

Kap. 5, del 1: Parsering nedentfra-opp
(Bottom-up parsering)
INF5110

Pluss oppgaver til kap. 4., se neste side.
De gjennomgås 21. febr.

14/2-2012
Stein Krogdahl
Ifi, UiO

Oppgaver til kapittel 4

Gjennomgås tirsdag 21. febr.

Oppgave 1: Sjekk om grammatikken " $S \rightarrow (S) S \mid \varepsilon$ " er LL(1) (etter formelle krav)

Oppgave 2: Gitt gram.: $\text{exp} \rightarrow \text{exp} + \text{exp} \mid (\text{exp}) \mid \text{if exp then exp else exp} \mid \text{var}$

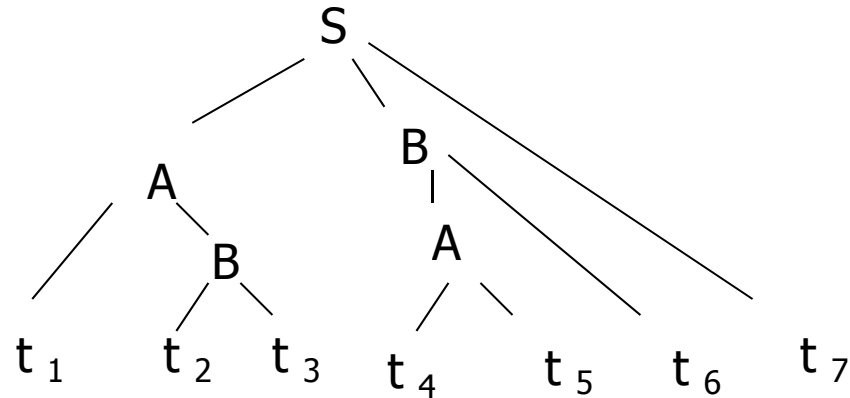
- Lag en entydig grammatikk for dette språket, der + skal være venstreassosiativ, og der "if x then y else z+u" skal bety "if x then y else (z+u)".
- Hvorfor får vi ikke noe "dangling else"-problem her?

Oppgave 3: (Mye repetisjon, også fra kap 3).

Gitt gram.: $\text{exp} \rightarrow \text{exp op exp} \mid (\text{exp}) \mid \text{num}$
 $\text{op} \rightarrow + \mid - \mid * \mid / \mid ** \mid < \mid =$

- Grammatikken over er opplagt flertydig. Lag en *entydig* grammatikk for språket ut fra at følgende tilleggsregler:
 - ** (opphøying) har presedens 3 (høyest) og er høyre-assosiativ
 - * og / har presedens 2, og er venstre-assosiativ
 - + og - har presedens 1 og er venstre-assosiativ
 - < og = har presedens 0, og er ikke-assosiativ
- Se på grammatikken du fant under a), og skriv et syntaksdiagram (med løkker der det passer) for hver ikke-terminal. Del opp "op"-terminalene på hensiktsmessig måte.
- Lag recursive-descent prosedyrer for å sjekke programmet (med while-setninger der det passer) ut fra grammatikken fra b). Du kan bruke både "match(token)" og "getToken()" fra boka (som begge setter neste symbol inn i variabelen "token").
- Ut fra svaret på c), legg til trebyggings-setninger i prosedyren som behandler en sekvens av ** slik at treet får riktig høyre-assosiativ form.
- Ta hele grammatikken fra a), og gjør den fri for venstreassosiativitet, og gjør all mulig venstrefaktorisering (men behold entydighet).
- Sjekk om grammatikken fra e) er LL(1).

"Bottom up" parsing (nedenfra-og-opp)



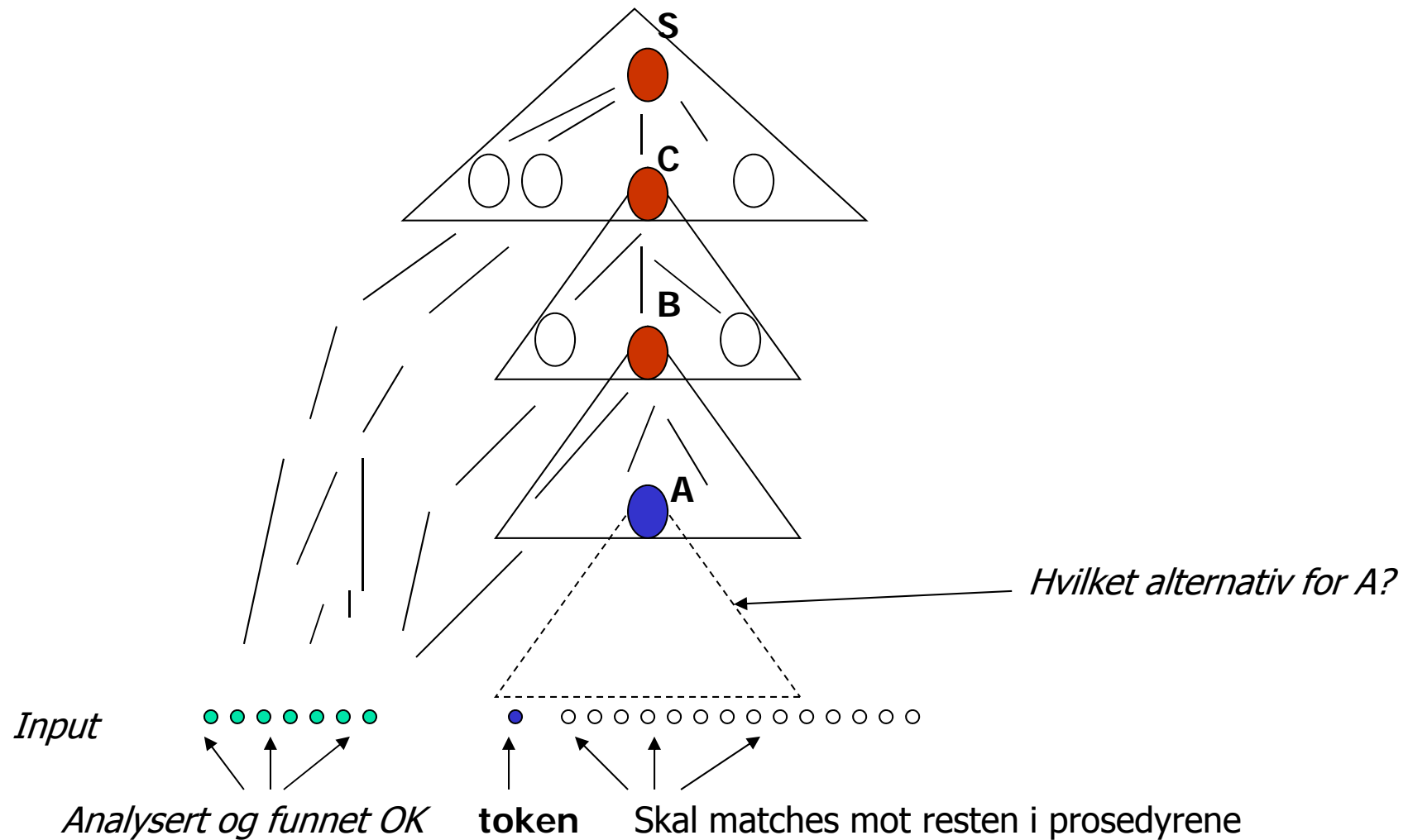
LR-parsing og grammatikker: Metodene angitt under kan i rekkefølge ta "vanskeligere og vanskeligere" grammatikker

- LR(0) Gir ca 300 "tilstander" på et vanlig programmeringsspråk
- SLR(1) -- Samme -- Vi skal se mest på denne
- LALR(1) -- Samme -- Og litt på denne
- LR(1) Gir ca 3000 "tilstander" på et vanlig programmeringsspråk

-Automatiserte verktøy som ut fra BNF-grammatikk leverer parserings-program:

-YACC, Bison, CUP (LALR(1))

Til sammenlikning: "Top down"-parsering



Datastrukturen for LR-parsing

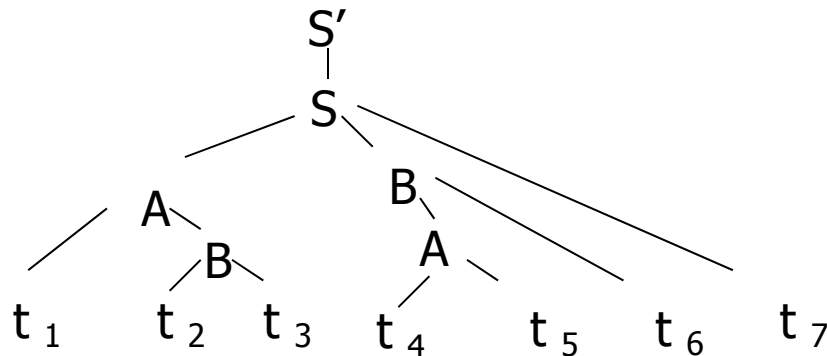
$S' \rightarrow S$

$S \rightarrow A B t_7 \mid \dots$

$A \rightarrow t_4 t_5 \mid t_1 B \mid$

$B \rightarrow t_2 t_3 \mid A t_6 \mid \dots$

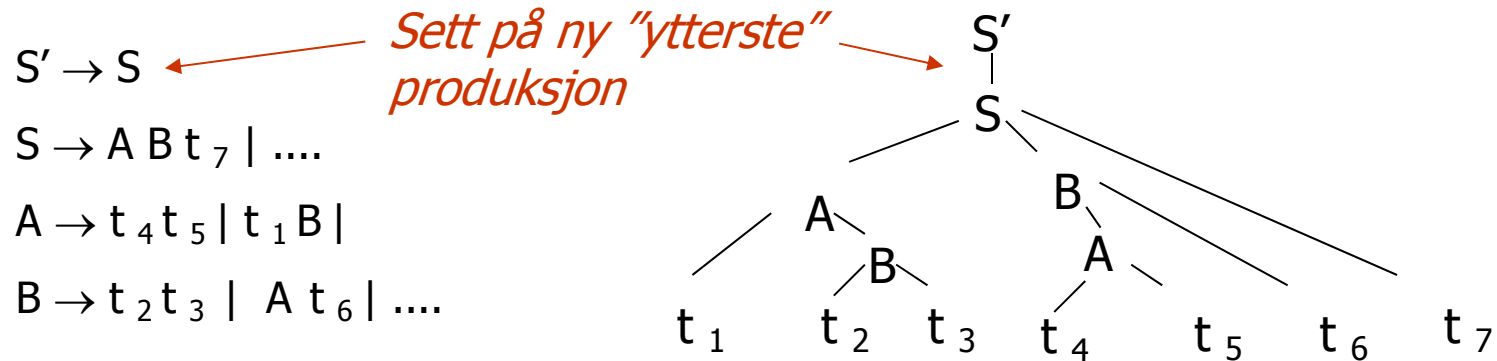
Anta at grammatikken er entydig, og at vi *kjenner syntaks-treet* for setningen:



LR-parsing :

- Ha en "stakk" for *det som er lest*
- Gjør "reduksjonen" av et subtre når det ligger "på toppen av" stakken

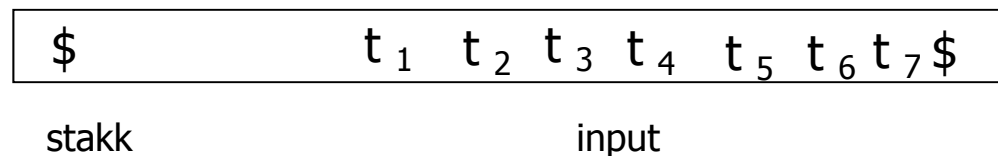
Prinsippet og datastrukturen LR-parsering II



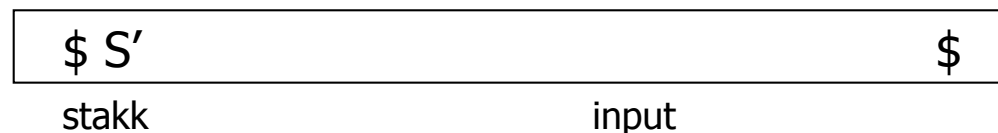
Anta at grammatikken er entydig, og anta (foreløpig!) at vi kjenner syntaks-treet for setningen:

- Ha en "stakk" for det som er lest (og er "reduert"!)
- Gjør "reduksjonen" av et subtre når subtreet ligger "på toppen av" stakken.
- Da erstatter vi dette med den ikke-terminalen som produserte dette subtreet.
- *En reduksjon:* En produksjon brukt "baklengs"

Start-situasjonen:

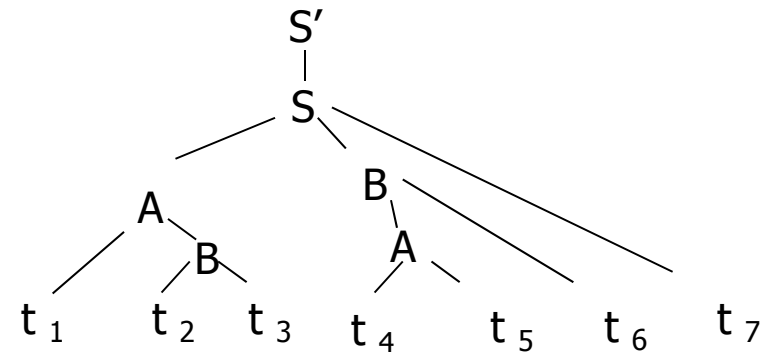


Slutt-situasjonen:



Prinsippet for LR-parsing og operasjonene: skift og redusèr

$S' \rightarrow S$
 $S \rightarrow A B t_7 \mid \dots$
 $A \rightarrow t_4 t_5 \mid t_1 B \mid$
 $B \rightarrow t_2 t_3 \mid A t_6 \mid \dots$



stakk	input
\$	t ₁ t ₂ t ₃ t ₄ t ₅ t ₆ t ₇ \$
\$ t ₁	t ₂ t ₃ t ₄ t ₅ t ₆ t ₇ \$
\$ t ₁ t ₂	t ₃ t ₄ t ₅ t ₆ t ₇ \$
\$ t ₁ t ₂ t ₃	t ₄ t ₅ t ₆ t ₇ \$
\$ t ₁ B	t ₄ t ₅ t ₆ t ₇ \$
\$ A	t ₄ t ₅ t ₆ t ₇ \$
\$ A t ₄	t ₅ t ₆ t ₇ \$
\$ A t ₄ t ₅	t ₆ t ₇ \$
\$ A A	t ₆ t ₇ \$
\$ A A t ₆	t ₇ \$
\$ A B	t ₇ \$
\$ A B t ₇	\$
\$ S	\$

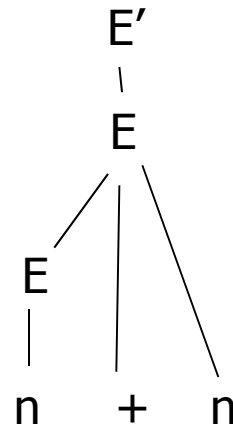
- Det blir to typer steg:
 - *Skift* ("les") neste input over til stakken
 - *Reduksjon* på toppen av stakken med bestemt regel $A \rightarrow \alpha$
 - Dersom man kjenner syntakstreet for den aktuelle setning, er det lett å angi de rette stegene.
- **MEN: Hvordan i all verden gjøre dette underveis uten å kjenne verken:**
 - resten av input
 - eller det fulle syntakstre

Husk at vi nå skal *redusere* input (bottom up) til startsymbolet S (eller S'), IKKE produsere input fra startesymbolet (slik vi gjorde "top-down" i kap 4)

Eksempel på LR-parsering (når vi kjenner treet!)

$$E' \rightarrow E$$

$$E \rightarrow E + n \mid n$$



Høyre-avledning (men ikke det tydeligste eksempel på det!)

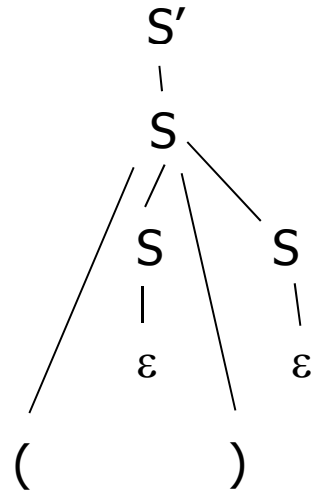
I boka, men *vi* stresser det ikke:

- Den neste reduksjonen som skal gjøres, kalles situasjonens **"handle"** (håndtak).
- **"stakk + input"** utgjør et stadium i en **høyre-avledning**, men de kommer i **omvendt rekkefølge** av en høyre-avledning

	Parsing stack	Input	Action
1	\$	n + n \$	shift
2	\$ n	+ n \$	reduce $E \rightarrow n$
3	\$ E	+ n \$	shift
4	\$ E +	n \$	shift
5	\$ E + n	\$	reduce $E \rightarrow E + n$
6	\$ E	\$	reduce $E' \rightarrow E$
7	\$ E'	\$	accept

Eksempel på LR-parsering

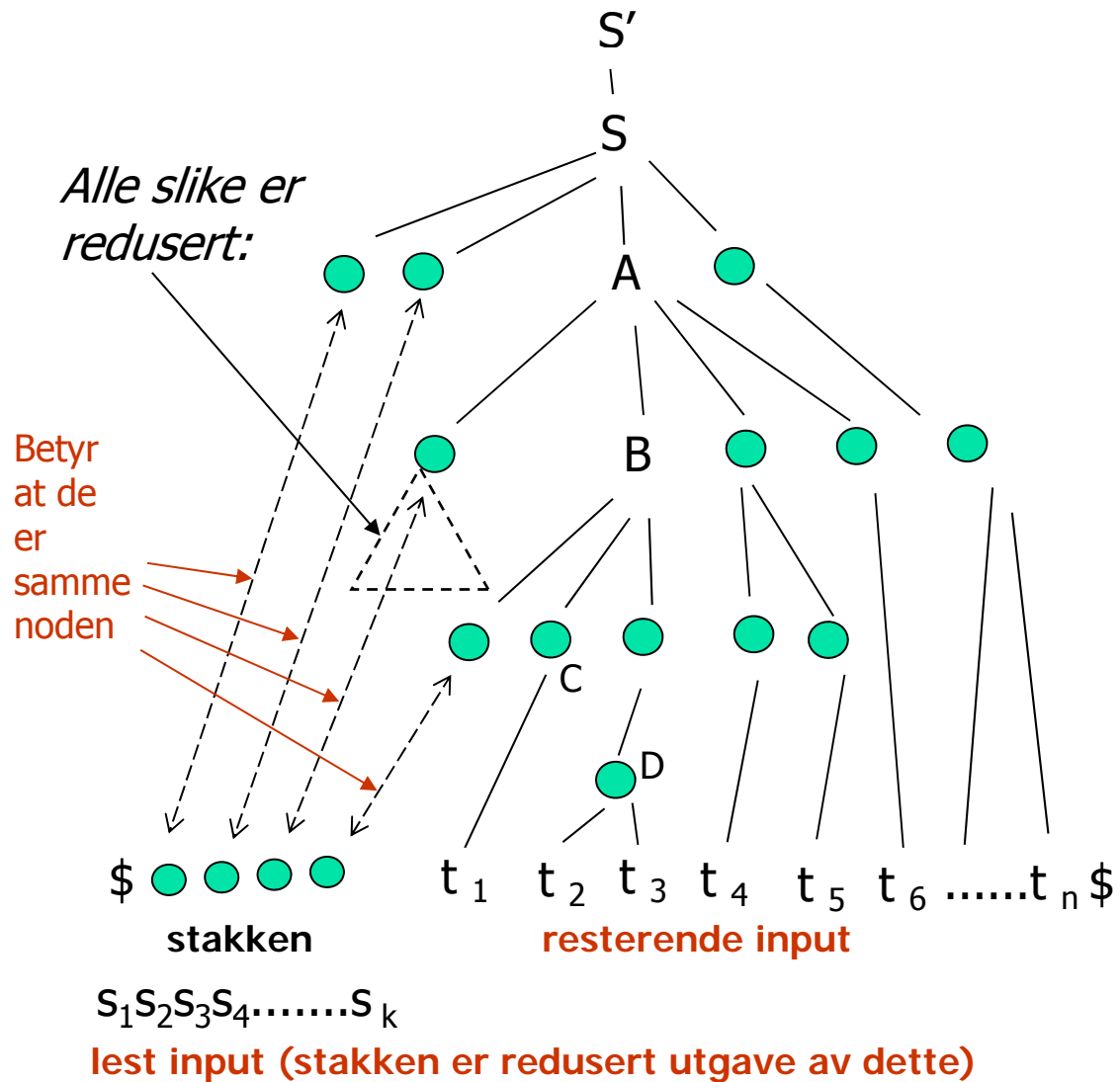
$S' \rightarrow S$
 $S \rightarrow (S) S \mid \varepsilon$



NB: $S \rightarrow \varepsilon$ blir litt "rar":
 Ved reduksjon med denne dukker en ikke-termial bare opp på enden av stakken, her og her.

	Parsing stack	Input	Action
1	\$	()\$	shift
2	\$ ()\$	reduce $S \rightarrow \varepsilon$
3	\$ (S)\$	shift
4	\$ (S)	\$	reduce $S \rightarrow \varepsilon$
5	\$ (S) S	\$	reduce $S \rightarrow (S) S$
6	\$ S	\$	reduce $S' \rightarrow S$
7	\$ S'	\$	accept

Typisk situasjon under LR-parsering



Det neste reduksjonen som blir gjort er reduksjon med

$C \rightarrow t_1$

Deretter (etter noen skift):

$D \rightarrow t_2 t_3$

Gitt en entydig grammatikk G.

Denne behandles som følger for å lage en LR-parser:

(Ikke alt er beskrevet i boka, og bare det i boka er pensum)

- Vi ser på *alle* mulig stakker som kan opptre, og ser disse som strenger over alfabetet: $\{\text{terminaler}\} \cup \{\text{ikke-terminaler}\}$. Vi ser altså på:
 $\{S \mid S \text{ utgjør stakken på ett eller annet stadium i LR-parsering av en setning i } L(G)\}$
- Dette språket viser seg å være **regulært**, og kan beskrives av en NFA der **alle tilstander er aksepterende**:
 - Tilstandene angis av "itemer" av formen: $A \rightarrow XY \cdot Z$ (mer presist: "LR(0)-itemer")
 - Kantene kan beskrives nokså greit (kommer snart)
- Denne NFA-en gjør vi om til en DFA på vanlig måte (subset construction fra kap. 2)
- Tilstandene i DFA-en er mengder av NFA-tilstander, og vil her altså være **mengder av itemer**
- **LR-parseringens hovedsetning (uformelt)**:
 - Se på en situasjon under LR-parsering av en setning, der vi har strengen S på stakken.
 - Strengen S vil da aksepteres av DFA-en angitt over. Anta at den bringer oss til DFA-tilstanden T.
 - Da vil mengden av *itemer* i T angi de mulige "lokale forhold" vi har ved det punktet vi nå er i i parseringen.
 - Selv om grammatikken er entydig er det ofte flere muligheter/valg (ett for hvert item). Dette kommer av at vi ennå ikke har sett på resten av input.



Altså "Itemer": Hver produksjon i BNF-grammatikken gir opphav til et antall slike.

- For produksjonen :

$A \rightarrow XYZ$

- Lager vi itemer med punktum på alle plasser (punktumet er et Meta-symbol):

$A \rightarrow .XYZ$

$A \rightarrow X.YZ$

$A \rightarrow XY.Z$

$A \rightarrow XYZ.$

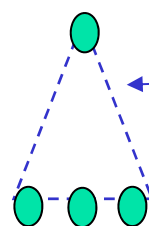
- Punktumet angir grovt sett hvor vi er hen i "dypeste aktive del" av parseringen i øyeblikket:
 - Det til venstre for punktum er lest fra input, enten bare lest eller at deler av det er redusert til en ikke-terminal
 - Det til høyre for punktum har vi enda ikke sett/lest

Kantene i NFA-en

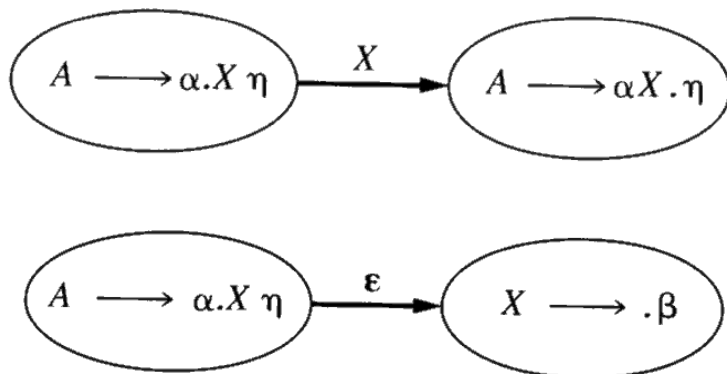
Dette er ikke fullt forklart i boka, og er ikke pensum

Skal lage en NFA som aksepterer alle slike "linjer", og intet annet

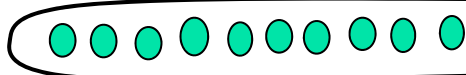
Gitt en grammatikk G og en situasjon under passering av en setning s i $L(G)$.



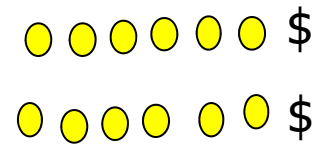
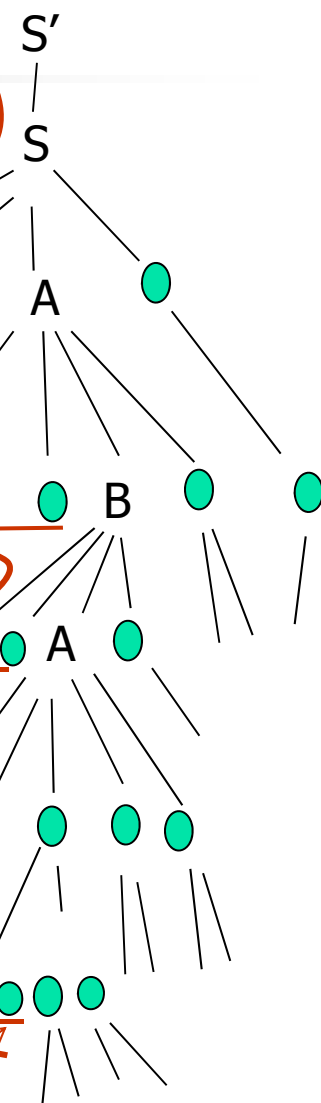
Husk: Ingen slik på venstre side. De er redusert!



Stakk:



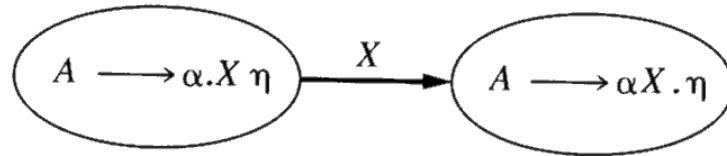
Opprinnelig input:



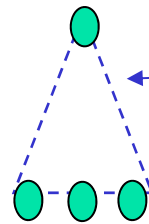
Mer: Kantene i NFA-en

Dette er ikke fullt forklart i boka, og er ikke pensum

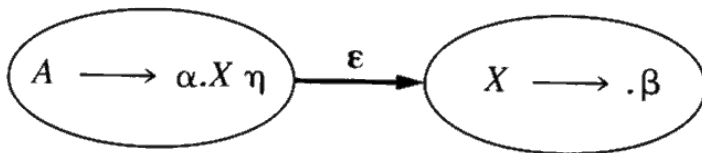
Gitt en grammatikk G og en situasjon under passering av en setning s i L(G).



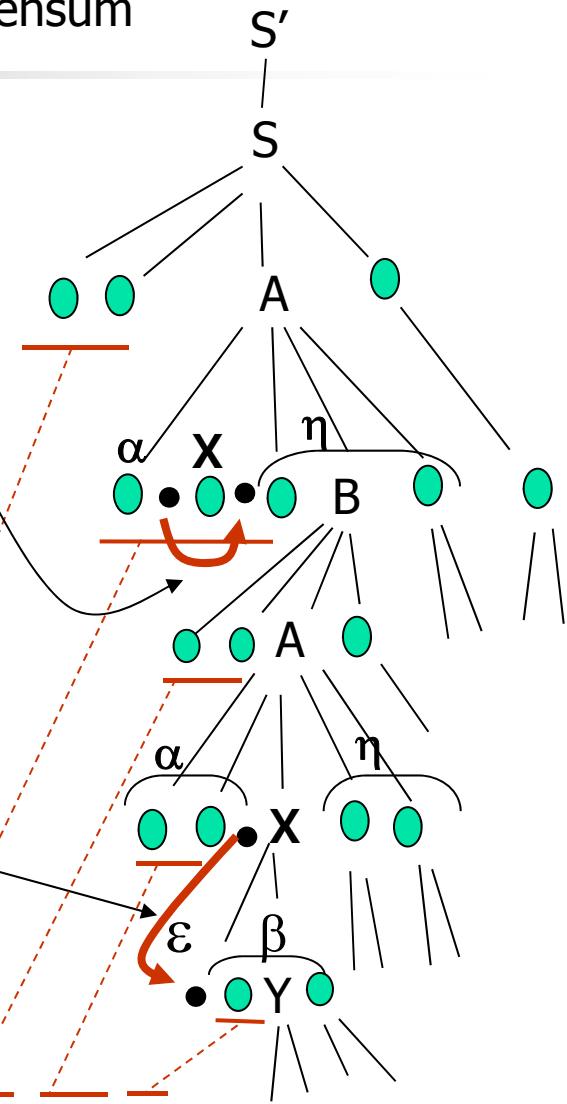
tilsvare



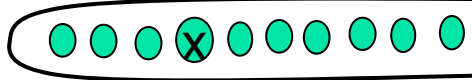
Husk: Ingen slik på venstre side. De er redusert!



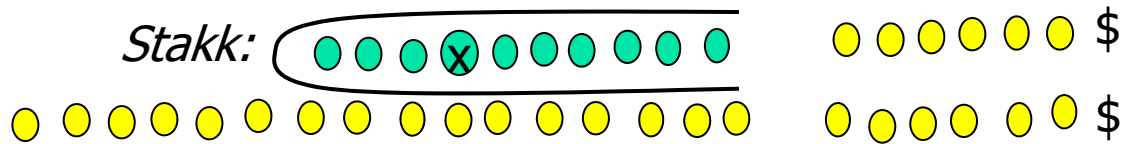
tilsvare



Stakk:



Opprinnelig input:

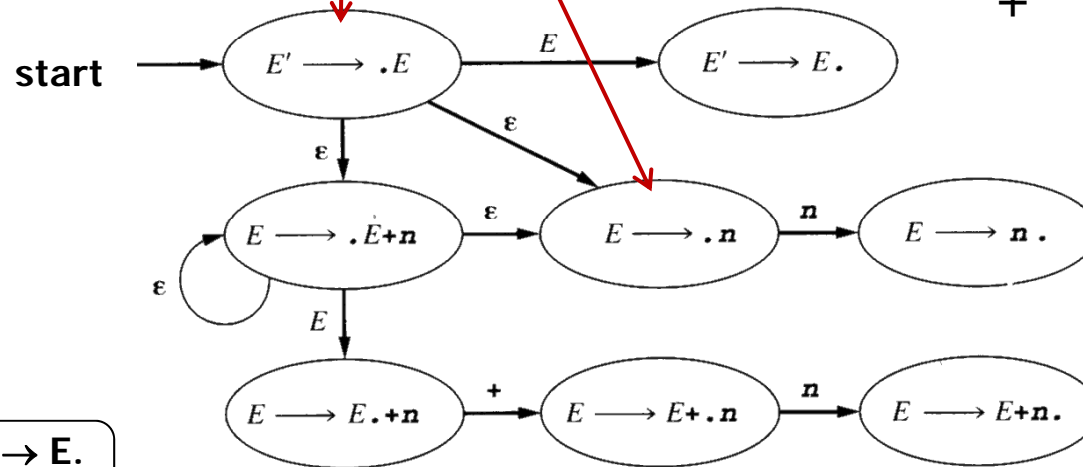
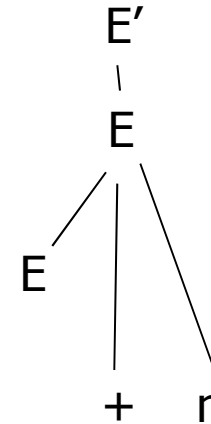


NFA'en som beskriver alle mulige stakker

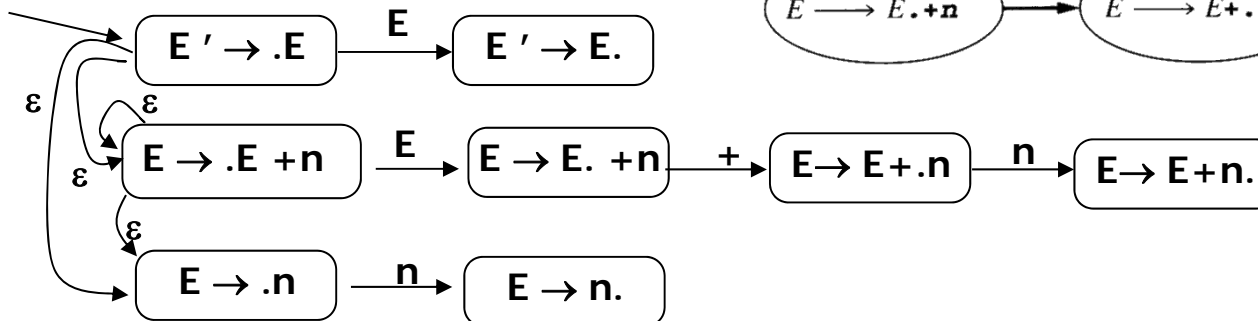
$E' \rightarrow E$
 $E \rightarrow E + n$
 $E \rightarrow n$

$E' \rightarrow .E$
 $E' \rightarrow E.$
 $E \rightarrow .E + n$
 $E \rightarrow E. + n$
 $E \rightarrow E + .n$
 $E \rightarrow E + n.$
 $E \rightarrow .n$
 $E \rightarrow n.$

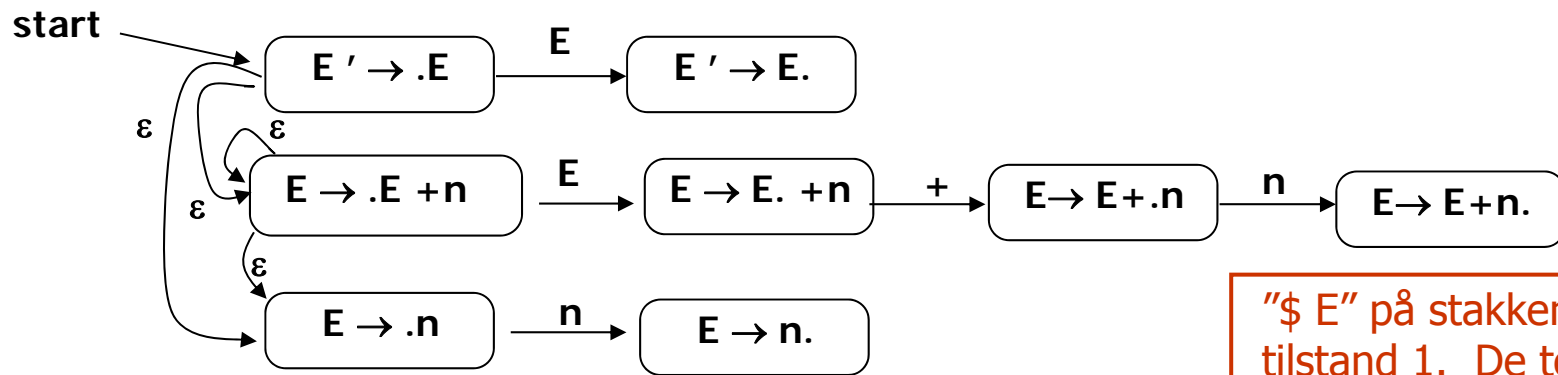
Itemer (LR(0)-itemer)



Den samme, i litt ryddigere utgave:



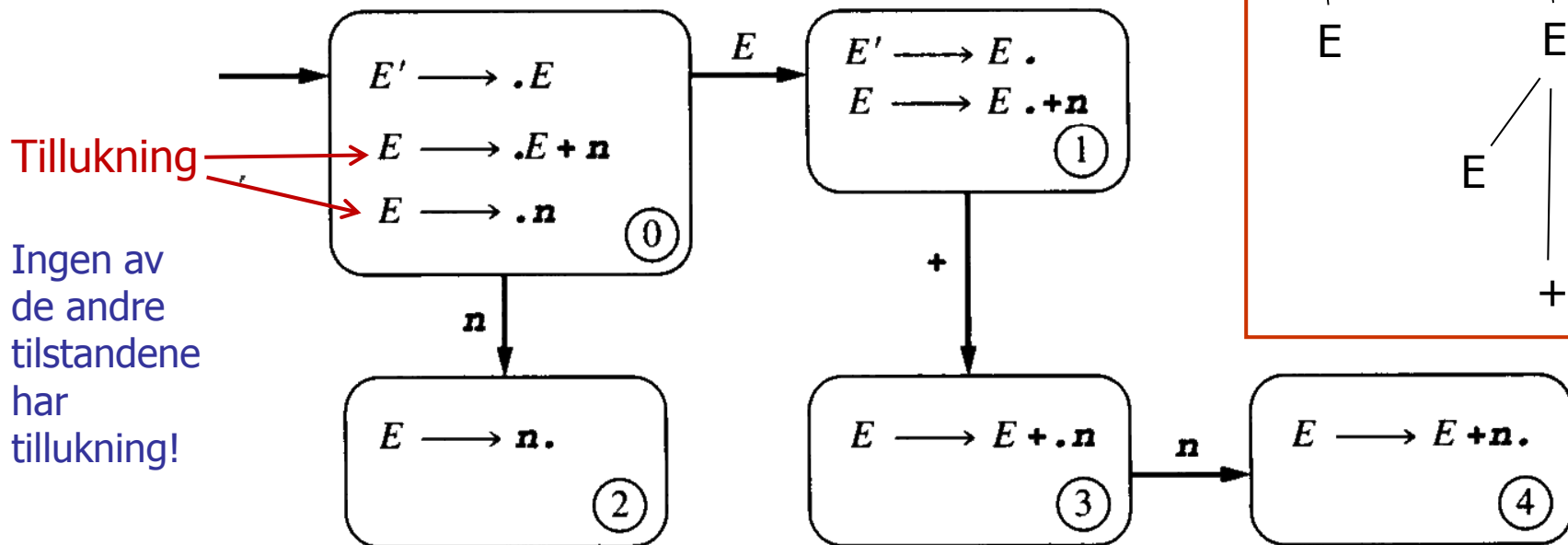
LR(0)-NFA (Altså litt mer systematisk enn i boka):



"\$ E" på stakken gir tilstand 1. De to mulighetene kan da eksemplifiseres ved:

The first parse tree shows E' deriving E . The second parse tree shows E' deriving E , which then derives E , $+$, and n .

LR(0)-DFA, laget ved subset-konstruksjonen (kap 2):



Hvordan sette opp LR(0)-DFA'en direkte fra grammatikken? (Rett fram bruk av subset-konstruksjonen. Bør trenes på!)

- Tillukning av item-mengde I:

- Dersom:

- $A \rightarrow \alpha \cdot B \gamma$ er med i I
- B er en ikke-terminal
- $B \rightarrow \beta_1 \mid \beta_2 \mid \dots$

Da er også Itemene:

- $B \rightarrow \cdot \beta_2$
- $B \rightarrow \cdot \beta_1$
- ...

med i I

- Start-tilstanden til LR(0)-DFA'en er:

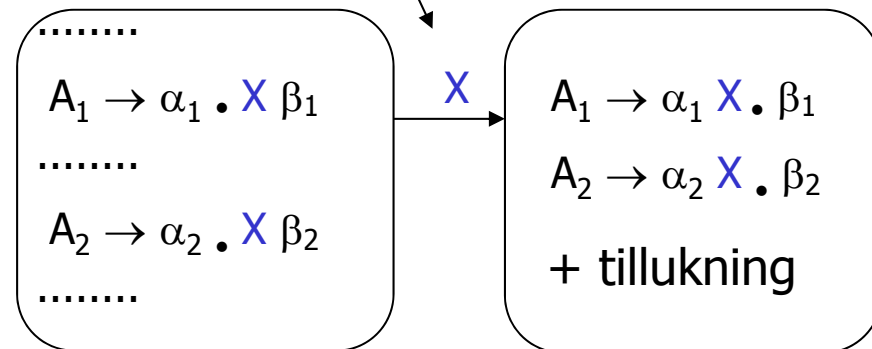
$S' \rightarrow \cdot S$

+ tillukning

- Tilstandsovergang ved symbol X fra tilstand I

- X er terminal eller ikke-terminal (behandles likt!)

Tilstand I:



Få med *alle* av formen $A \rightarrow \alpha \cdot X \beta$

Hvordan sette opp LR(0)-DFA'en direkte fra grammatikken - II

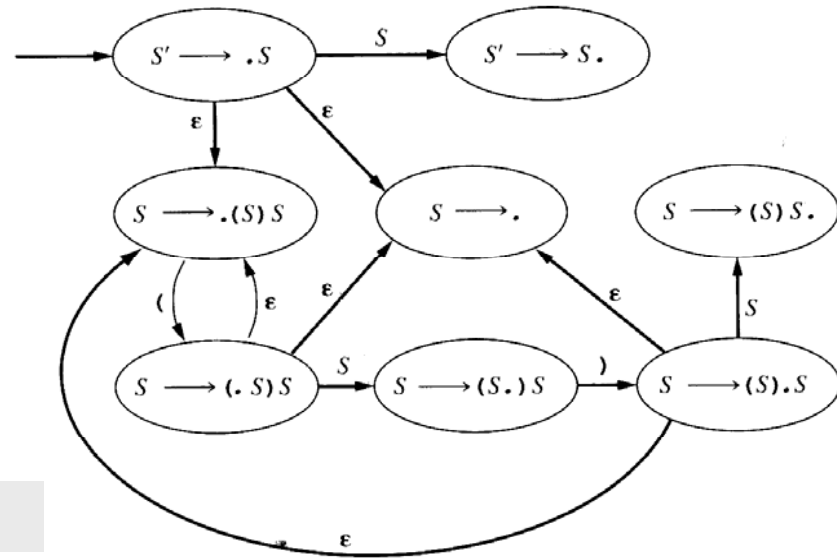
$S' \rightarrow S$
 $S \rightarrow (S) S \mid \epsilon$

Merk at $S \rightarrow \epsilon$ bare gir ett item, nemlig: $S \rightarrow \bullet$

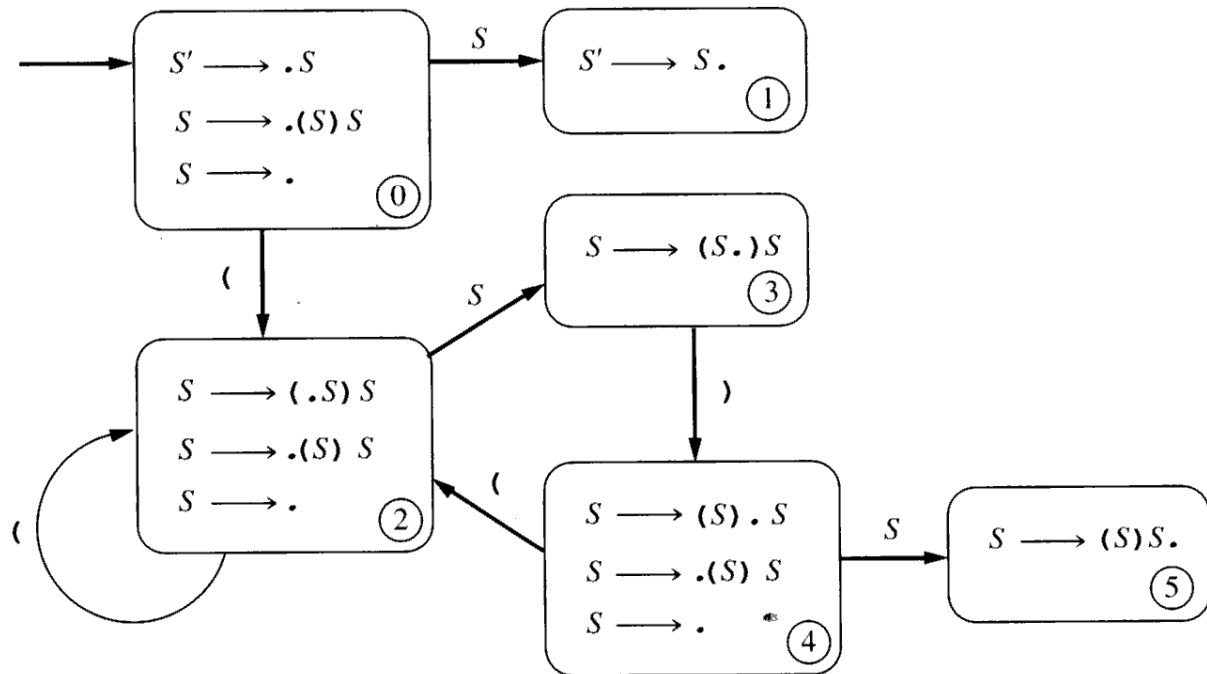
(Altså ikke: ~~$S \rightarrow \bullet \epsilon$~~
og ~~$S \rightarrow \epsilon \bullet$~~)

- $S' \rightarrow \bullet S$
- $S' \rightarrow S \bullet$
- $S \rightarrow \bullet (S) S$
- $S \rightarrow (\bullet S) S$
- $S \rightarrow (S \bullet) S$
- $S \rightarrow (S) \bullet S$
- $S \rightarrow (S) S \bullet$
- $S \rightarrow \bullet$

LR(0) – NFA:



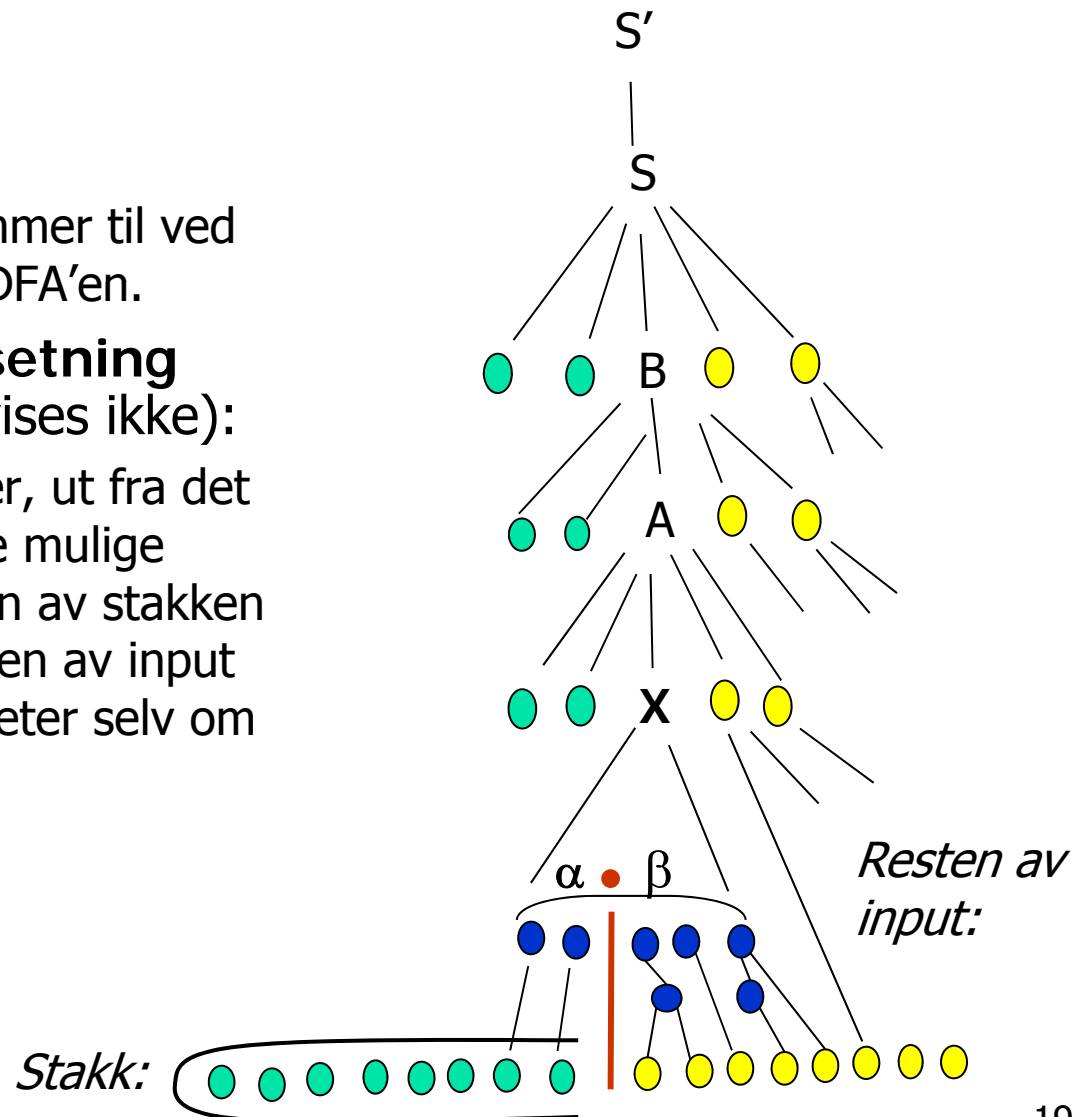
LR(0) – DFA:



Hva sier "topp-tilstanden"?

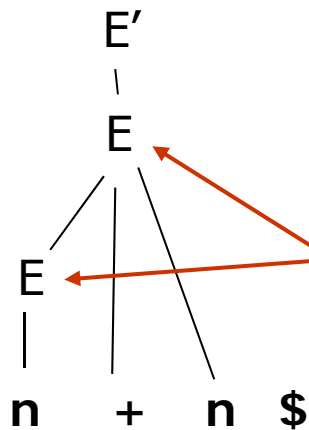
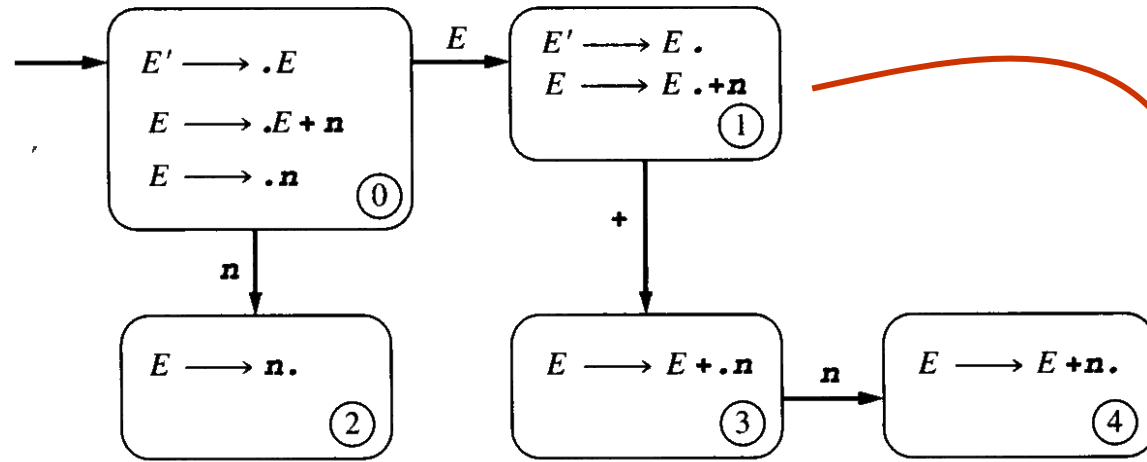
- **Topp-tilstanden:**
 - Den DFA-tilstanden vi kommer til ved å la stakken gå gjennom DFA'en.
- **LR-parseringens hovedsetning** (En gang til, litt løselig, bevises ikke):
 - Itemene i topp-tilstanden er, ut fra det vi har lest av input, *alle* de mulige "lokale forhold" ved toppen av stakken (Siden vi ikke har lest resten av input kan det være flere muligheter selv om grammatikken er entydig)

Dersom itemet $X \rightarrow \alpha \bullet \beta$ er med i topp-tilstanden kan situasjonen altså være som angitt til høyre.



Hva sier topp-tilstanden? Eksempel:

$E' \rightarrow E$
 $E \rightarrow E + n \mid n$



\$		$n + n$	\$
\$	n	$+ n$	\$
\$	E	$+ n$	\$
\$	$E +$	n	\$
\$	$E + n$		\$
\$	E		\$
\$	E'		\$

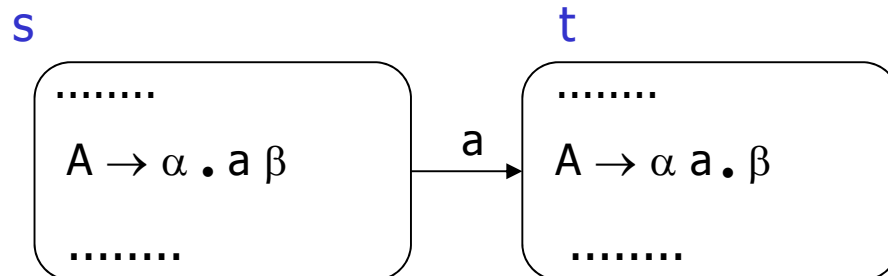
Skal skifte (between \$E\$ and \$E+\$)
 Skal redusere med $E' \rightarrow E$ (between \$E\$ and \$E'\$)

LR(0)-grammatikk: Topptilstanden bestemmer alltid entydig neste steg.

Derved: Denne grammatikken er ikke LR(0). For stakk "E" er neste skritt ubestemt. 20

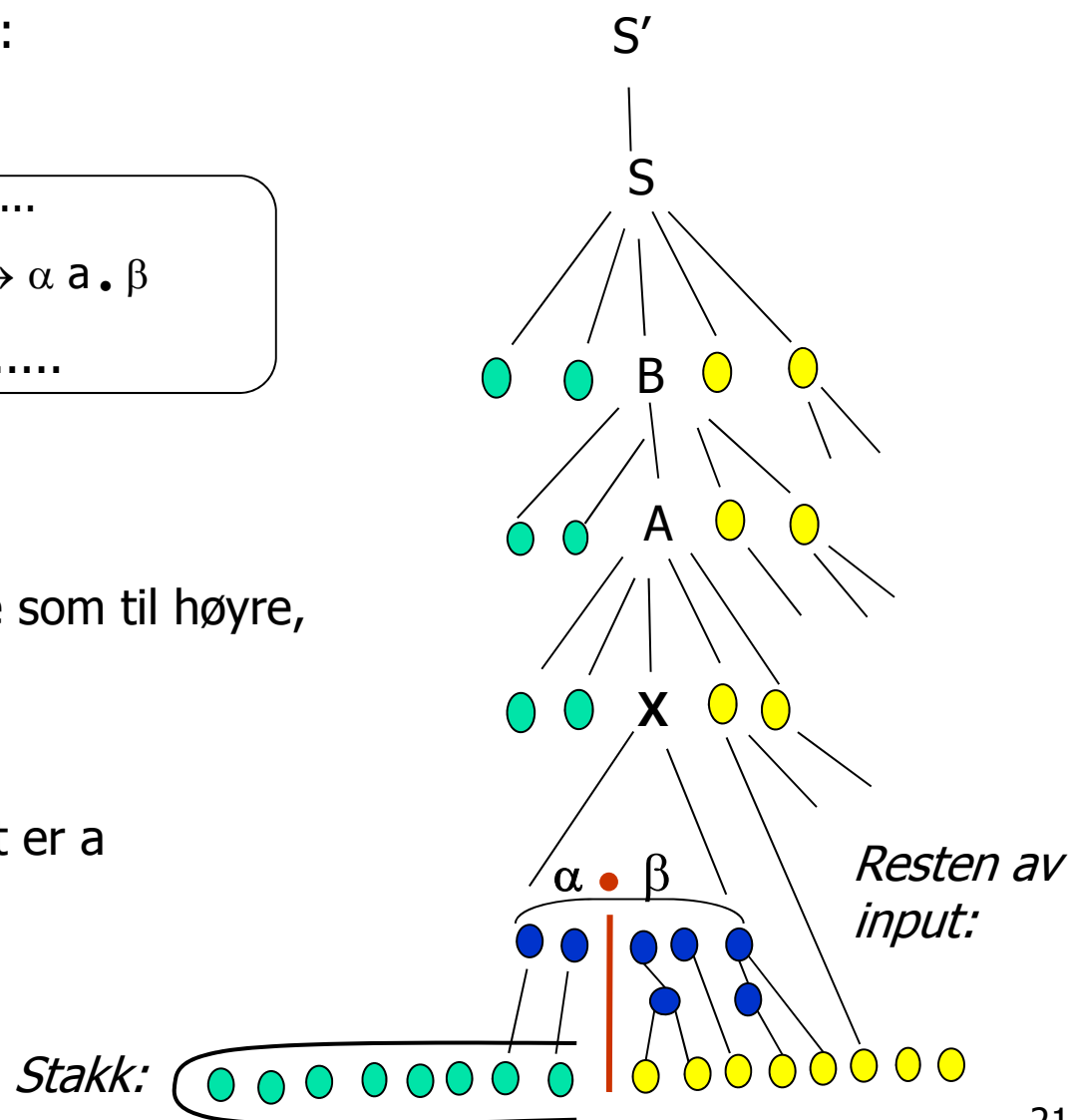
Når er "skift" en mulighet?

Anta vi er i tilstand s under:



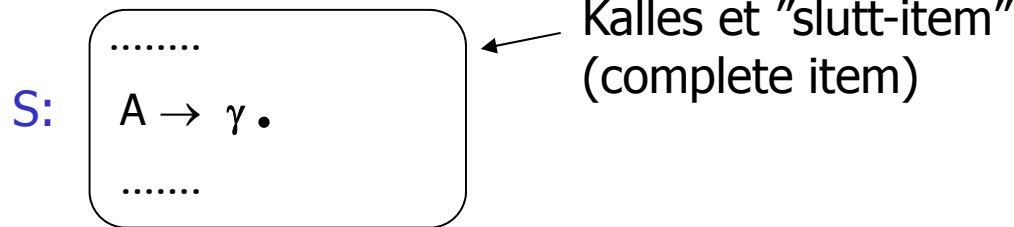
Angir at situasjonen *kan* være som til høyre, og derved:

- Neste skritt kan være skift
- Skift er lovlig om neste input er a
- Ny topptilstand blir i så fall t



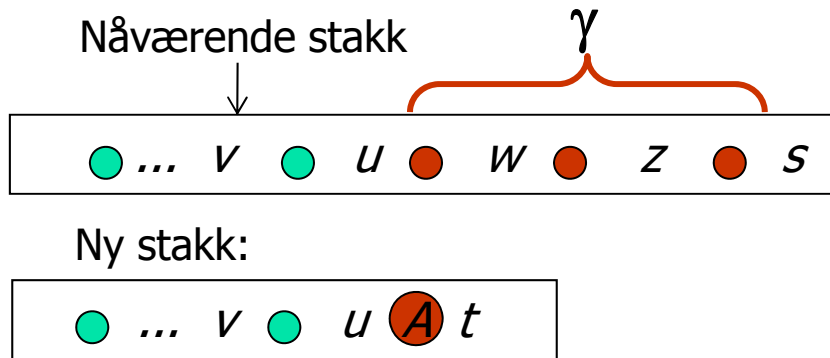
LR(0)-grammatikker og LR(0)-algoritmen:

Anta at toppstilstanden er s , og at den er slik:

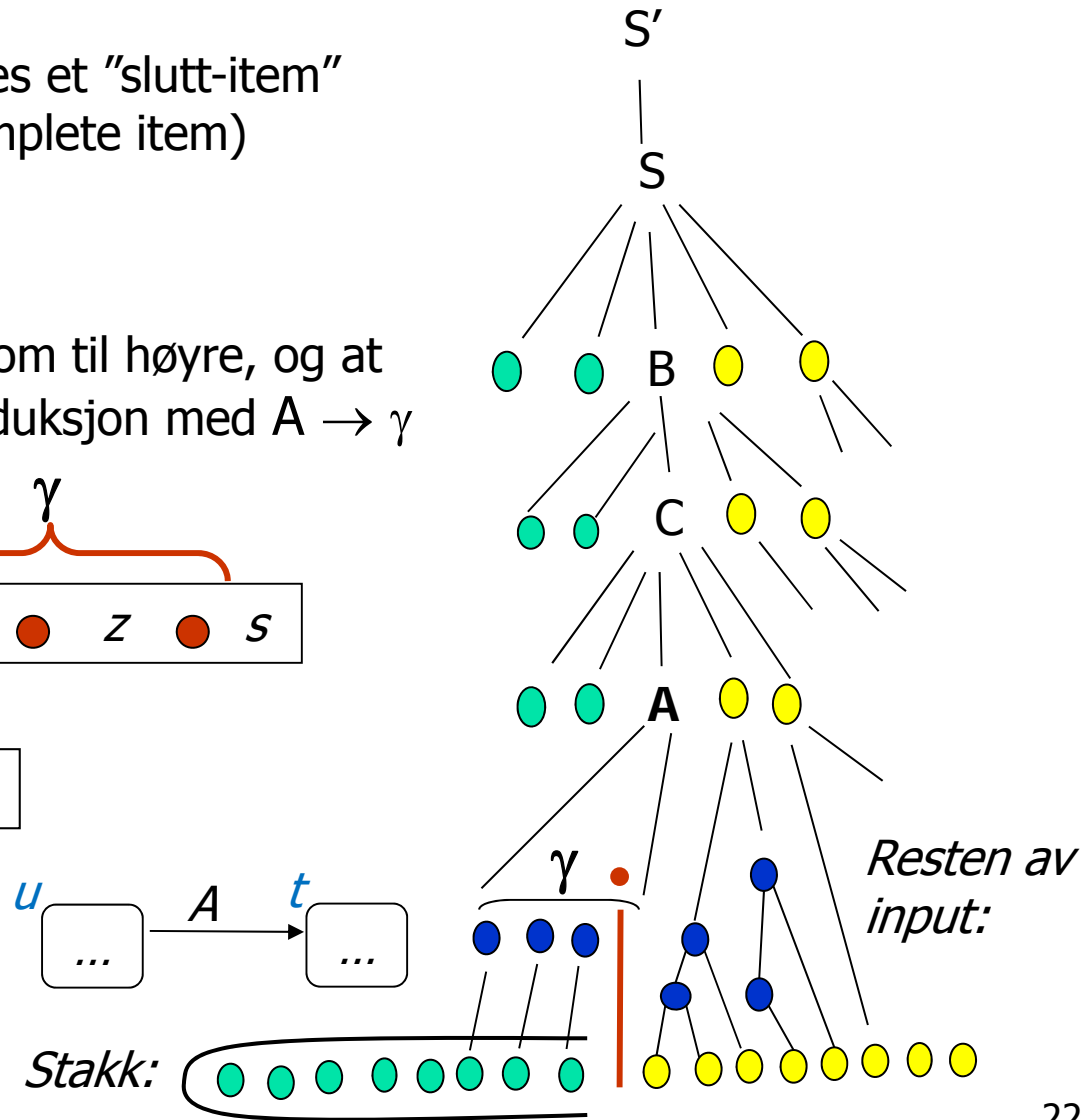


Angir at situasjonen *kan* være som til høyre, og at neste skritt derved kan være reduksjon med $A \rightarrow \gamma$

NB: Vi holder tilsandene mellom stakk-symbolene



Reduser-steget: Pop av det som tilsvarer γ (og mellomliggende tilstander), og push på A , og finn riktig ny toppstilstand ut fra u og A



Er eksempel-grammatikkene i Kap 5 LR(0) ?

Ser på tre grammatikker:

$$A \rightarrow (A) \mid a$$

$$S' \rightarrow S$$

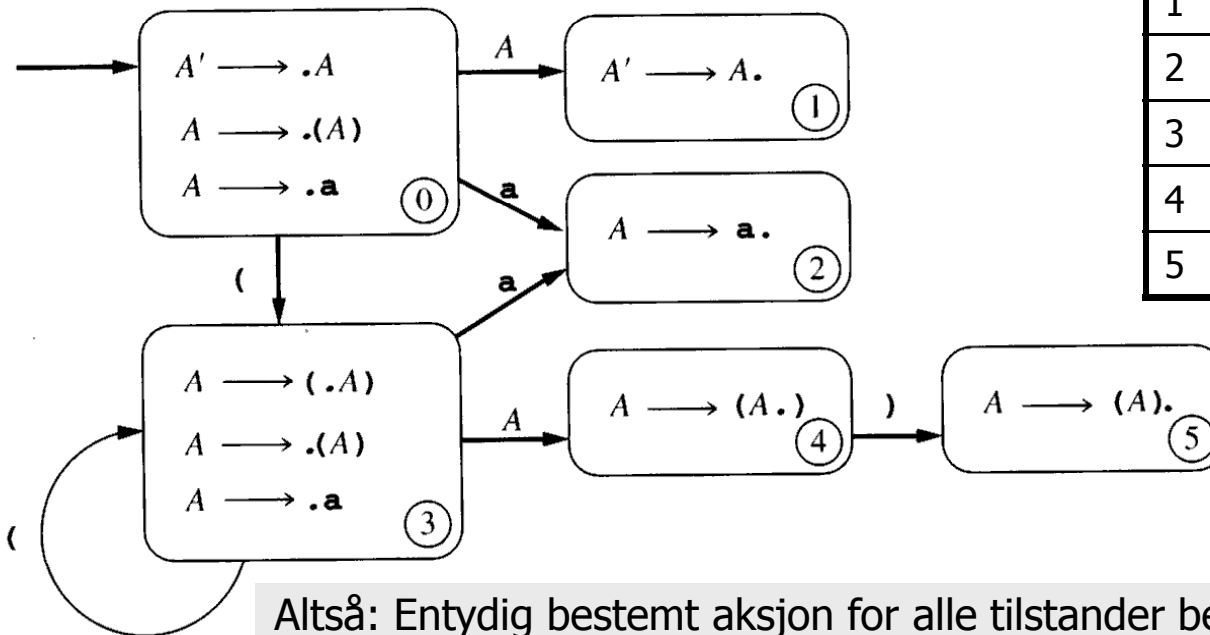
$$S \rightarrow (S) S \mid \varepsilon$$

$$E' \rightarrow E$$

$$E \rightarrow E + n \mid n$$

Har vi sett på:
Er ikke LR(0)

A gir følgende LR(0) – DFA:



Tilst.	Mulig aksjoner:
0	Bare skift mulig, for "(" og "a"
1	Bare red. mulig, med $A' \rightarrow A$
2	Bare red. mulig, med $A \rightarrow a$
3	Bare skift mulig, for "(" og "a"
4	Bare skift mulig, for ")"
5	Bare red. mulig, med $A \rightarrow (A)$

Altså: Entydig bestemt aksjon for alle tilstander betyr: **Grammatikken er LR(0)**

MERK: Der det er reduksjon må det ikke være tvil om med hvilken produksjon!

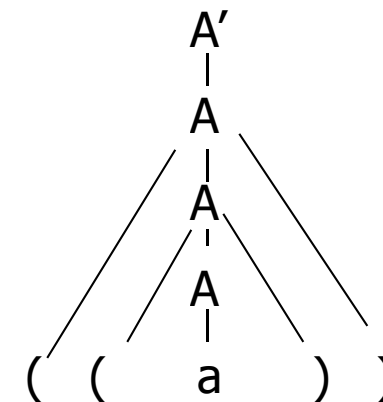
Tabell-oppsett for en LR(0) - grammatikk:

State	Action	Rule	Input			Goto
			(a)	
0	shift		3	2		1
1	reduce	$A' \rightarrow A$				0
2	reduce	$A \rightarrow a$				0
3	shift		3	2		4
4	shift				5	
5	reduce	$A \rightarrow (A)$				0

Hvis en reduksjon bringer oss tilbake til tilstand 0 eller 3, sier Goto hvilken tilstand A gir.

Parsering av setningen: ((a))

	Parsing stack	Input	Action
1	\$ 0	((a)) \$	shift
2	\$ 0 (3	(a)) \$	shift
3	\$ 0 (3 (3	a)) \$	shift
4	\$ 0 (3 (3 a 2)) \$	reduce $A \rightarrow a$
5	\$ 0 (3 (3 A 4)) \$	shift
6	\$ 0 (3 (3 A 4) 5) \$	reduce $A \rightarrow (A)$
7	\$ 0 (3 A 4) \$	shift
8	\$ 0 (3 A 4) 5	\$	reduce $A \rightarrow (A)$
9	\$ 0 A 1	\$	accept



Skal her redusere med $A' \rightarrow A$, og input tom: Ferdig

Parsering av noen gale strenger for: $A \rightarrow (A) \mid a$

State	Action	Rule	Input			Goto
			(a)	
0	shift		3	2		1
1	reduce	$A' \rightarrow A$				
2	reduce	$A \rightarrow a$				
3	shift		3	2		4
4	shift				5	
5	reduce	$A \rightarrow (A)$				

\$ 0	((a) \$
\$ 0 (3	(a) \$
\$ 0 (3	(a) \$
\$ 0 (3 (3	a) \$
\$ 0 (3 (3 a) \$
\$ 0 (3 (3 A) 5	\$
\$ 0 (3 (3 A 4	\$

\$ 0	() \$
\$ 0 (3) \$

*Viktig: Skifter
aldri noe ulovlig
inn på stakken!*

En (litt ruskete) formulering av LR(0)-kravet: Dersom følgende gir en entydig algoritme er grammatikken LR(0):

s er en DFA-tilstand med flere itemer

The LR(0) parsing algorithm. Let s be the current state (at the top of the parsing stack). Then actions are defined as follows:

1. If state s contains any item of the form $A \rightarrow \alpha.X\beta$, where X is a terminal, then the action is to shift the current input token onto the stack. If this token is X , and state s contains item $A \rightarrow \alpha.X\beta$, then the new state to be pushed on the stack is t , where $s \xrightarrow{X} t$ ~~the state containing the item $A \rightarrow \alpha X \beta$.~~ If this token is not X for some item in state s of the form just described, an error is declared.
2. If state s contains any complete item (an item of the form $A \rightarrow \gamma$), then the action is to reduce by the rule $A \rightarrow \gamma$. A reduction by the rule $S' \rightarrow S$, where S is the start state, is equivalent to acceptance, provided the input is empty, and error if the input is not empty. In all other cases, the new state is computed as follows. Remove the string γ and all of its corresponding states from the parsing stack (the string γ must be at the top of the stack, according to the way the DFA is constructed). Correspondingly, back up in the DFA to the state from which the construction of γ began (this must be the state u uncovered by the removal of γ). Again, by the construction of the DFA, this state u must contain an item of the form $B \rightarrow \alpha.A\beta$. Push A onto the stack, and push (as the new state) the state t , where $u \xrightarrow{A} t$ ~~containing the item $B \rightarrow \alpha A \beta$.~~ (Note that this corresponds to following the transition on A in the DFA, which is indeed reasonable, since we are pushing A onto the stack.)

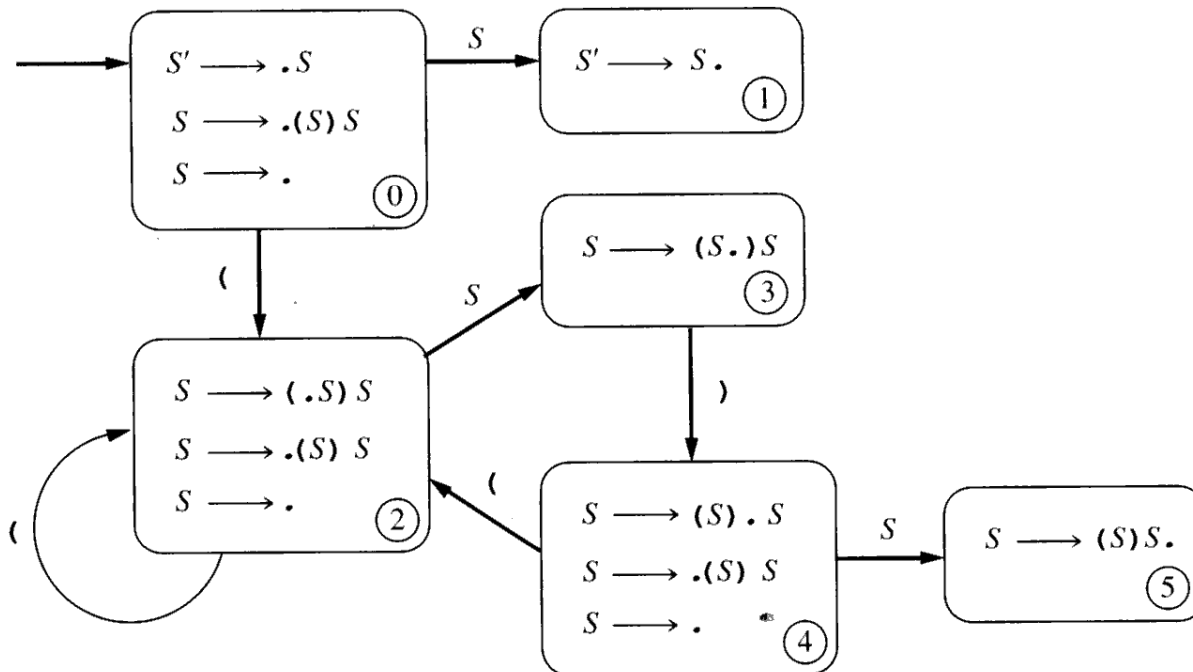
Kan forekomme i flere tilstander

Avslutning

Er denne grammatikken LR(0)?
 Nei, pga. tilstandene 0, 2 og 4

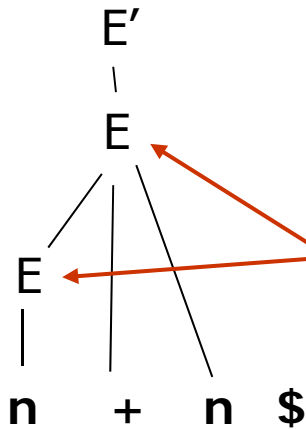
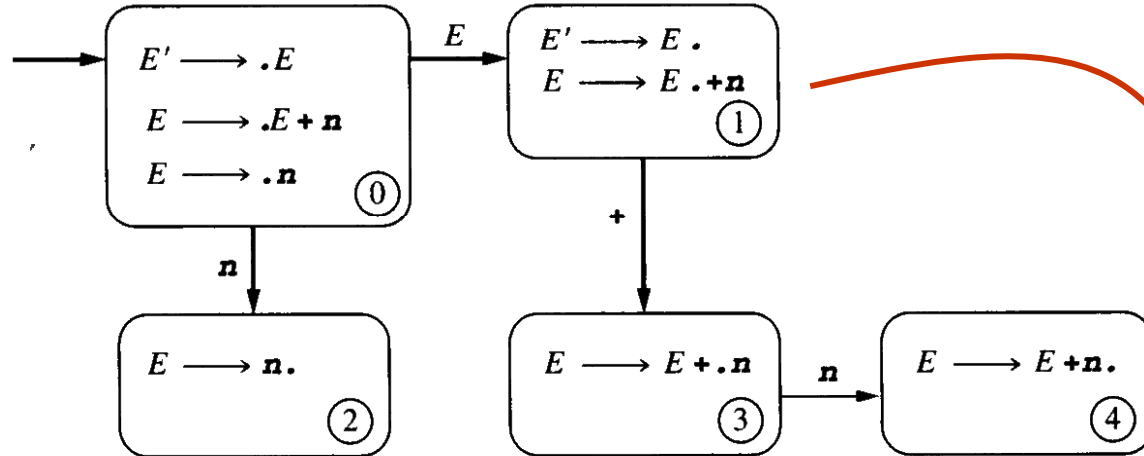
$S' \rightarrow S$
 $S \rightarrow (S) S \mid \varepsilon$

LR(0) – DFA:



Er den LR(0)? Nei, pga. tilstand 1: (har vi sett på tidligere)

$E' \rightarrow E$
 $E \rightarrow E + n \mid n$



\$	n + n \$
\$ n	+ n \$
\$ E	+ n \$
\$ E +	n \$
\$ E + n	\$
\$ E	\$
\$ E'	\$

Skal skifte

Skal redusere med $E' \rightarrow E$

Derved: Denne grammatikken er ikke LR(0): For stakk "E" er neste skritt ubestemt 28