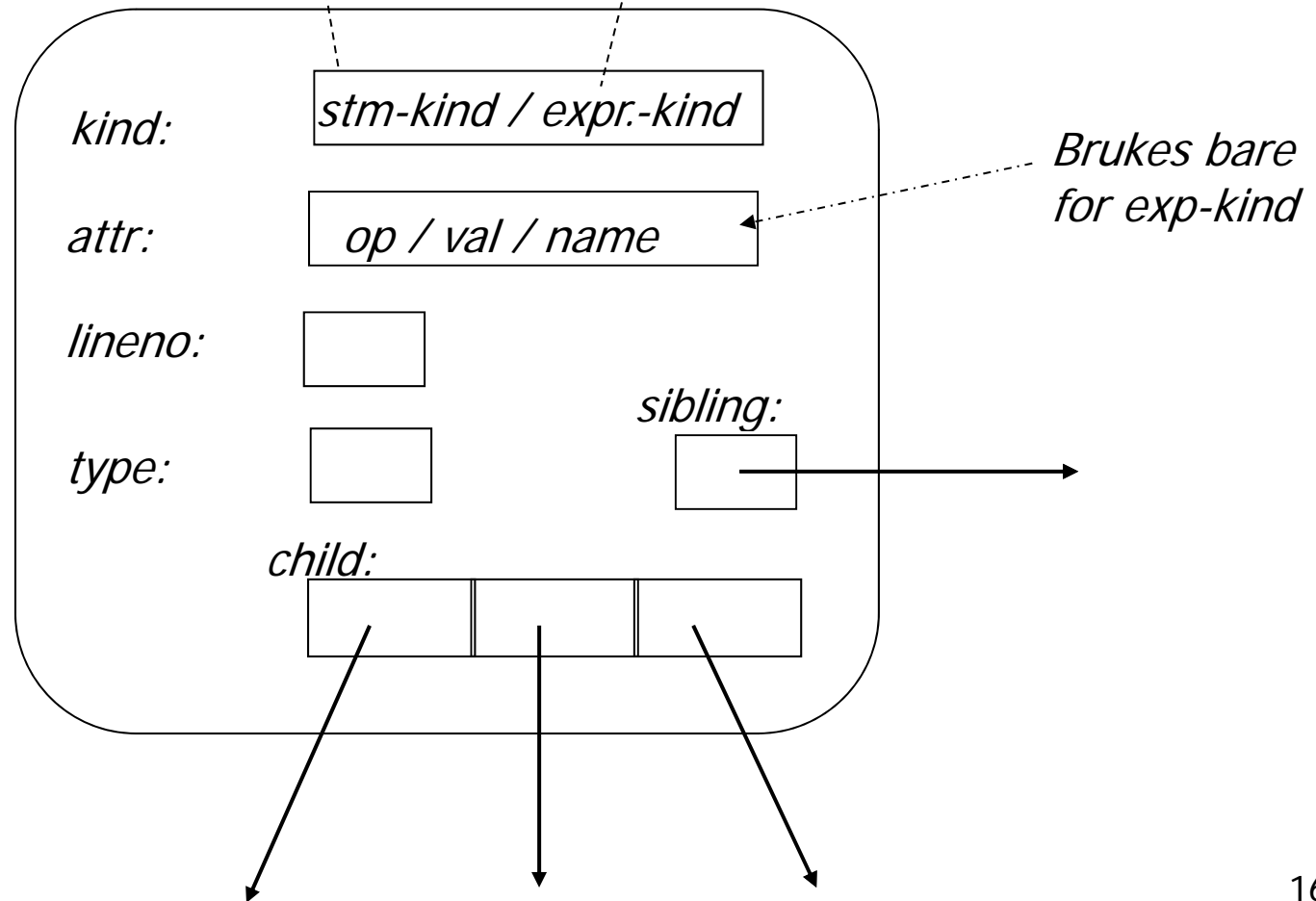
Finn et klassehierarki for nodeklassene til en obj.-ori. utgave av dette treet!



Nodestruktur i C for Tiny

If, Repeat, Assign, Read, Write - tegnes:

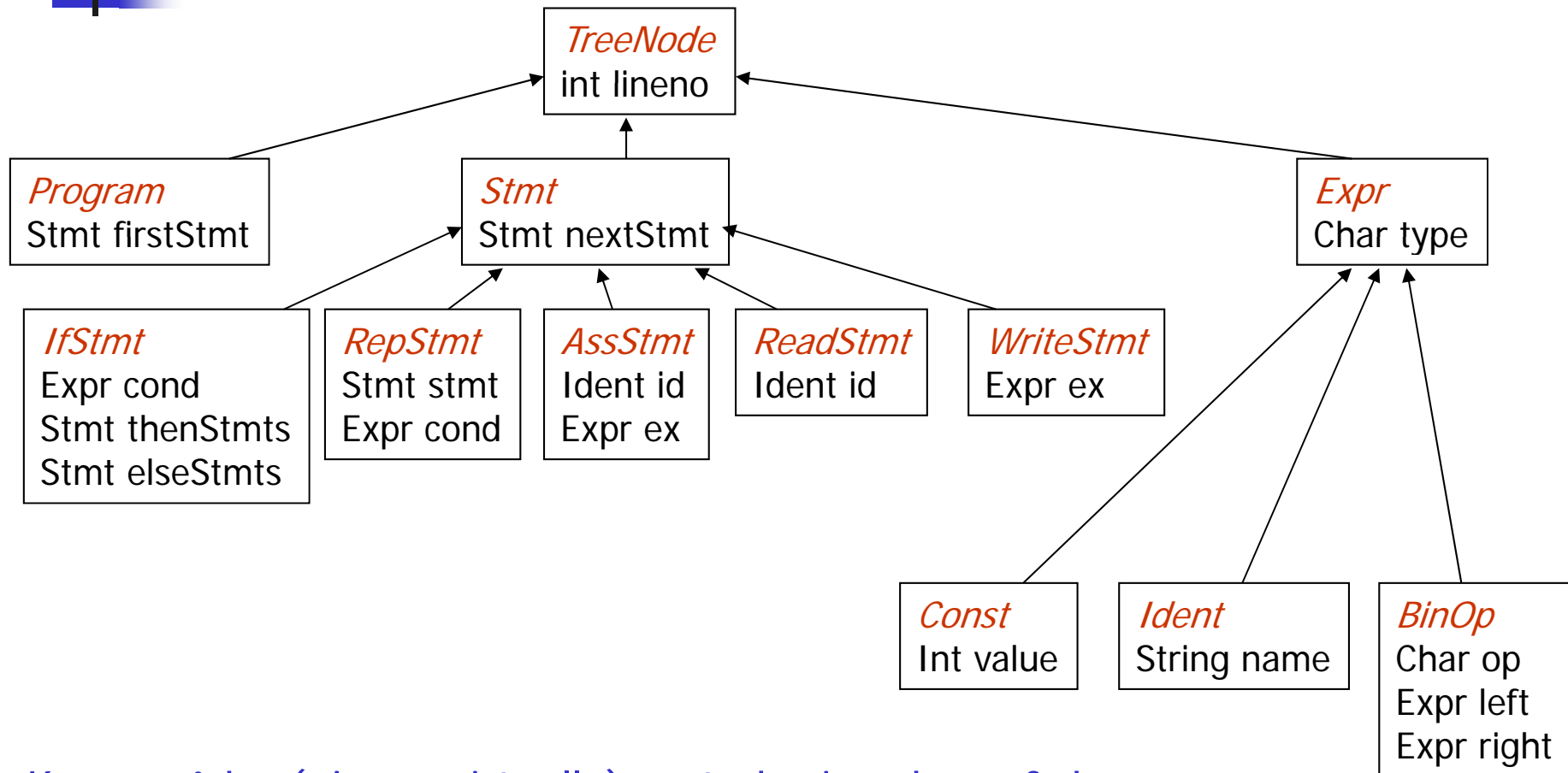
Op, Const, id - tegnes:



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

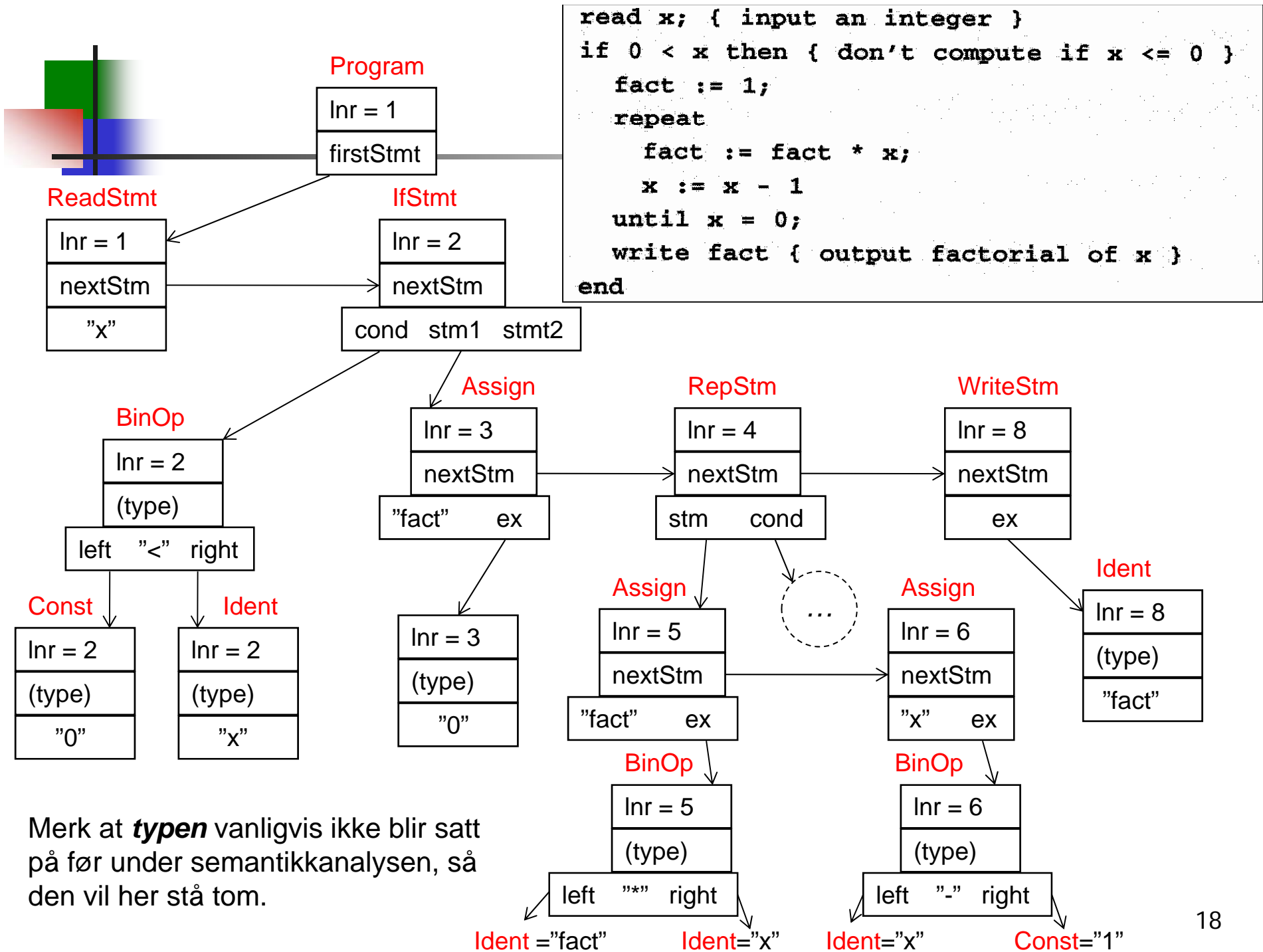
Nodeklasser for OO-utgave av abst.-synt.-tre for Tiny-språket

Merk: Dette er altså en fast subklasse-struktur i kompilatoren, og må ikke forveksles med det abstrakte syntaks-treet for et gitt program



Kan også ha (gjerne virtuelle) metoder i nodene, f.eks:

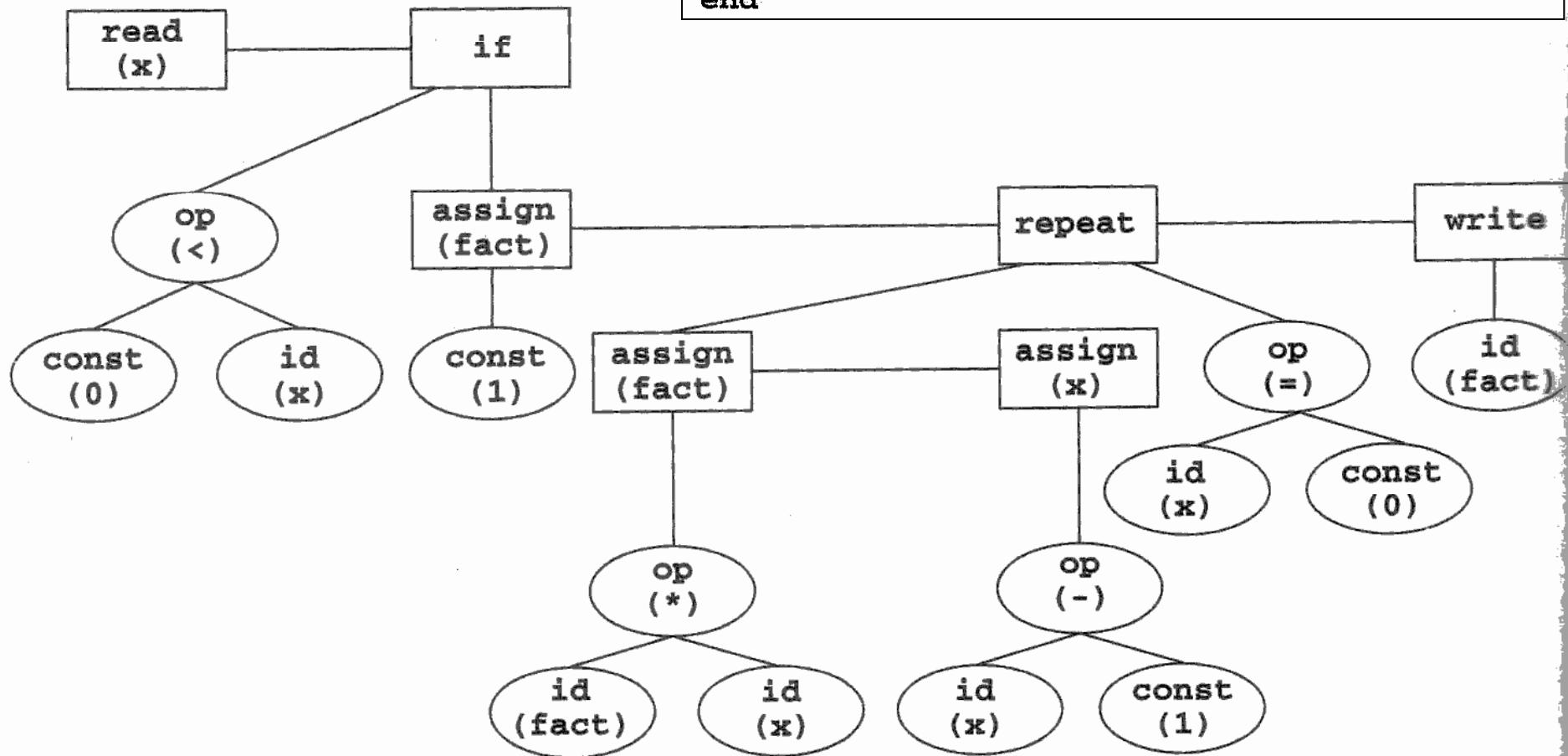
- `doSemAnalyses()`; Gjør semantisk analyse av noden, og av subtreet det er rot i
- `generateCode()`; Genererer kode for noden, og for subtreet det er rot i



Merk at **typen** vanligvis ikke blir satt på før under semantikkanalysen, så den vil her stå tom.

Abstrakt syntakstre for Tiny-programmet:

```
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
```





Entydig grammatikk for uttrykk med eksponensiering

Uten eksponensiering:

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

Med eksponensiering. Vi treneger en ny ikke-terminal "expon":

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

Høyre-assosiativ!



Algoritme for å beregne mengden av utnullbare ikke-terminaler (uten å beregne First)

Vi skal ha en mengde "UNB" av ikke-terminaler:

- Er fra starten tom
- Skal få nye elementer etter en regel 1 og 2 under
- Når den ikke øker mer er vi ferdig

1. Først legges alle ikke-terminaler A som har en produksjon $A \rightarrow \epsilon$, inn i UNB

2. (Gjentas til det ikke blir mer forandring):

Om en ikke-terminal A har et alternativ:

$$A \rightarrow \dots \mid B C \dots G \mid \dots$$

der alle $B C \dots G$ er ikke-terminaler som allerede er med i UNB, så legg også A inn i UNB.