# LR Bottom-up Parsing

# Roadmap

Last class

– Name analysis

Previous-ish last class

– LL(1)

Today's class

– LR Parsing

• SLR(1)

# Lecture Outline

Bottom-Up parsing
- Talk about the language class / theory
- Describe the state that it keeps / intuition
- Show how it works
- Show how it is built

# LL(1) Not Powerful Enough for all PL

Left-recursion

Not left factored

Doesn't mean LL(1) is bad
- Right tool for simple parsing jobs



```
stmtList  ::= stmtList stmt
          |  /* epsilon */
          ;
```

# We Need a *Little* More Power

Could increase the lookahead

– Up until the mid 90s, this was considered impractical

Could increase the runtime complexity

– CYK has us covered there

Could increase the memory complexity

– i.e. more elaborate parse table

# LR Parsers

Left-to-right scan of the input file

Reverse rightmost derivation

Advantages
- Can recognize almost any programming language
- Time and space $O(n)$ in the input size
- More powerful than the corresponding LL parser i.e. LL(1) < LR(1)

Disadvantages
- More complex parser generation
- Larger parse tables

# LR Parser Power

Let $S \Longrightarrow \alpha_1 \Longrightarrow \alpha_2 \Longrightarrow \ldots \Longrightarrow w$ be a rightmost derivation, where $\omega$ is a terminal string

Let $\alpha A\gamma \Longrightarrow \alpha\beta\gamma$ be a step in the derivation

- So $A \longrightarrow \beta$ must have been a production in the grammar
- $\alpha\beta\gamma$ must be some $\alpha_i$ or w

- A grammar is LR(k) if for every derivation step, $A \longrightarrow \beta$ can be inferred using only a scan of $\alpha\beta$ and at most k symbols of $\gamma$

Much like LL(1), you generally just have to go ahead and try it
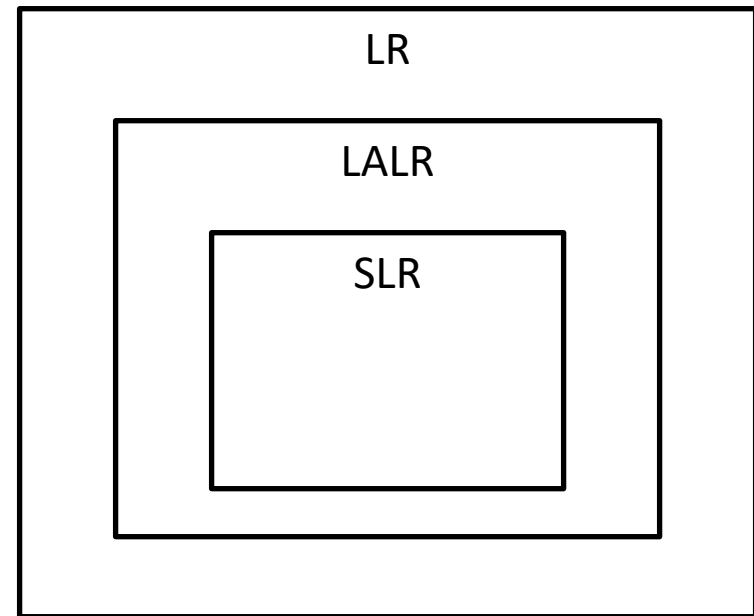
# LR Parser types

LR(1)

– Can recognize any DCFG

– Can experience blowup in parse table size

LALR(1)

SLR(1)

– Both proposed at the same time to limit parse table size

**Recognizable by a deterministic PDA**

LR

LALR

SLR

# Which parser should we use?

Different variants mostly differ in how they build the parse table, we can still talk about all the family in general terms

- Today we'll cover SLR
- Pretty easy to learn LALR from there

LALR(1)

- Generally considered a good compromise between parse table size and expressiveness
- Class for Java CUP, yacc, and bison

# How does Bottom-up Parsing work?

Already seen 1 such parser: CYK
- Simultaneously tracked every possible parse tree
- LR parsers work in a similar same way

Contrast to top-down parser
- We know exactly where we are in the parse
- Make predictions about what's next

# Parser State

## Top-down parser state
- Current token
- Stack of symbols
  - Represented what we expect in the rest of our descent to the leaves
- Worked down and to the left through tree

## Bottom-up parser state
- Also maintains a stack and token
  - Represents summary of input we've seen
- Works upward and to the right through the tree
- Also has an auxiliary state machine to help disambiguate rules

Grammar

$S$ ::= **ε**
  | **(** $S$ **)**
  | **[** $S$ **]**

Stack

| [ |
| S |
| ] |
| ) |
| EOF |

Current

↓

[

# LR Derivation Order

Let's remember derivation orders again

| Reverse | | Rightmost derivation |
|---------|---|---|
| 8 | 1 | $E \Longrightarrow E + T$ |
| 7 | 2 | $\Longrightarrow E + T * F$ |
| 6 | 3 | $\Longrightarrow E + T * id$ |
| 5 | 4 | $\Longrightarrow E + F * id$ |
| 4 | 5 | $\Longrightarrow E + id * id$ |
| 3 | 6 | $\Longrightarrow T + id * id$ |
| 2 | 7 | $\Longrightarrow F + id * id$ |
| 1 | 8 | $\Longrightarrow id + id * id$ |

# Parser Operations
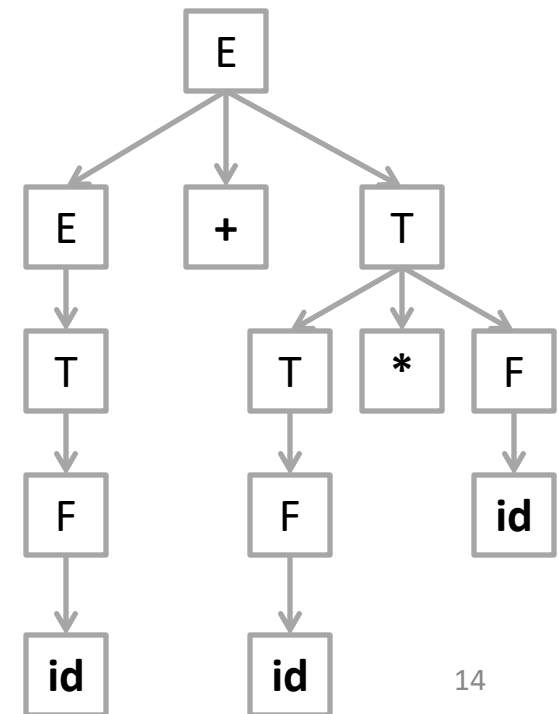
Top-down parser

- *Scan* the next input token
- *Push* a bunch of RHS symbols
- *Pop* a single symbol

Bottom-up parser

- *Shift* an input token into a stack item
- *Reduce* a bunch of stack items into a new parent item (on the stack)

# Parser Actions: Simplified view

| Stack | Input | Action |
|---|---|---|
| | id + id * id EOF | shift(id) |
| id | + id * id EOF | reduce by F $\rightarrow$ id |
| F | + id * id EOF | reduce by T $\rightarrow$ F |
| T | + id * id EOF | reduce by E $\rightarrow$ T |
| E | + id * id EOF | shift + |
| E + | id * id EOF | shift id |
| E + id | * id EOF | reduce by F $\rightarrow$ id |
| E + F | * id EOF | reduce by T $\rightarrow$ F |
| E + T | * id EOF | shift * |
| E + T * | id EOF | shift id |
| E + T * id | EOF | reduce by F $\rightarrow$ id |
| E + T * F | EOF | reduce by T $\rightarrow$ T * F |
| E + T | EOF | reduce by E $\rightarrow$ E + T |
| E | EOF | accept |

# Stack Items

Note that the previous slide was called "simplified"

Stack elements are representative of symbols
- Actually known as items
  - Indicate a production and a position within the production

$$X \longrightarrow \alpha \,.\, B\, \beta$$

  - Means
    - we are in a production of X
    - We believe we've parsed (arbitrary) symbol string $\alpha$
    - We could handle a production of B
    - After that we'll have $\beta$

# Stack Item Examples

Example 1

    *PList* → ( . IDList )

Example 2

    *PList* → ( IDList . )

Example 3

    *PList* → ( IDList ) .

Example 4

    *PList* → . ( IDList )

# Stack Item State

You may not know exactly which item you are parsing

LR Parsers actually track the set of states that you *could* have been in

**Grammar snippet**
S → A
A → B
  | C
B → D **id**
C → **id** E
D → **id** E

{S → . A, A → . B, A → . C, ...}

# LR Parser FSM

**Grammar G**
S' → *PList*
*PList* → ( *IDList* )
*IDList* → **id**
*IDList* → *IDList* **id**

I₀
S' → . *PList*
*PList* → . **(** IDList **)**

*PList*

I₁
S' → *PList* .

**(**

I₂
*PList* → **(** . IDList **)**
*IDList* → . **id**
*IDList* → . IDList **id**

I₅
*PList* → **(** IDList **)** .

*IDList*

I₃
*PList'* → **(** IDList . **)**
*IDList* → IDList . **id**

**)**

**id**

I₄
*IDList* → **id** .

**id**

I₆
*IDList* → IDList **Id** .

18

# Automaton as a table



- *Shift* corresponds to taking a terminal edge
- *Reduce* corresponds to taking a nonterminal edge

**Action table**   **GoTo table**

|   | ( | ) | id | eof | *PList* | *IDList* |
|---|---|---|----|-----|---------|----------|
| 0 | S 2 |   |    |     | 1 |   |
| 1 |   |   |    |     |   |   |
| 2 |   |   | S 4 |    |   | 3 |
| 3 |   | S 5 | S 6 |   |   |   |
| 4 |   |   |    |     |   |   |
| 5 |   |   |    |     |   |   |
| 6 |   |   |    |     |   |   |

Shift and go to state 6

# How do we know when to reduce?

**Action table**

**GoTo table**

| | ( | ) | id | eof | PList | IDList |
|---|---|---|---|---|---|---|
| 0 | S 2 | | | | 1 | |
| 1 | | | | | | |
| 2 | | | S 4 | | | 3 |
| 3 | | S 5 | S 6 | | | |
| 4 | | R ❸ | R ❸ | | | |
| 5 | | | | R ❷ | | |
| 6 | | R ❹ | R ❹ | | | |

**Grammar G**
❶ *S' → PList*
❷ *PList → ( IDList )*
❸ *IDList → id*
❹ *IDList → IDList id*

Only see terminals in the input

Actually do reduce steps in 2 phases

– Action table will tell us when to reduce (and how much)

– GoTo will tell us where to… go to

# How do we know we're done?

**Action table**

| | ( | ) | id | eof |
|---|---|---|---|---|
| 0 | S 2 | | | |
| 1 | | | | ☺ |
| 2 | | | S 4 | |
| 3 | | S 5 | S 6 | |
| 4 | | R ❸ | R ❸ | |
| 5 | | | | R ❷ |
| 6 | | R ❹ | R ❹ | |

**GoTo table**

| PList | IDList |
|---|---|
| 1 | |
| | |
| | 3 |
| | |
| | |
| | |
| | |

Add an accept token

Any other cell is an error

**Grammar G**
❶ S' → PList
❷ PList → ( IDList )
❸ IDList → **id**
❹ IDList → IDList **id**

# Full Parse Table Operation

```
Initialize stack
a = scan()
do forever
    t = top-of-stack (state) symbol
    switch action[t, a] {
       case shift s:
          push(s)
          a = scan()
       case reduce by A → alpha:
          for i = 1 to length(alpha) do pop() end
        t = top-of-stack symbol
          push(goto[t, A])
       case accept:
          return( SUCCESS )
       case error:
          call the error handler
          return( FAILURE )
    }
end do
```

# Example Time

**I₀**
$S' \rightarrow . PList$
$PList \rightarrow . ( IDList )$

**PList** →

**I₁**
$S' \rightarrow PList .$

**(**

**I₂**
$PList \rightarrow ( . IDList )$
$IDList \rightarrow . id$
$IDList \rightarrow . IDList id$

**IDList**

**I₅**
$PList \rightarrow ( IDList ) .$

**I₃**
$PList \rightarrow ( IDList . )$
$IDList \rightarrow IDList . id$

**)**

**id**

**I₄**
$IDList \rightarrow id .$

**id**

**I₆**
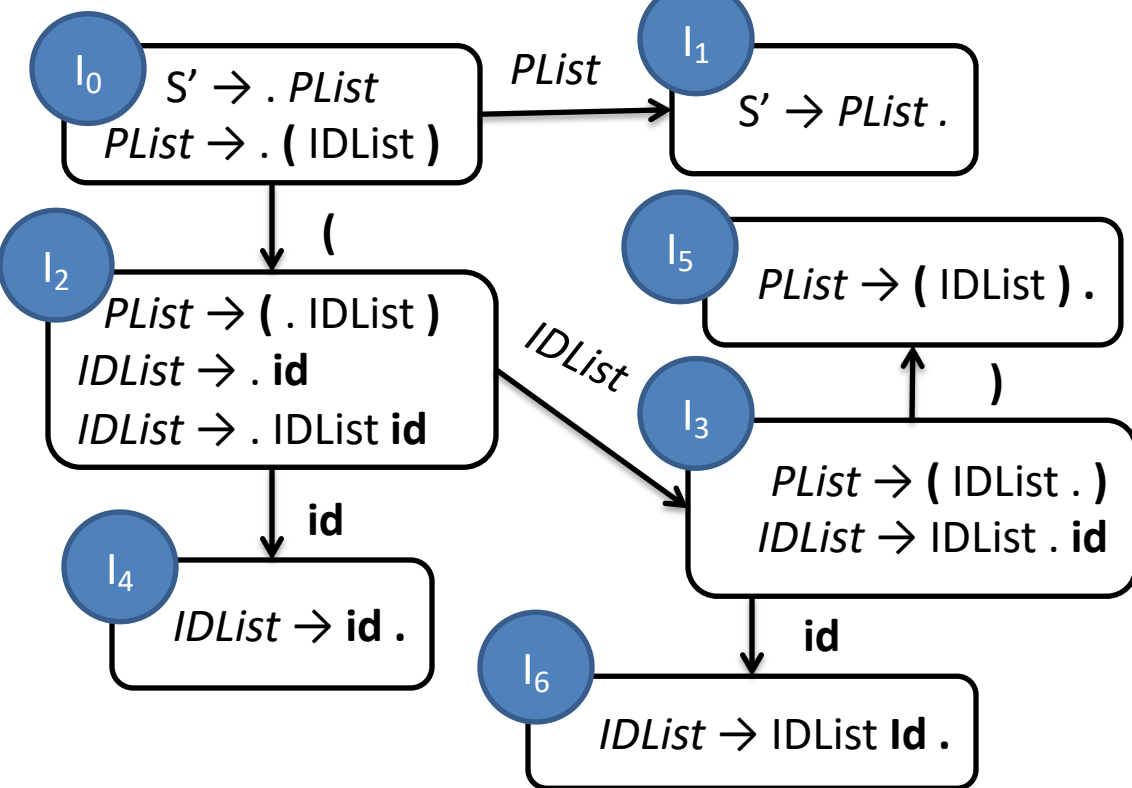$IDList \rightarrow IDList Id .$

current element

**( id id id ) eof**

**Grammar G**
❶ $S' \rightarrow PList$
❷ $PList \rightarrow ( IDList )$
❸ $IDList \rightarrow id$
❹ $IDList \rightarrow IDList id$

[I₅]
[I₃]
[I₁]
[I₀]

| | ( | ) | id | eof | *PList* | *IDList* |
|---|---|---|---|---|---|---|
| 0 | S 2 | | | | 1 | |
| 1 | | | | ☺ | | |
| 2 | | | S 4 | | | 3 |
| 3 | | S 5 | S 6 | | | |
| 4 | | R ❸ | R ❸ | | | |
| 5 | | | | R ❷ | | |
| 6 | | R ❹ | R ❹ | | | |

# Seems that LR Parser works great
# What could possible go wrong?

# LR Parser State Explosion

Tracking sets of states can cause the size of the FSM to blow up

The SLR and LALR variants exist to combat this explosion

Slight modification to item and table form

# Building the SLR Automaton

Uses 2 sets
- Closure(I)
  - What is the set of items we could be in?
  - Given I: what is the set of items that could be mistaken for I (reflexive)
- Goto(s,X)
  - If we are in state I, where might we be after parsing X?

Vaguely reminiscent of FIRST and FOLLOW

# Closure Sets

Put I itself into Closure(I)


While there exists an item in Closure(I) of form

      X $\longrightarrow$ α . B β

              such that there is a production B $\longrightarrow$ γ


      and B $\longrightarrow$ . γ is not in Closure(I)

add B $\longrightarrow$ . γ to Closure(I)

# GoTo Sets

Goto(I, X) =

Closure({ A $\longrightarrow$ $\alpha$ X . B | A $\longrightarrow$ $\alpha$ . X $\beta$ is in I })

States in FSM:

- $I_0$: $S' \to .\ PList$ ; $PList \to .\ ($ IDList $)$
- $I_1$: $S' \to PList\ .$
- $I_2$: $PList \to ($ . IDList $)$ ; $IDList \to .\ \mathbf{id}$ ; $IDList \to .$ IDList $\mathbf{id}$
- $I_3$: $PList' \to ($ IDList . $)$ ; $IDList \to$ IDList . $\mathbf{id}$
- $I_4$: $IDList \to \mathbf{id}$ .
- $I_5$: $PList \to ($ IDList $)$ .
- $I_6$: $IDList \to$ IDList $\mathbf{Id}$ .

Transitions: $I_0 \xrightarrow{PList} I_1$ ; $I_0 \xrightarrow{(} I_2$ ; $I_2 \xrightarrow{IDList} I_3$ ; $I_2 \xrightarrow{\mathbf{id}} I_4$ ; $I_3 \xrightarrow{)} I_5$ ; $I_3 \xrightarrow{\mathbf{id}} I_6$

**Grammar G**

$S' \to PList$

$PList \to ($ IDList $)$

$IDList \to \mathbf{id}$

$IDList \to IDList\ \mathbf{id}$

**GoTo(I, X)**

Put items Closure(I) items

Repeat for X.β s.t. A → α.Xβ ∈ I

   X → α.Bβ ∈ Closure(I) s.t.

    ∃B →γ, add B →.γ to Closure(I)

GoTo(I₀, PList)

all Items A → α PList .β)

  all Items A → α IDList .β

[1] PList S → . PList )

  [1] PList → ( IDList . )

those where A → α PList β ∈ I₀)

  [2] IDList Closure . id . PList} = {

for [1] PList S → ( IDList is in I₀

  those where A → α IDList β ∈ I₂

set to closure of the following):

for [1] PList → ( . IDList ) is in I₂

{ PList S → ( . PList )}

for [2] PList PList IDList IDList . ) in I₂

Items add nothing where IDList → γ ∈ G

set to closure is

{ IDList → . id PList .}

{ PList → ( IDList id . ), IDList → ( IDList . ) }

IDList → . IDList id }

Only terminals after . so closure done

Done with closure, and GoTo

**Parse Table Construction**

1: Add new start S' and $S' \to S$
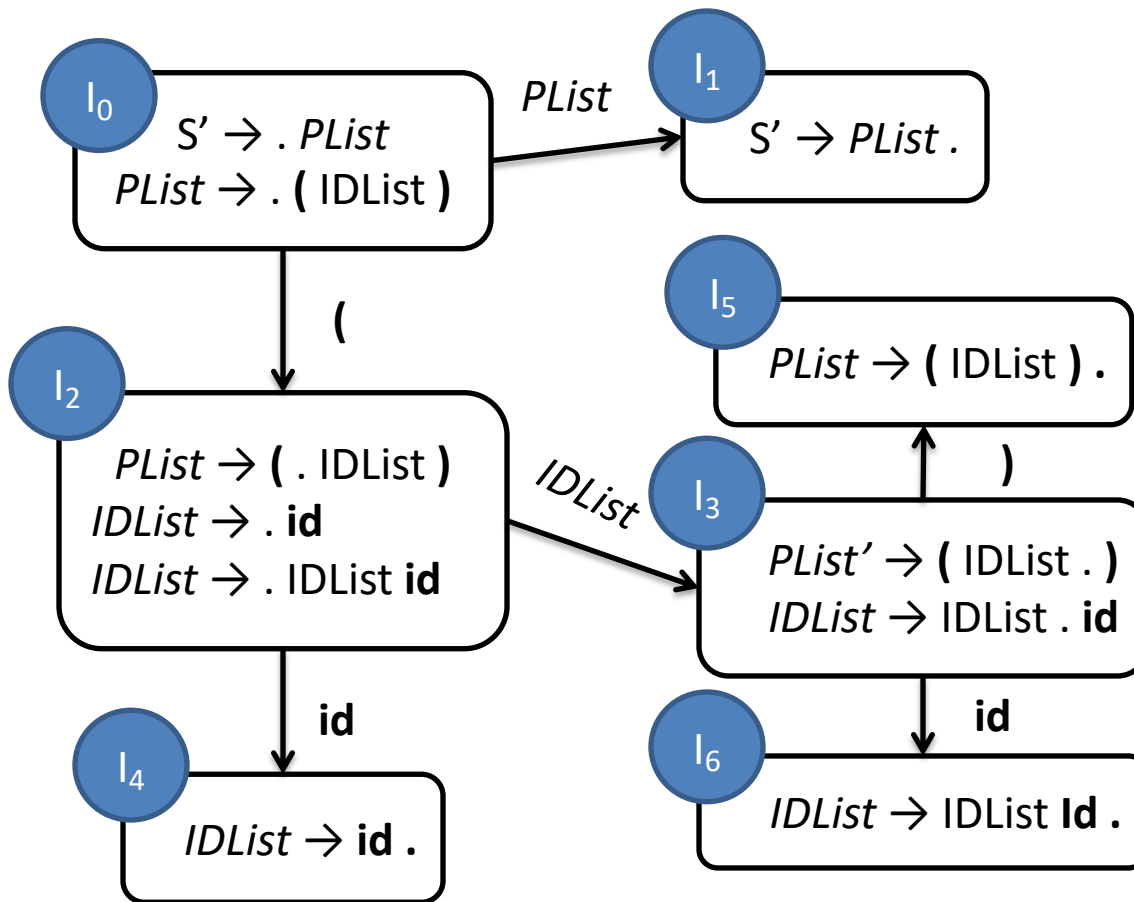
2: Build State $I_0$ for Closure( $\{ S' \to .\ S \}$ )

3: Saturate FSM:

  for each symbol X s.t. there is a item in state j containing . X

    add transition from state j to state for GoTo(j, X)

# From FSM to parse table(s)



Need to connect the FSM back to the grammar

**Grammar G**
❶ S' → *PList*
❷ *PList* → ( *IDList* )
❸ *IDList* → **id**
❹ *IDList* → *IDList* **id**

# Can Now Build Action and GoTo Tables



$I_0$
$S' \rightarrow .\ PList$
$PList \rightarrow .\ ( \text{IDList} )$

PList

$I_1$
$S' \rightarrow PList\ .$

(

$I_2$
$PList \rightarrow ( .\ \text{IDList} )$
$IDList \rightarrow .\ \textbf{id}$
$IDList \rightarrow .\ \text{IDList}\ \textbf{id}$

IDList

$I_3$
$PList' \rightarrow ( \text{IDList} .\ )$
$IDList \rightarrow \text{IDList} .\ \textbf{id}$

$I_5$
$PList \rightarrow ( \text{IDList} )\ .$

)

id

$I_4$
$IDList \rightarrow \textbf{id}\ .$

id

$I_6$
$IDList \rightarrow \text{IDList}\ \textbf{Id}\ .$

# Building the GoTo Table



$I_0$
$S' \rightarrow . \; PList$
$PList \rightarrow . \; ( \; IDList \; )$

$I_1$
$S' \rightarrow PList \; .$

PList

(

$I_2$
$PList \rightarrow ( \; . \; IDList \; )$
$IDList \rightarrow . \; id$
$IDList \rightarrow . \; IDList \; id$

IDList

$I_5$
$PList \rightarrow ( \; IDList \; ) \; .$

$I_3$
$PList' \rightarrow ( \; IDList \; . \; )$
$IDList \rightarrow IDList \; . \; id$

)

id

$I_4$
$IDList \rightarrow id \; .$

$I_6$
$IDList \rightarrow IDList \; Id \; .$

id

For every nonterminal $X$
　if there is an (i,j) edge on $X$
　　set GoTo[i,$X$] = j

|   | PList | IDList |
|---|-------|--------|
| 0 | 1     |        |
| 1 |       |        |
| 2 |       | 3      |
| 3 |       |        |
| 4 |       |        |
| 5 |       |        |
| 6 |       |        |

33

# Building the Action Table

If state i includes item A → α . **t** β

- where **t** is a terminal
- and there is an (i,j) transition on **t**
- set Action[i,**t**] = shift j

If state i includes item A → α .

- where A is not S'
- for each t in FOLLOW(A):
- set Action[i,**t**] = reduce by A → α

If state i includes item S → S .

- set Action[i, **eof**] = accept

All other entries are error actions

# Action Table: Shift



$I_0$: $S' \rightarrow$ . $PList$
$PList \rightarrow$ . ( IDList )

$I_1$: $S' \rightarrow PList$ .

$I_2$: $PList \rightarrow$ ( . IDList )
$IDList \rightarrow$ . id
$IDList \rightarrow$ . IDList id

$I_4$: $IDList \rightarrow$ id .

$I_5$: $PList \rightarrow$ ( IDList ) .

$I_3$: $PList \rightarrow$ ( IDList . )
$IDList \rightarrow$ IDList . id

$I_6$: $IDList \rightarrow$ IDList Id .
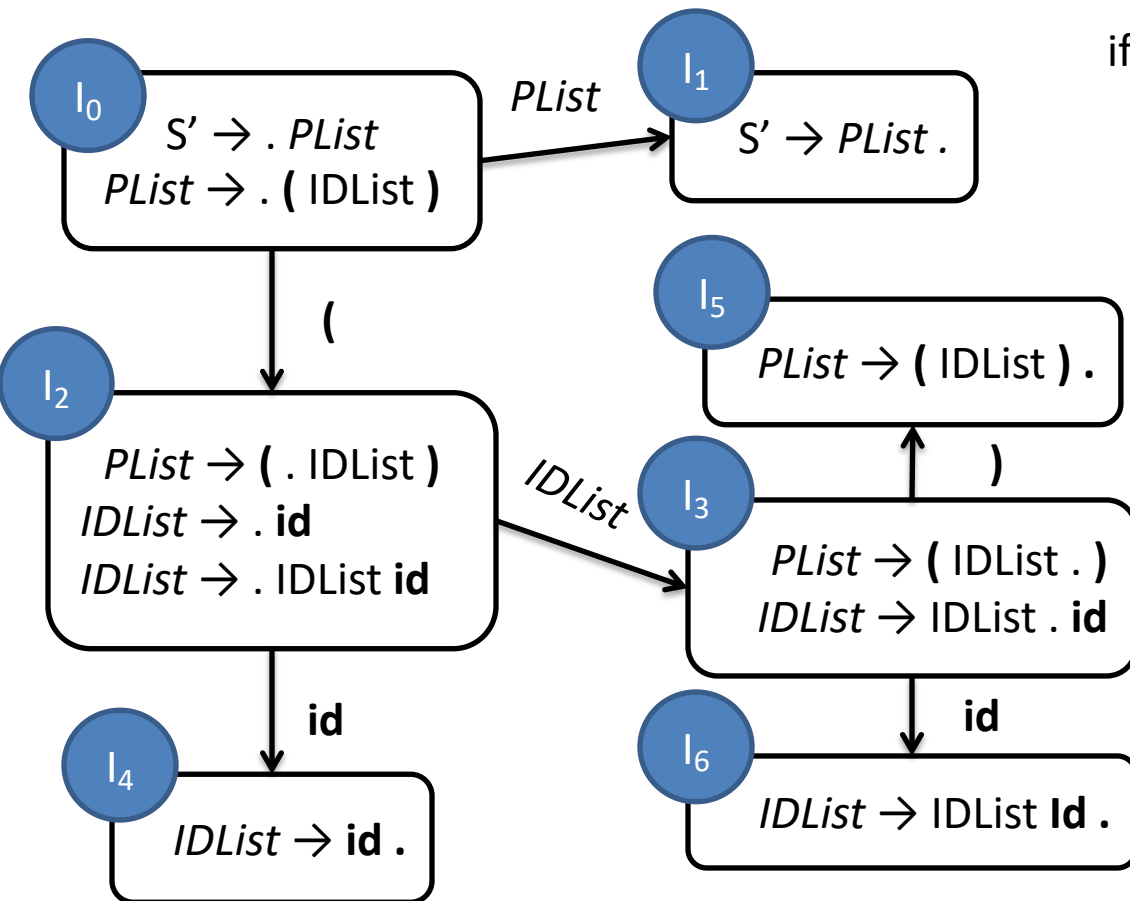
if state i includes item $A \rightarrow \alpha$ . **t** $\beta$
where **t** is a terminal
and there is an (i,j) transition on **t**
set Action[i,**t**] = shift j

| | ( | ) | id | eof |
|---|---|---|---|---|
| 0 | S 2 | | | |
| 1 | | | | |
| 2 | | | S 4 | |
| 3 | | S 5 | S 6 | |
| 4 | | | | |
| 5 | | | | |
| 6 | | | | |

35

# Action Table: Reduce



$I_0$
S' → . PList
PList → . **(** IDList **)**

PList → $I_1$
S' → PList **.**

**(**

$I_2$
PList → **(** . IDList **)**
IDList → . **id**
IDList → . IDList **id**

IDList

$I_5$
PList → **(** IDList **) .**

**)**

$I_3$
PList → **(** IDList **.** **)**
IDList → IDList **.** **id**

**id**

$I_4$
IDList → **id .**

**id**

$I_6$
IDList → IDList **Id .**

**Grammar G**
❶ S' → PList
❷ PList → ( IDList )
❸ IDList → **id**
❹ IDList → IDList **id**

if state i includes item A → α .
  where A is not S'
    for each t in FOLLOW(A):
      set Action[i,**t**] = reduce by A → α

FOLLOW(IDList) = { **)**, **id** }
FOLLOW(PList) = { **eof** }

| | **(** | **)** | **id** | **eof** |
|---|---|---|---|---|
| 0 | S 2 | | | |
| 1 | | | | |
| 2 | | | S 4 | |
| 3 | | S 5 | S 6 | |
| 4 | | R ❸ | R ❸ | |
| 5 | | | | R ❷ |
| 6 | | R ❹ | R ❹ | |

36

# Action Table: Accept

if state i includes item S' → S .
set Action[i,**eof**] = accept

I₀: S' → . PList
PList → . **(** IDList **)**

PList → I₁: S' → PList .

**(**

I₂: PList → **(** . IDList **)**
IDList → . **id**
IDList → . IDList **id**

IDList → I₃: PList → **(** IDList . **)**
IDList → IDList . **id**

**)** → I₅: PList → **(** IDList **)** .

**id** → I₄: IDList → **id** .

**id** → I₆: IDList → IDList **Id** .

**Grammar G**
❶ S' → PList
❷ PList → ( IDList )
❸ IDList → **id**
❹ IDList → IDList **id**

| | **(** | **)** | **id** | **eof** |
|---|---|---|---|---|
| 0 | S 2 | | | |
| 1 | | | | ☺ |
| 2 | | | S 4 | |
| 3 | | S 5 | S 6 | |
| 4 | | R ❸ | R ❸ | |
| 5 | | | | R ❷ |
| 6 | | R ❹ | R ❹ | |

# Some Final Thoughts on LR Parsing

A bit complicated to build the parse table

– Fortunately, algorithms exist

Still not as powerful as CYK

– Shift/reduce: action table cell includes S and R

– Reduce/reduce: cell include > 1 R rule

SDT similar to LL(1)

– Embed SDT action numbers in action table

– Fire off on reduce rules