

Kap.4 del 2
Top Down Parsing
INF5110 – v2006

Stein Krogdahl
Ifi, UiO

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:

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

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



Alternativ def. av LL(1) grammatikker

Sier at alle alternativer for A skal ha disjunkte startmengder (og dermed høyst ett utnullbart alternativ, ellers ville ε være med i startmengden for flere alternativer)

A grammar in BNF is **LL(1)** if the following conditions are satisfied.

1. For every production $A \rightarrow \alpha_1 \mid \alpha_2 \mid \dots \mid \alpha_n$, $\text{First}(\alpha_i) \cap \text{First}(\alpha_j)$ is empty for all i and j , $1 \leq i, j \leq n$, $i \neq j$.
2. For every nonterminal A such that $\text{First}(A)$ contains ε , $\text{First}(A) \cap \text{Follow}(A)$ is empty.

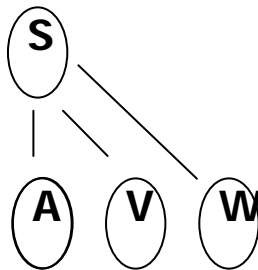
Dersom A er utnullbar, så må $\text{First}(A)$ og $\text{Follow}(A)$ være disjunkte.

At 1. og 2. til sammen er ekvivalent med det opprinnelige LL(1)-kravet krever et lite bevis.

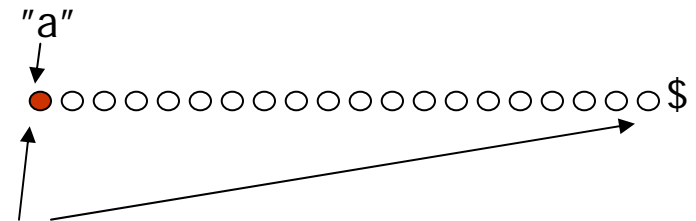
Kap 4.2: LL(1) –parsering med eksplisitt stakk

- Dette kalles i boka bare "LL(1)-parsering"
- Gjør logisk sett nøyaktig det samme som rek. desc. parsering for ren BNF.
- Bruker LL(1)-tabellen $M[T,N]$ til å styre parseringen

Start-situasjonen: (S er start-symbolet)



input:

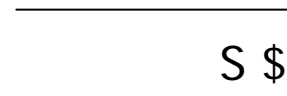


hele input, med "a" som første token

Anta at $a \in \text{First}(S)$ og at $M[S,a] = "S \rightarrow AVW"$, da bruker vi denne produksjonen ved å erstatte S på stakken med AVW.

Se det generelle steget på neste foil

Stakk etter init.:



Stakk etter ett steg:

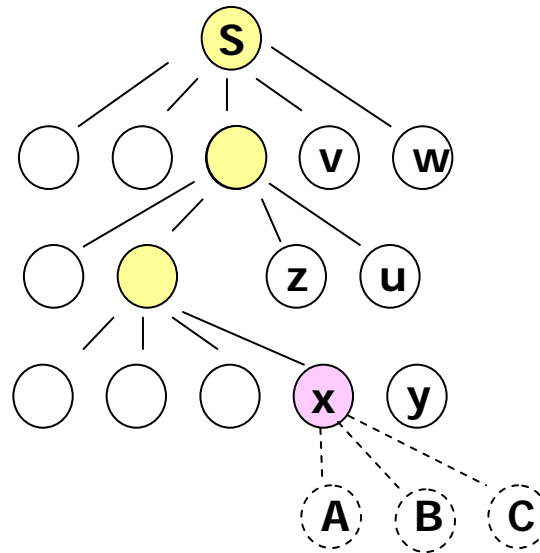


Starter med S og slutt-symbolet på stakken (End-Of-File EOF).

På stakken har vi det vi ikke er "ferdig med"

Steget i LL(1)- parsing med stakk

La X være på toppen av stakken og "a" er nåværende token (neste på input)



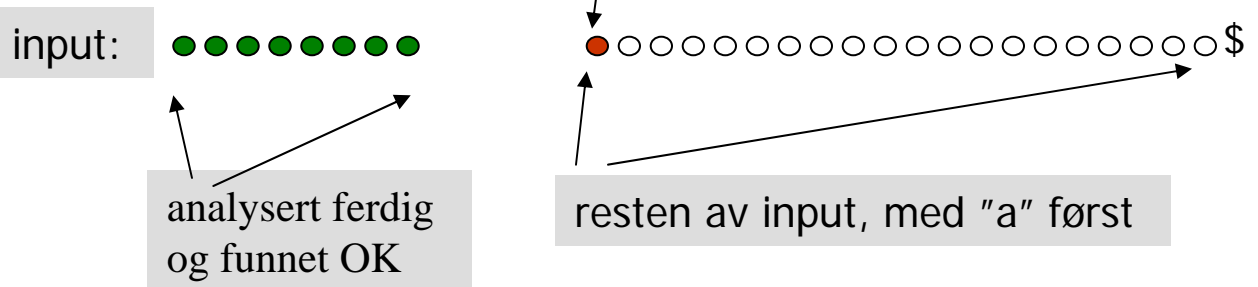
STEGET:

Om X er en terminal:

Da må $X == a$. Fjern X (pop) og les a (ikke vist i figur)

Om X ikke-terminal:

Velg aksjon $M[X,a]$. Om ikke "error", så fjern X (pop X). Anta $M[X,a]$ er "X \rightarrow ABC" og legg da (Push) ABC på stakken (N.B. vi har ikke lest "a" nå)



Stakk før:



Stakk etter:



Dette er en setningsform som oppstår under venstre-avledning fra startesymbolet til input-setningen.

Utførelse av LL(1)-persering med tabell $M[N,T]$ og eksplisitt stakk . Merk: noen feil må rettes.

```
(* assumes $ marks the bottom of the stack and the end of the input *)
push the start symbol onto the top of the parsing stack ;
while the top of the parsing stack  $\neq$  $ and the next input token  $\neq$  $ do
  if the top of the parsing stack is terminal a
    and the next input token = a
  then (* match *)
    pop the parsing stack ;
    advance the input ;
  else if the top of the parsing is nonterminal A
    and the next input token is terminal a
    and parsing table entry  $M[A, a]$  contains
      production  $A \rightarrow X_1 X_2 \dots X_n$ 
  then (* generate *)
    pop the parsing stack ;
    for  $i := n$  downto 1 do
      push  $X_i$  onto the parsing stack ;
  else error ;
if the top of the parsing stack  $\neq$  $
  and the next input token = $
then accept
else error ;
```

Det ville stoppe for tidlig (se eksempler)

"top" er terminal

"top" er ikke-terminal

Pusk baklengs på stakken (slik at første symbol kommer på toppen)

Dette vet vi holder

Aksept hvis vi greide å tolke hele input som S
Hit hvis input er "for lang"

i
 g

$M[N, T]$	if	other	else	0	1	\$
<i>statement</i>	<i>statement</i> → <i>if-stmt</i>	<i>statement</i> → other				
<i>if-stmt</i>	<i>if-stmt</i> → if (<i>exp</i>), <i>statement</i> <i>else-part</i>					
<i>else-part</i>			1 <i>else-part</i> → else <i>statement</i> 2 <i>else-part</i> → ϵ			<i>else-part</i> → ϵ
<i>exp</i>				<i>exp</i> → 0	<i>exp</i> → 1	

LL(1)-parsering med en stakk ut fra en tabell som ikke er LL(1).
Må gjøre valg i tvetydige situasjoner

Parsing stack	Input	Action
\$ S	i (0) i (1) o e o \$	$S \rightarrow I$
\$ I	i (0) i (1) o e o \$	$I \rightarrow i (E) S L$
\$ L S) E (i	i (0) i (1) o e o \$	match
\$ L S) E ((0) i (1) o e o \$	match
\$ L S) E	0) i (1) o e o \$	$E \rightarrow 0$
\$ L S) 0	0) i (1) o e o \$	match
\$ L S)) i (1) o e o \$	match
\$ L S	i (1) o e o \$	$S \rightarrow I$
\$ L I	i (1) o e o \$	$I \rightarrow i (E) S L$
\$ L L S) E (i	i (1) o e o \$	match
\$ L L S) E ((1) o e o \$	match
\$ L L S) E	1) o e o \$	$E \rightarrow 1$
\$ L L S) 1	1) o e o \$	match
\$ L L S)) o e o \$	match
\$ L L S	o e o \$	$S \rightarrow o$
\$ L L o	o e o \$	match
\$ L L	e o \$	$L \rightarrow e S$
\$ L S e	e o \$	match
\$ L S	o \$	$S \rightarrow o$
\$ L o	o \$	match
\$ L	\$	$L \rightarrow \epsilon$
\$	\$	accept



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.



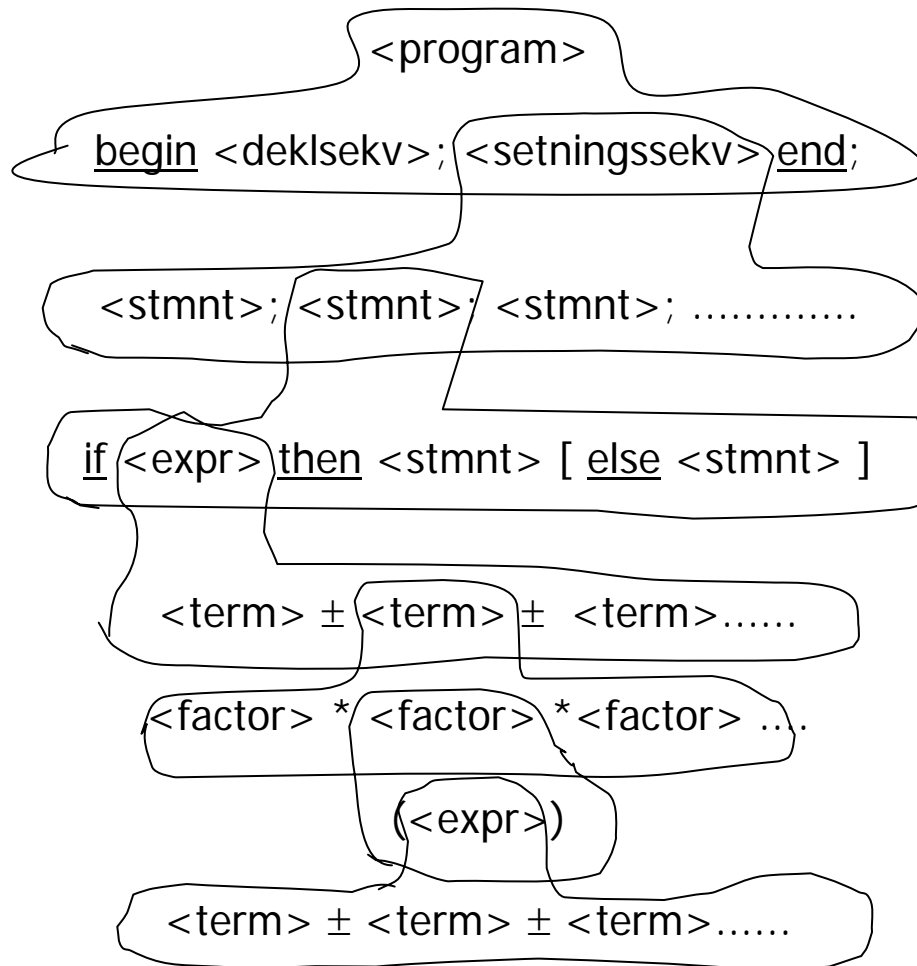
Spesielt viktig ved syntaksfeil:

- 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" parsering.

Metode: "Panic mode" og synkroniserings-mengde



Synchset (stakk eller parameter):

end

; First(stmt)

↙ navn if while for ...

then First(stmt) else

+ - First(term)

↙ (tall navn

* First(factor)

)

+ - (tall navn



Syntaksfeil ved "rec. descent" – 2

Ut fra skissen er det greit å finne:

- hvem som skal ta opp tråden
- "hvor" denne skal gjøre av det

Algoritme:

For hvert input-symbol framover:

Let gjennom stakken topp → bunn

Om finnes:

Den tilsvarende proc" skal ta opp tråden

Den vet selv hvor den skal fortsette, med dette symbolet

Om input-symbol ikke er i stakken:

- gå til neste input-symbol

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 ( { } ) ); ← Hvorfor ikke også "synchset"?
        match ( ( );
    number :
      match ( number );
    else error ;
    end case ;
    checkinput ( synchset, { (, number } ); *
  end if ;
end factor ;
```

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

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



Pensum i kap. 4

Blir lagt ut på hjemmesidens pensumliste!

Hele kapittelet, med følgende modifikasjoner:

- side 159-160 : Ikke detaljene i "Case 3", men man skal vite at en slik algoritme finnes.
- Ikke avsnitt 4.2.4 side 166-167
- Sidene 168 og 173: Vi bruker litt mer intuitive definisjoner på First og Follow (se foilene og trykkfeillisten, men også diskusjonen på punktene 1. og 2. øverst på s. 178 der koblingen til de intuitive definisjonene er gjort)
- Ikke avsnittene 4.5.2 og 4.5.3 side 186 – 189