



NTNU – Trondheim
Norwegian University of
Science and Technology

LR(0) parser construction

Bottom-up parsing

- Recapping the last lecture, what we need for a bottom-up parsing scheme is
 - An internal stack to shift and reduce symbols on
 - An automaton that tell us what to do when, and uses the stacked history to backtrack its footsteps
 - A grammar with one and only one initial production



The LR(0) automaton

- The overall construction is similar to the NFA \rightarrow DFA idea, in that
 - We're tracking all the different things that can happen throughout a derivation (there will be closures of related things)
 - We're introducing states whenever it is necessary to cover something new
 - The states represent all the different paths that may have led to them
 - Some states are *reducing*, that means
 - Pop body, push head
(but we won't draw the stack in, just remember that this happens)
 - Revert to where we started recognizing the present production
 - ...and the production which led to that one, if it is now finished...
 - ... and the production which led to that other one...
 - ... until the start symbol is all that's left
 - Transitions *shift* symbols
 - They are what moves us ahead while working toward a reduction



LR(0) *items*

- Since we have to track how far we have come along in a production's body, the notation for productions extends to a bunch of *LR(0) items*
 - All this means is that we add a marker ('.') to denote it
 - It's not punctuation in the language, just a position-tracking dot
 - The production
 $S \rightarrow aBc$
gives us the four items
 - $S \rightarrow . a B c$ (I'm just starting on an S)
 - $S \rightarrow a . B c$ (I'm working on S, already saw an a)
 - $S \rightarrow a B . c$ (I'm working on S, have found a and B)
 - $S \rightarrow a B c .$ (We have seen a whole S at this point, toss its history and put S)
 - We won't need every item all the time, but they are the things we will combine into states, so this is what they mean

(Augmenting the grammar with) A Place to Begin

- The point of rewriting the grammar

$S \rightarrow iCtSz \mid iCtSeSz$

into

$S' \rightarrow S$

$S \rightarrow iCtSz \mid iCtSeSz$

is that when we're looking at items,

$S' \rightarrow \cdot S$

uniquely says that we're just starting out, and

$S' \rightarrow S \cdot$

uniquely says that we're finished.

- It's just a convenience, to avoid the corner cases of
 - “Begin the construction with every item from the start symbol”, and
 - “Accept if any of the start symbol's productions reduce”
 – This way, the rule can just be “start at the beginning, stop at the end”, so to speak



(Let's keep the grammar handy, for reference) →

$$\begin{array}{l} D \rightarrow E F \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

Closures of items

- Closures of items are the sets of all items that can be obtained from the grammar without moving the position-marker
 - Consider this grammar
$$\begin{array}{l} D \rightarrow E F \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$
and start from the item $D \rightarrow \cdot E F$



$$\begin{array}{l} D \rightarrow E F \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

What can happen here?

- Without moving the marker, applying productions can give us

$$D \rightarrow \cdot E F$$

$$E \rightarrow \cdot y$$

$$E \rightarrow \cdot F$$

(those were the E-productions, we've made an F in the process)

$$F \rightarrow \cdot x$$

(this comes from repeating the procedure with the new items we found)

- What this closure represents, is that derivations from D can begin as

$$D \rightarrow \cdot E F$$

$$D \rightarrow \cdot E F \rightarrow \cdot y F$$

$$D \rightarrow \cdot E F \rightarrow \cdot F F \rightarrow \cdot x F$$

without moving the marker past any symbols.

- The notation implies all this (transitively)



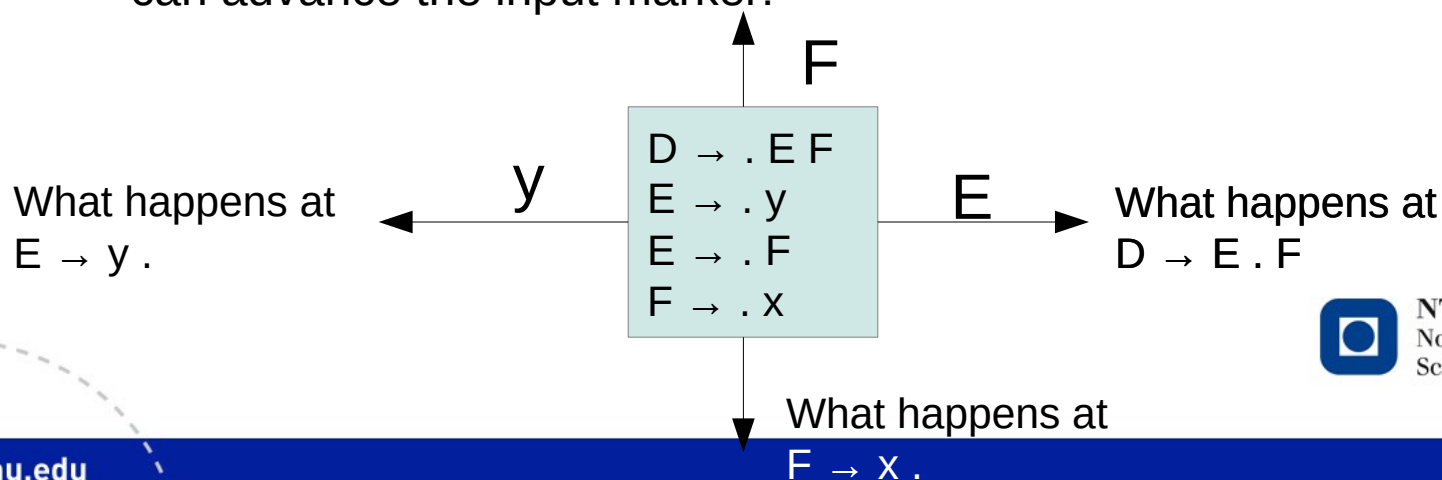
$$\begin{array}{l} D \rightarrow E F \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

Closures are states

- The closure gives us an automaton state, labeled with what it represents

$$\begin{array}{l} D \rightarrow \cdot E F \\ E \rightarrow \cdot y \\ E \rightarrow \cdot F \\ F \rightarrow \cdot x \end{array}$$

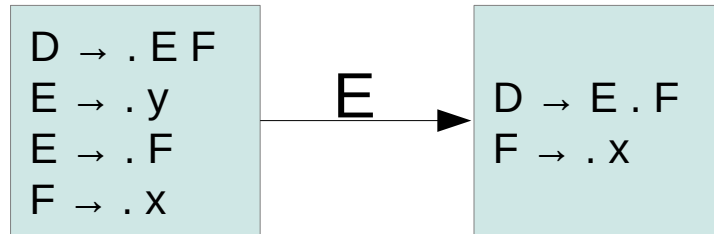
- The transitions out of this state represent all the different ways we can advance the input marker:



Advancing input (*Shift* actions)

$$\begin{array}{l} D \rightarrow EF \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

- *Selecting* one of these transitions gives us a new item to find the closure of, and hence, the destination state

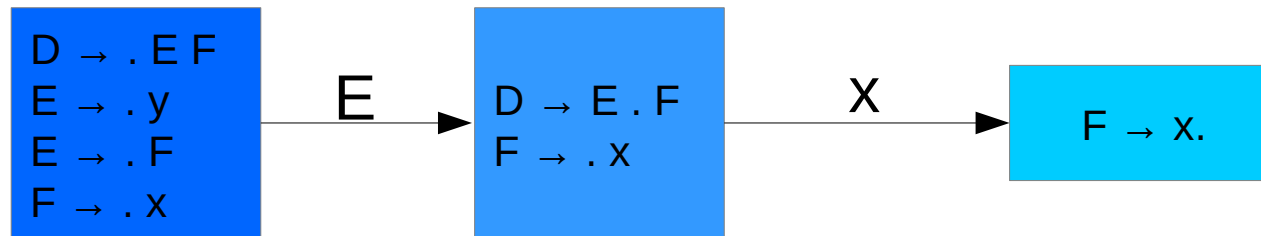


$$\begin{array}{l} D \rightarrow E F \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

Matching productions

(Reduce actions)

- When we reach an item that has the marker at the end, we've gone through states that encode a sequence which can appear in a derivation

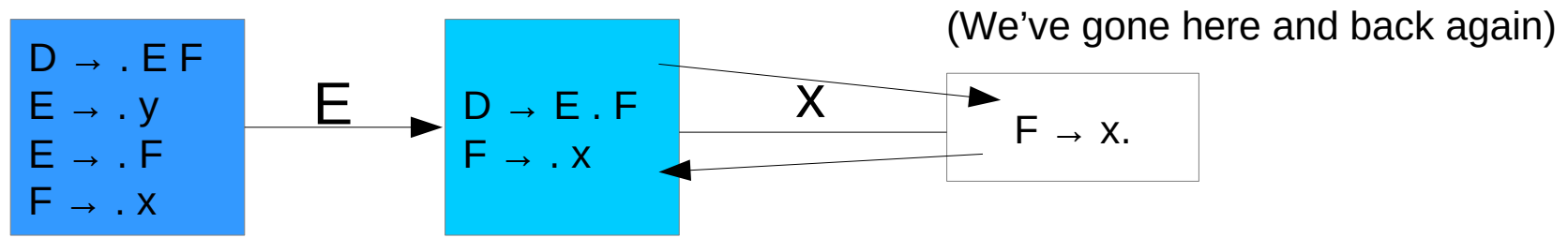


- $D \rightarrow \cdot E F$ (is) $E \cdot F \rightarrow E \cdot X$ (is) $E X \cdot$
 - As far as the grammar is concerned, $D \rightarrow E F \rightarrow E X$
 - We just decorated it with a position in our stack-based reasoning

D	→	E F
E	→	y F
F	→	x

This is where we backtrack

- Reducing at state $F \rightarrow x \cdot$ means
 - we have a stack with an x on top of it
 - Remove it, and replace it with an F
- That returns us to the stage where we were about to shift x

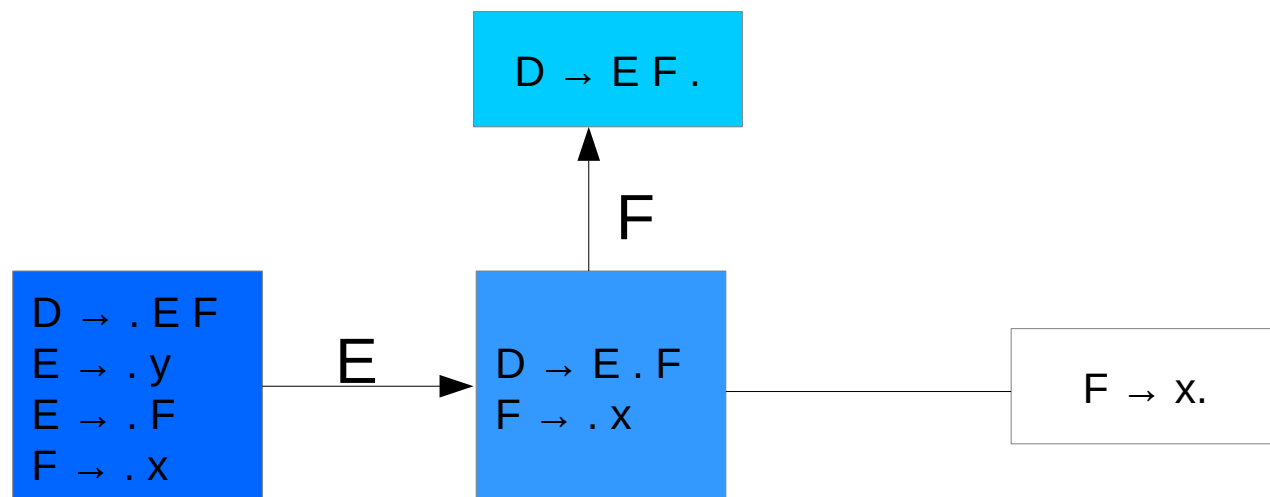


- Next thing on stack is the F we just created a moment ago

$$\begin{array}{l} D \rightarrow EF \\ E \rightarrow y | F \\ F \rightarrow x \end{array}$$

...and again

- Below the F is the E that brought us here...

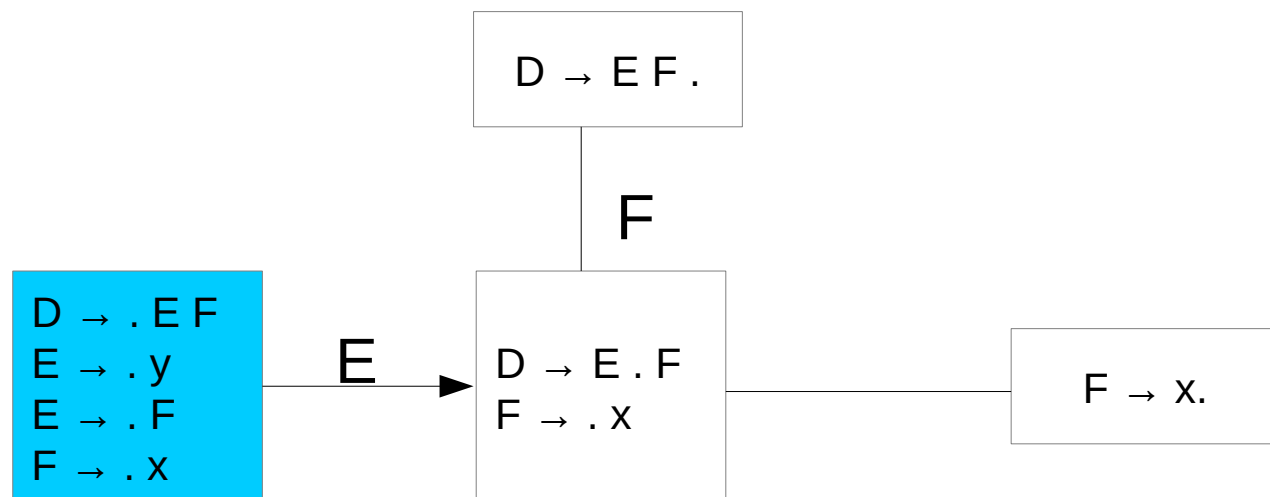


- $D \rightarrow EF \cdot$ is another reducing state, take the E and F off stack, put in a D instead, and return to before the E was shifted...

$$\begin{array}{l} D \rightarrow EF \\ E \rightarrow y \mid F \\ F \rightarrow x \end{array}$$

That was one traversal

- We've been through $D \rightarrow EF \rightarrow Ex$



- We started by shifting an E
 - that one must have been produced by a similar traversal elsewhere in the automaton, which shifted y , reduced $E \rightarrow y$, and thus gave us an E to shift

Making the LR(0) automaton

Start with the designated start item

– The one that looks like $X' \rightarrow \cdot X$

- 1) Find its closure, make a state
- 2) Follow all the transitions
- 3) Repeat from 1

until you reach the reduction $X' \rightarrow X$ at the other end.

(and don't duplicate states when you get an item you've already made a state from)



(Going by hand)

- If you don't respect some strict depth-first traversal ordering, the final reduction can be the first thing you find
 - It gets hard to remember where you were after a few branches and backtracks
- In this case, it is necessary to stare at the automaton for a bit, to convince yourself that you've visited it everywhere
 - That's not so mathematically rigorous, but it's OK for what we're after
- We're doing this to understand how it works
 - Homespun LR(0) parsing is a waste of time, there are splendid generator programs

A simple grammar to try it on

This one models nested*, comma-separated lists:

$$S \rightarrow (L) \mid x$$

$$L \rightarrow S \mid L, S$$

i.e. statements like

$$(x,x,(x,x))$$

$$S \rightarrow (L) \rightarrow (L,S) \rightarrow (L,S,S) \rightarrow (S,S,S) \rightarrow^* (x,x,S) \rightarrow (x,x,(L)) \rightarrow (x,x,(L,S)) \\ \rightarrow^* (x,x,(x,x))$$

or

$$(x,(x,x),x,(x,x))$$

and similar

**(Confer w. regex versus nested parentheses - Context-Free Grammars are more powerful...)*



Blackboard time

- There's a summary of the whole development on the next few slides, but I think it's worth going through at a more leisurely pace, so I'll draw it step by step.
- Follow as it fits yourself, you can review the slides

From the beginning

- Augmenting the grammar with a production that has a single rule, we get somewhere to start:

$$S' \rightarrow S$$

- Taking that as the basis for a first state, its closure is

$$S' \rightarrow .S$$

$$S \rightarrow .(L)$$

$$S \rightarrow .X$$

so we make an automaton state out of that

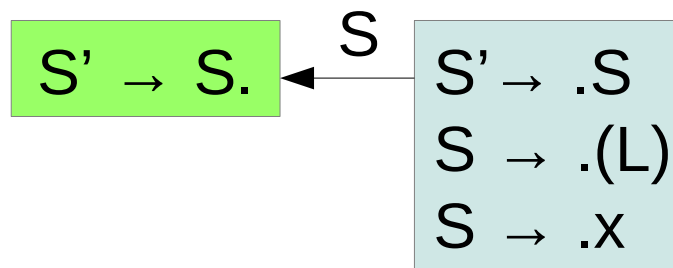
$$\begin{array}{l} S' \rightarrow .S \\ S \rightarrow .(L) \\ S \rightarrow .X \end{array}$$



$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow (L) \\ S &\rightarrow X \\ L &\rightarrow S \\ L &\rightarrow L, S \end{aligned}$$

From the end

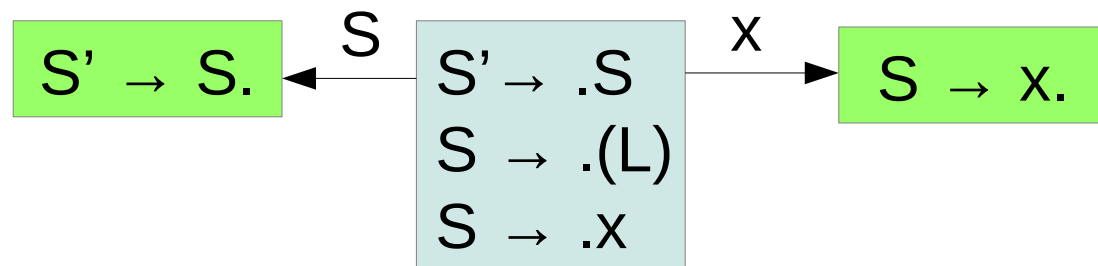
- The last thing that will happen is that we have parsed an S, and can shift it to complete parsing
- That sends us to a state where we can reduce our artificial start symbol, and declare victory



$$\begin{array}{l} S' \rightarrow S \\ S \rightarrow (L) \\ S \rightarrow X \\ L \rightarrow S \\ L \rightarrow L, S \end{array}$$

Another thing can happen

- We might shift an 'x' terminal

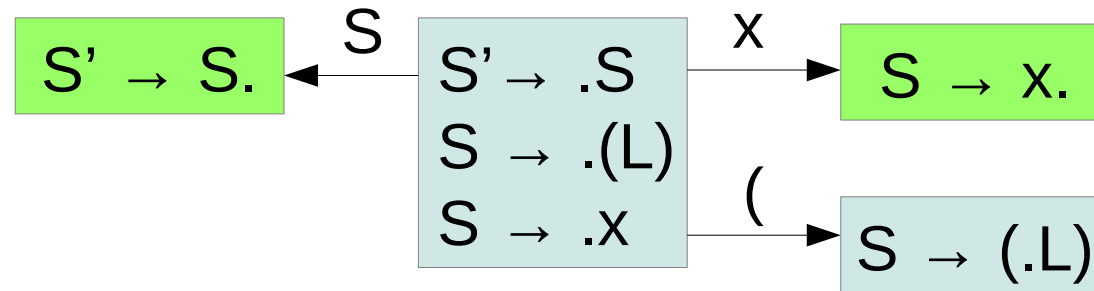


- This completes another production (. is at the end), so it is also a state where we have found a reduction

$$\begin{array}{l}
 S' \rightarrow S \\
 S \rightarrow (L) \\
 S \rightarrow X \\
 L \rightarrow S \\
 L \rightarrow L, S
 \end{array}$$

The final thing that can happen

- We might shift a '(' terminal

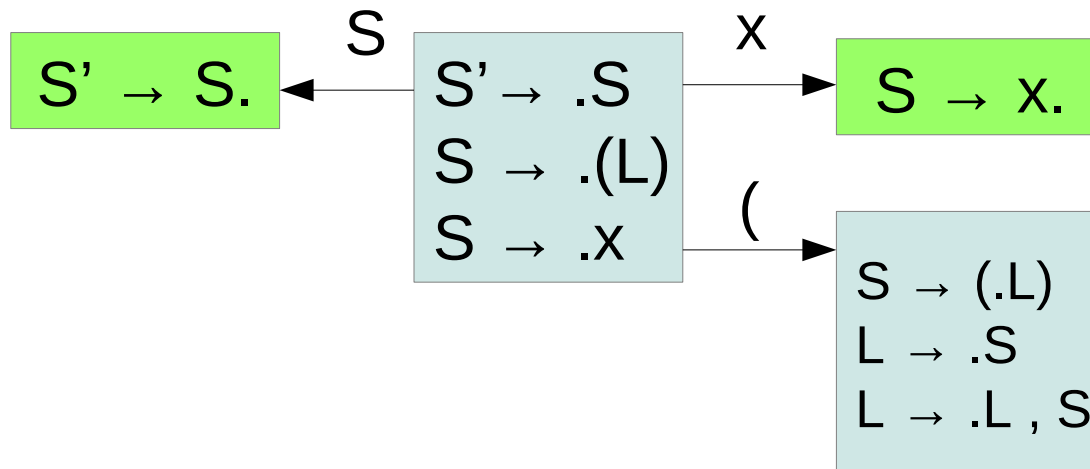


- This doesn't complete any productions, so we'll have to build more states
- Start over, with the closure at the destination state

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow (L) \\ S &\rightarrow X \\ L &\rightarrow S \\ L &\rightarrow L, S \end{aligned}$$

The closure at $S \rightarrow (.L)$

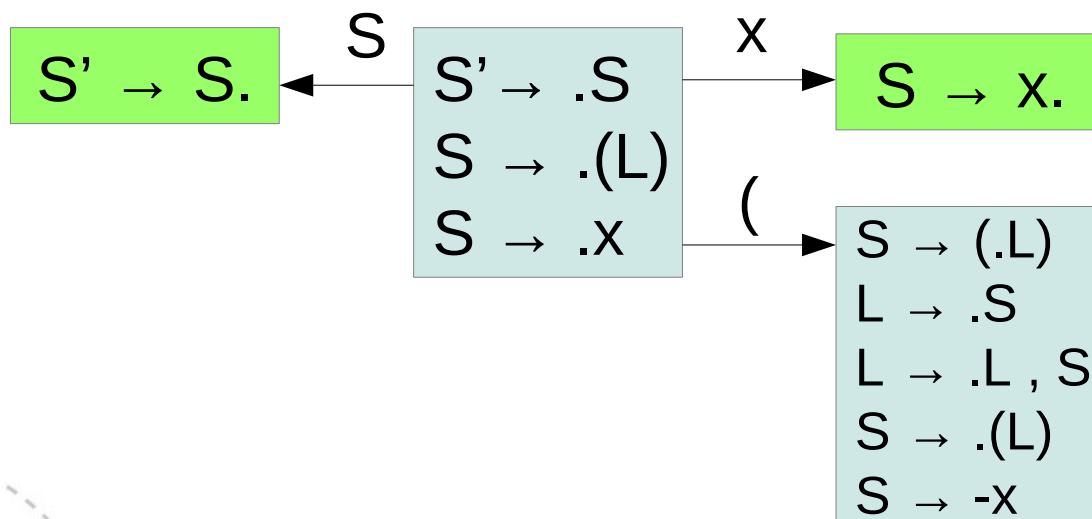
- A non-terminal follows the position marker
 - That can expand into
 - $L \rightarrow S$
 - $L \rightarrow L, S$



$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow (L) \\ S &\rightarrow x \\ L &\rightarrow S \\ L &\rightarrow L, S \end{aligned}$$

We're not done yet

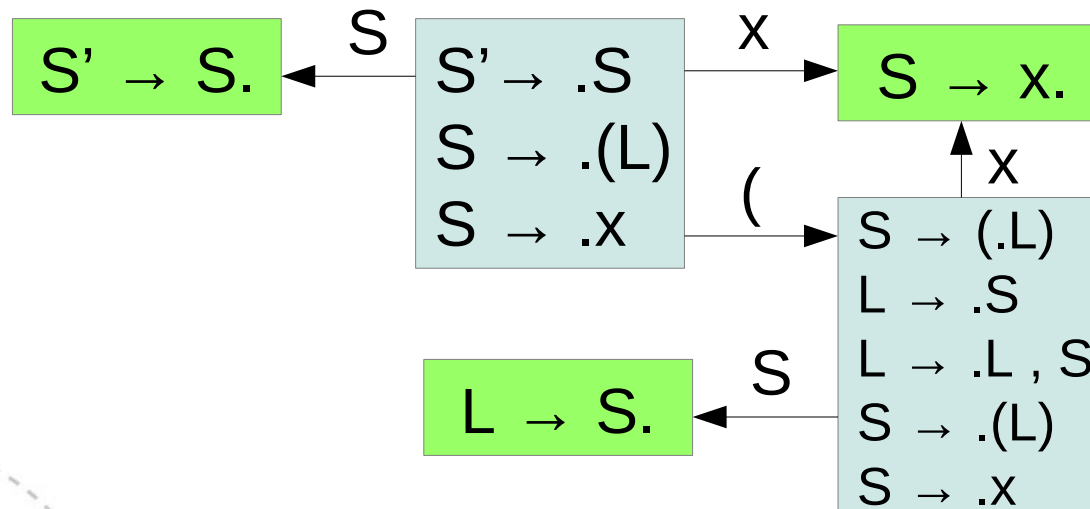
- The expansion put the position marker ahead of another nonterminal (that is, S)
 - S can expand into
 - $S \rightarrow (L)$
 - $S \rightarrow x$



$$\begin{aligned}
 S' &\rightarrow S \\
 S &\rightarrow (L) \\
 S &\rightarrow x \\
 L &\rightarrow S \\
 L &\rightarrow L, S
 \end{aligned}$$

Time to shift more symbols

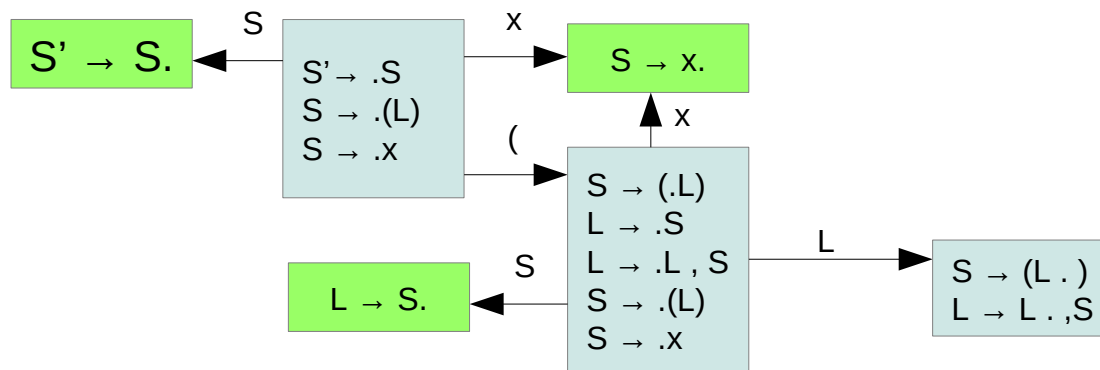
- Reducing states are easy to find, shifting S or x completes a production
 - We already have a state for $S \rightarrow x$.
 - $L \rightarrow S$ is the other reducing state we can reach in one shift



$$\begin{aligned}
 S' &\rightarrow S \\
 S &\rightarrow (L) \\
 S &\rightarrow x \\
 L &\rightarrow S \\
 L &\rightarrow L, S
 \end{aligned}$$

More states, same work

- The remaining two items both suggest shifting an L
- That's one transition, but both items can be a reason to take it

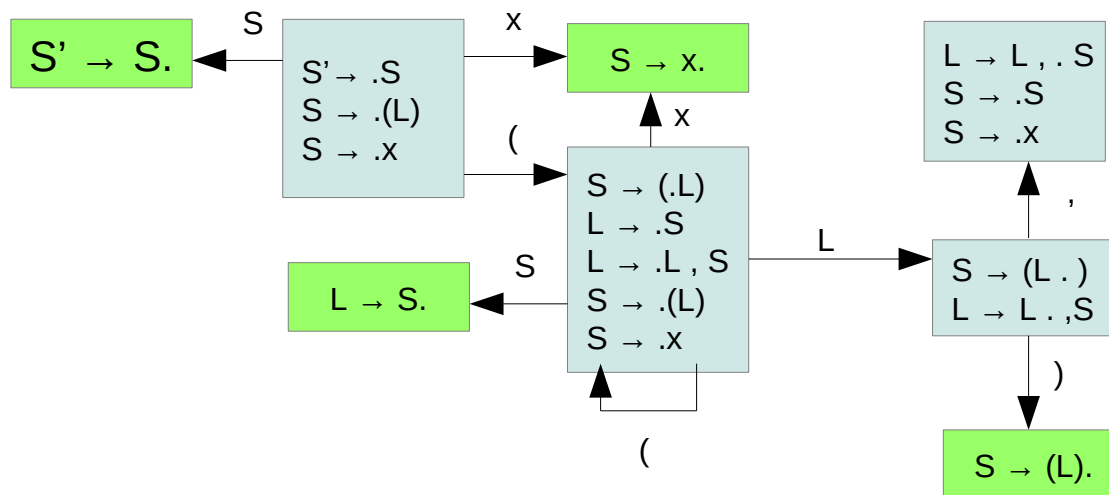


- The position marker doesn't precede any nonterminals, so this is all the closure we need

$$\begin{aligned} S' &\rightarrow S \\ S &\rightarrow (L) \\ S &\rightarrow x \\ L &\rightarrow S \\ L &\rightarrow L, S \end{aligned}$$

Two alternatives

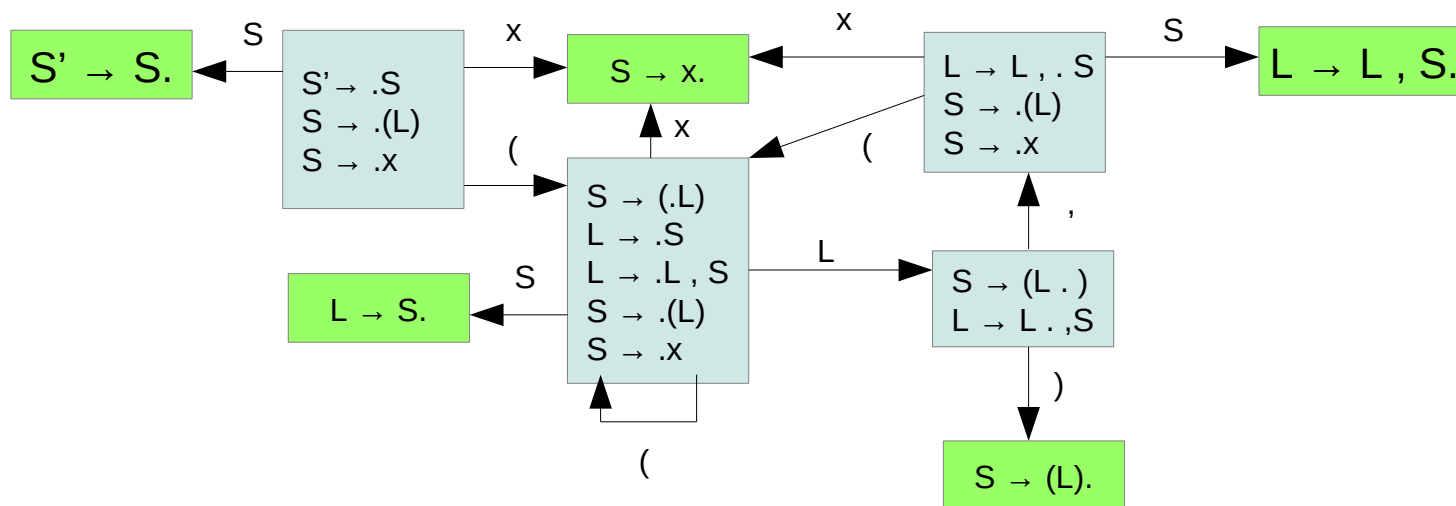
- Shift a ')' terminal, or a ',' terminal
 - The ')' leads to a reducing state
 - ',' leads to an item we haven't treated yet, so make a state, and find the closure it represents



$$\begin{aligned}
 S' &\rightarrow S \\
 S &\rightarrow (L) \\
 S &\rightarrow x \\
 L &\rightarrow S \\
 L &\rightarrow L, S
 \end{aligned}$$

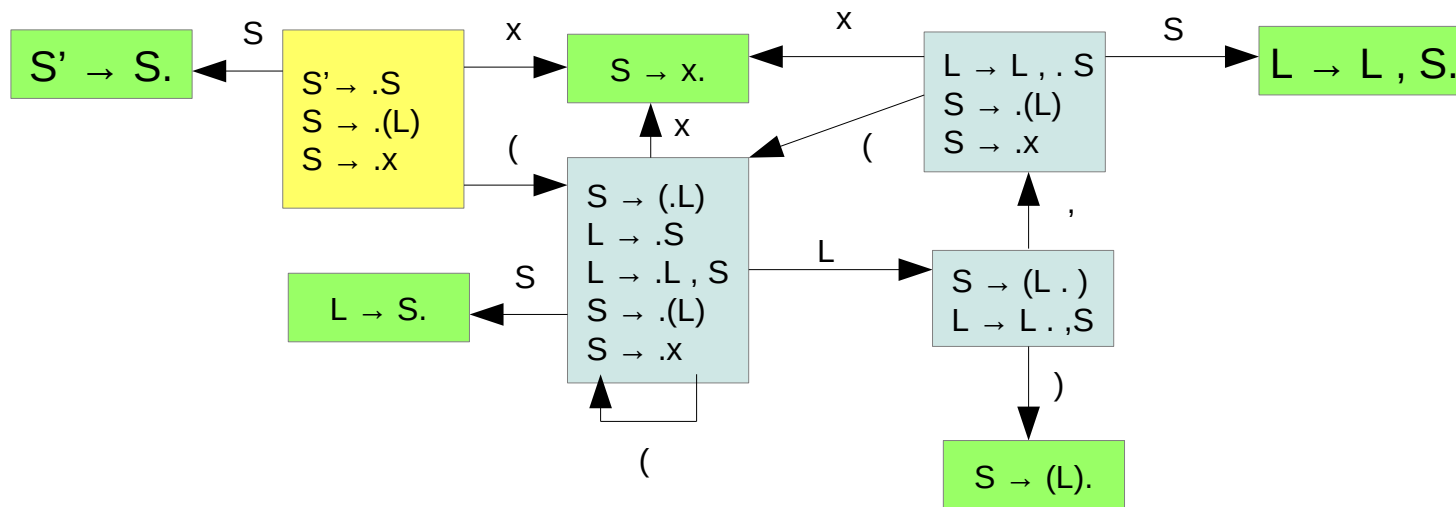
There's one more state

- The only new thing here is that we can finalize parsing of the production $L \rightarrow L, S$
 - 'x' and '(' lead to states we have already created



That's our LR(0) automaton

- There is a bit of window-dressing to transform it into a table, we can look at that next time



Just a final footnote

- As you may have noticed
 - Wrt. the Kleene closure of regex (r^*), we saw that it's an infinite set
 - The epsilon closures have the added constraint that no character should be matched, so they become finite sets of states
 - Today's closures of items have the same constraint, for much the same reason
- There is a pattern here
 - The way we find FOLLOW sets in top-down parsing is the exact same principle at work too, I just didn't say it out loud at the time



The General Approach

(in a very pseudo-code way)

- 1) Initialize some number of sets
 - 2) Update them so that they satisfy all constraints
 - 3) Record whether any of them changed because of step 2
 - 4) If any did, repeat from step 2
 - 5) If none did, declare victory
- This is “*iteration to a fixed point*” (victory is the “fixed point”)
 - Calling it by that name foreshadows something deeper, every program can be rewritten as a constraint problem
 - We’re moving outside compiler construction here, so never mind
 - I mention it still
 - ...because recognizing this pattern on sight might make it easier to remember all our different variants.

