



中山大學
SUN YAT-SEN UNIVERSITY

计算机学院 (软件学院)

SCHOOL OF COMPUTER SCIENCE AND ENGINEERING

Compilation Principle

编译原理

第13讲：语法分析(10)

张献伟

xianweiz.github.io

DCS290, 4/6/2023



中山大學
SUN YAT-SEN UNIVERSITY

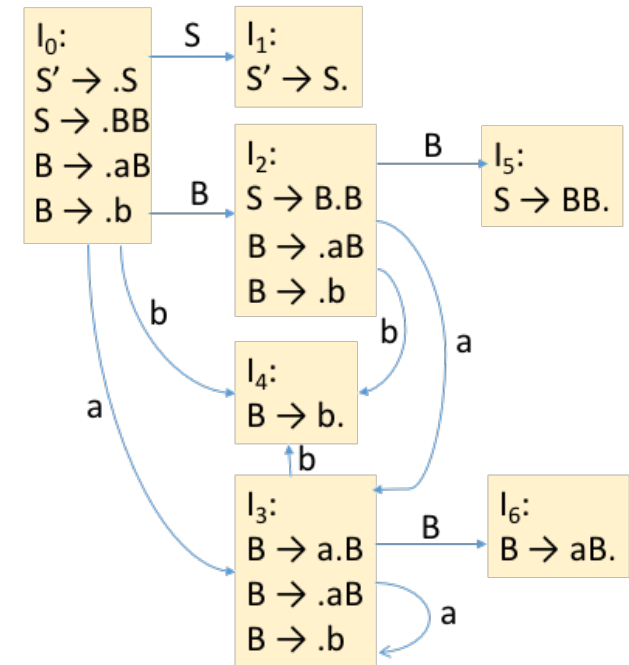


Review Questions

- Action table entries can be si and rj , what do they mean?
 - si : shift the input symbol and move to state i
 - rj : reduce by production numbered j
- Item/configuration: what does $A \rightarrow XYZ \bullet$ mean?
 - We have seen the body XYZ and it is time to reduce XYZ to A
- State: why we put the items into a configuration set?
 - Closure: we hope to see one symbol in $FIRST(Y)$
 - $Y \rightarrow u|w$
 - $A \rightarrow X \bullet YZ$
 - $Y \rightarrow \bullet u$
 - $Y \rightarrow \bullet w$
- What is augmented grammar?
 - Add one extra rule $S' \rightarrow S$ to guarantee only one 'acc' in the table
- What are the possible items of $S' \rightarrow S$?
 - $S' \rightarrow \bullet S$: initial item, haven't seen any input symbol
 - $S' \rightarrow S \bullet$: accept item, have reduced the input string to start symbol

LR(0) Parsing

- Construct LR(0) automaton from the Grammar [由文法构建自动机]
- Idea: assume
 - Input buffer contains α [但buffer不止有 α]
 - Next input is t [α 后是 t]
 - DFA on input α terminates in state s
 - α 处理完毕后处于状态 s
- Next: **reduce** by $X \rightarrow \beta$ if [归约]
 - s contains item $X \rightarrow \beta \cdot$
- Or, **shift** if [移进]
 - s contains item $X \rightarrow \beta \cdot tw$
 - Equivalent to saying s has a transition labeled t



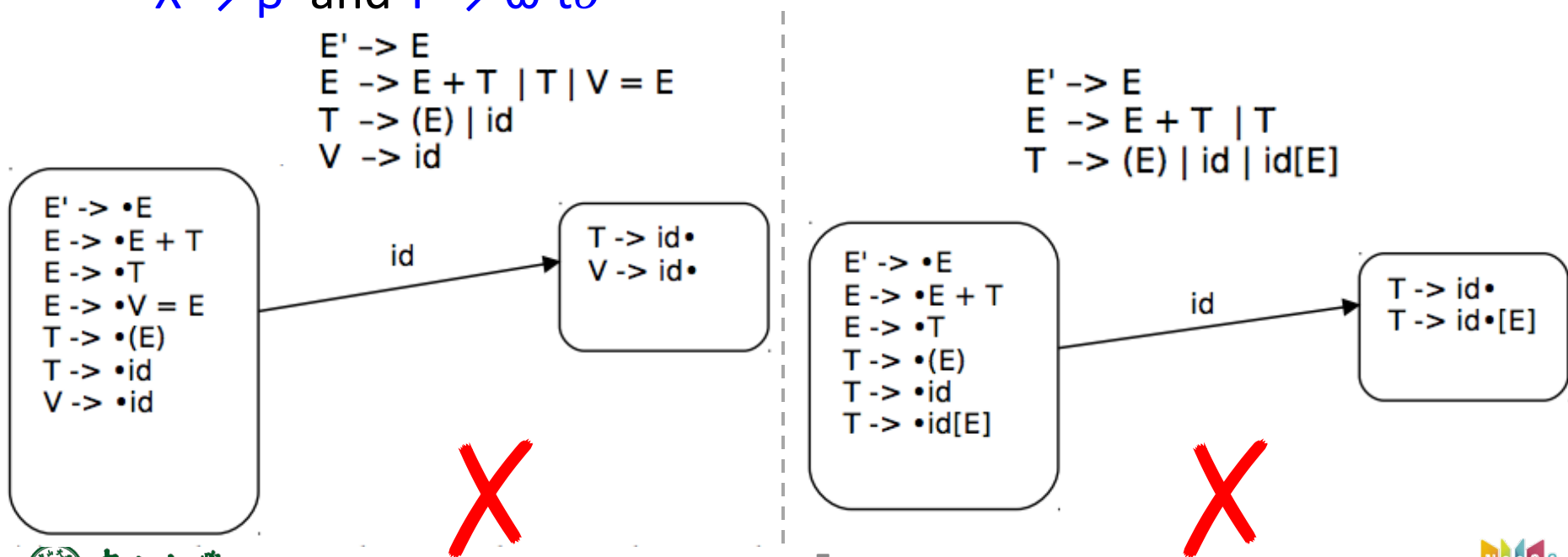
LR(0) Parsing (cont.)

- The parser must be able to determine what action to take in each state without looking at any further input symbols [没有展望就可以决定动作]
 - i.e. by only considering what the parsing stack contains so far
 - 🖱️ This is the '0' in the parser name
- In a LR(0) table, each state must only shift or reduce [确定性移进或归约]
 - Thus an LR(0) configuring set can only have exactly one reduce item [每个归约item自成一个状态]
 - cannot have both shift and reduce items; otherwise, conflicts

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

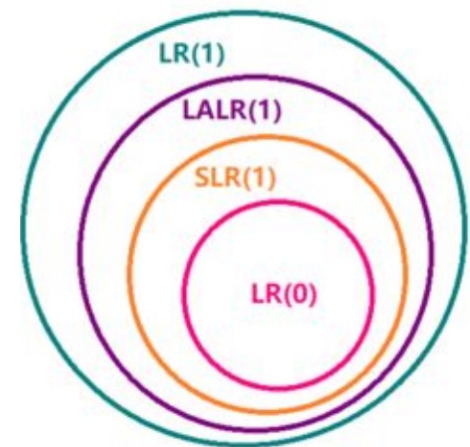
LR(0) Conflicts[冲突]

- LR(0) has a **reduce/reduce** conflict[归约-归约冲突] if:
 - Any state has two reduce items:
 - $X \rightarrow \beta \cdot$ and $Y \rightarrow \omega \cdot$
- LR(0) has a **shift/reduce** conflict[移进-归约冲突] if:
 - Any state has a reduce item and a shift item:
 - $X \rightarrow \beta \cdot$ and $Y \rightarrow \omega \cdot \tau$



LR(0) Summary[小结]

- LR(0) is the simplest LR parsing[最简单]
 - Table-driven shift-reduce parser[表驱动]
 - ACTION table[s, a] + GOTO table[s, X]
 - Weakest LR, not used much in practice[实际不常使用]
 - Parses without using any lookahead[没有任何展望]
- Adding just one token of lookahead vastly increases the parsing power[考虑展望]
 - SLR(1): simple LR(1), use FOLLOW[归约用FOLLOW]
 - LR(1): use dedicated symbols[比FOLLOW更精细]
 - LALR(1): balance SLR(1) and LR(1)[折衷]



SLR(1) Parsing

- LR(0) conflicts are generally caused by **reduce** actions
 - If the item is complete ($A \rightarrow \alpha.$), the parser must choose to reduce[项目形式完整就归约]

- Is this always appropriate?

- The next upcoming token may tell us sth different

- What tokens may tell the reduction is not appropriate?

- Perhaps **FOLLOW(A)** could be useful here

- If the sequence on top of the stack could be reduced to the nonterminal A, what tokens do we expect to find as the next input?

State	ACTION			GOTO	
	a	b	\$	S	B
0	s3	s4		1	2
1			acc		
2	s3	s4			5
3	s3	s4			6
4	r3	r3	r3		
5	r1	r1	r1		
6	r2	r2	r2		

- **SLR** = Simple LR

- Use the same LR(0) configuring sets and have the same table structure and parser operation[表结构一致]

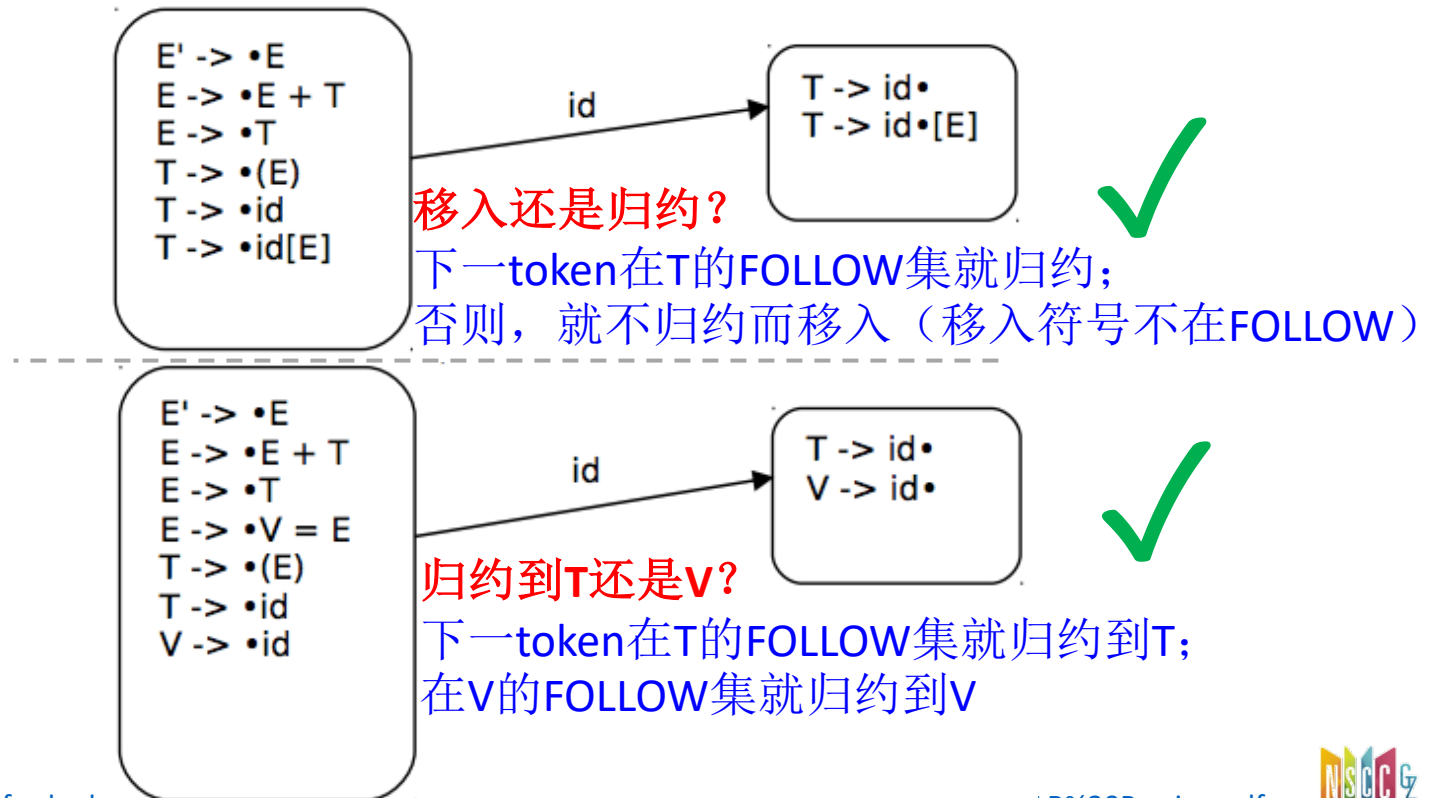
- The difference comes in assigning table actions[动作填充不同]

- Use one token of lookahead to help arbitrate among the conflicts

- Reduce only if the next input token is a member of the FOLLOW set of the nonterminal being reduced to[下一token在FOLLOW集才归约]

SLR(1) Parsing (cont.)

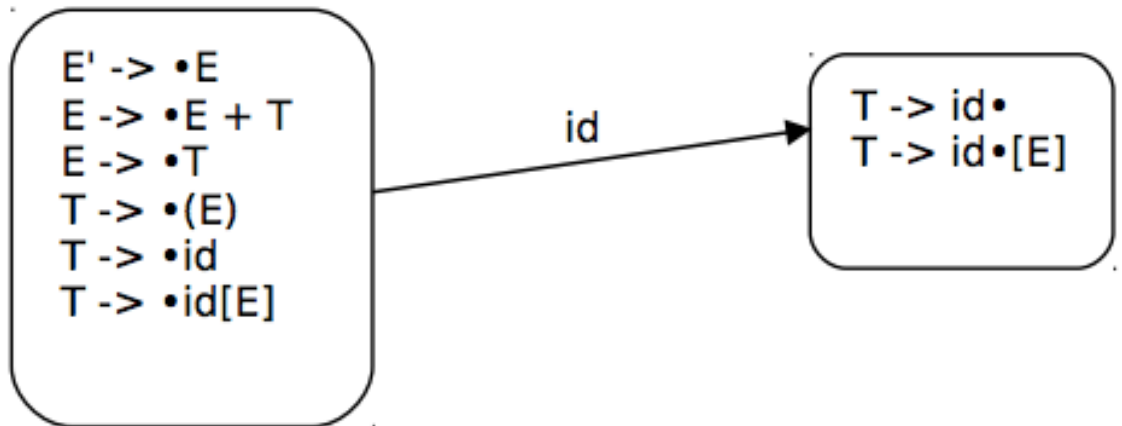
- In the SLR(1) parser, it is allowable for there to be both shift and reduce items in the same state as well as multiple reduce items
 - The SLR(1) parser will be able to determine which action to take as long as the FOLLOW sets are **disjoint**[可区分即可]



Example

- The first two LR(0) configurating sets entered if *id* is the first token of the input[用于识别id的前两个状态]
 - LR(0) parser: the set on the right side has a **shift-reduce conflict**
 - SLR(1) parser:
 - Compute FOLLOW(T) = { +,),], \$ }, i.e., only reduce on those tokens
 - FOLLOW(T) = FOLLOW(E) = {+,),], \$}
 - **id[id]**: next input is [, not in FOLLOW(T), **shift**
 - **id + id**: next input is +, in FOLLOW(T), **reduce**

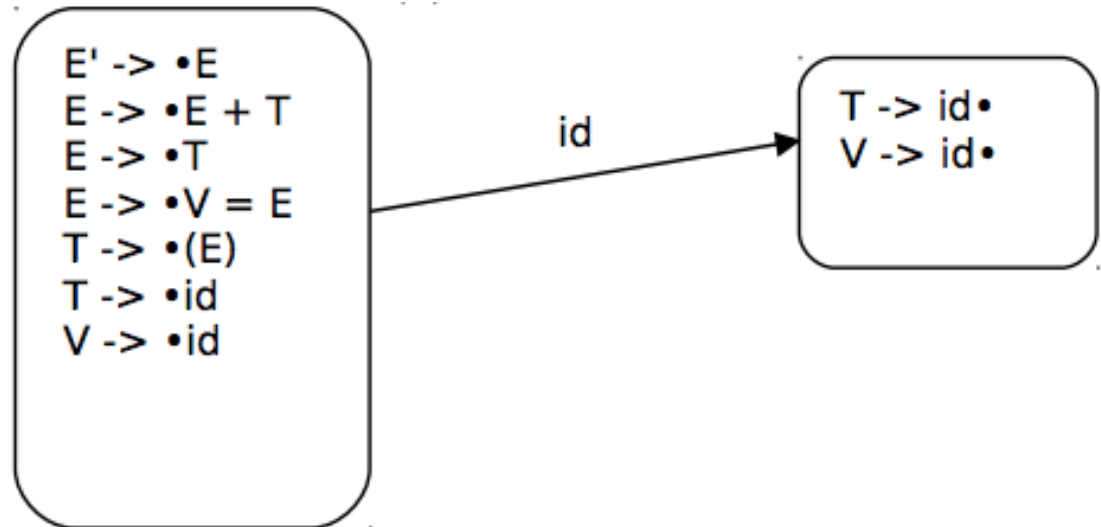
$E' \rightarrow E$
 $E \rightarrow E + T \mid T$
 $T \rightarrow (E) \mid id \mid id[E]$



Example (cont.)

- The first two LR(0) configurating sets entered if *id* is the first token of the input[用于识别id的前两个状态]
 - LR(0) parser: the right set has a **reduce-reduce conflict**
 - SLR(1) parser:
 - Capable to distinguish which reduction to apply depending on the next input token, **no conflict**
 - Compute FOLLOW(T) = { +,), \$ } and FOLLOW(V) = { = } id + id
id = id

$E' \rightarrow E$
 $E \rightarrow E + T \mid T \mid V = E$
 $T \rightarrow (E) \mid id$
 $V \rightarrow id$



SLR(1) Grammars[文法]

- A grammar is SLR(1) if the following two conditions hold for each configurating set[可区分]
- (1) For any item $A \rightarrow u \cdot x v$ in the set, with terminal x , there is no complete item $B \rightarrow w \cdot$ in that set with x in $\text{FOLLOW}(B)$ [无移入-归约冲突]
 - In the table, this translates no shift-reduce conflict on any state
- (2) For any two complete items $A \rightarrow u \cdot$ and $B \rightarrow v \cdot$ in the set, the follow sets must be disjoint, i.e. $\text{FOLLOW}(A) \cap \text{FOLLOW}(B)$ is empty[无归约-归约冲突]
 - This translates to no reduce-reduce conflict on any state
 - If more than one nonterminal could be reduced from this set, it must be possible to uniquely determine which using only one token of lookahead

SLR(1) Limitations[限制]

- SLR(1) vs. LR(0)
 - Adding just one token of lookahead and using the FOLLOW set greatly expands the class of grammars that can be parsed without conflict
- When we have a completed configuration (i.e., dot at the end) such as $X \rightarrow u \cdot$, we know that it is reducible[可归约]
 - We allow such a reduction whenever the next symbol is in $FOLLOW(X)$ [使用Follow集]
 - However, it may be that we should not reduce for every symbol in $FOLLOW(X)$, because the symbols below u on the stack preclude u being a handle for reduction in this case[Follow集不够]
 - In other words, SLR(1) states only tell us about the sequence on top of the stack, not what is below it on the stack
 - We may need to divide an SLR(1) state into separate states to differentiate the possible means by which that sequence has appeared on the stack[额外使用栈信息, FOLLOW是input buffer信息]



Example

- For input string: $id = id$, at I_2 after having reduced id_{Left} to L
 - Initially, at S_0
 - Move to S_5 , after shifting id to stack (S_5 is also pushed to stack)
 - Reduce, and back to S_0 , and further GOTO S_2
 - S_5 has a completed item, and next '=' is in FOLLOW(L)
 - S_5 and id are popped from stack, and L is pushed onto stack
 - $GOTO(S_0, L) = S_2$

$S' \rightarrow S$
 $S \rightarrow L = R$
 $S \rightarrow R$
 $L \rightarrow *R$
 $L \rightarrow id$
 $R \rightarrow L$

$I_0: S' \rightarrow \bullet S$
 $S \rightarrow \bullet L = R$
 $S \rightarrow \bullet R$
 $L \rightarrow \bullet *R$
 $L \rightarrow \bullet id$
 $R \rightarrow \bullet L$

$I_5: L \rightarrow id \bullet$

$I_6: S \rightarrow L = \bullet R$
 $R \rightarrow \bullet L$
 $L \rightarrow \bullet *R$
 $L \rightarrow \bullet id$

$I_1: S' \rightarrow S \bullet$

$I_7: L \rightarrow *R \bullet$

$I_2: S \rightarrow L \bullet = R$
 $R \rightarrow L \bullet$

$I_8: R \rightarrow L \bullet$

$I_3: S \rightarrow R \bullet$

$I_9: S \rightarrow L = R \bullet$

$I_4: L \rightarrow * \bullet R$
 $R \rightarrow \bullet L$
 $L \rightarrow \bullet *R$
 $L \rightarrow \bullet id$

Example (cont.)

- Choices upon seeing = coming up in the input:

```

S' -> S
S -> L = R
S -> R
L -> *R
L -> id
R -> L
    
```

- Action[2, =] = s6
 - Move on to find the rest of assignment
- Action[2, =] = r5
 - = ∈ FOLLOW(R): S => L=R => *R = R

- Shift-reduce conflict

- SLR parser fails to remember enough info
- Reduce using R -> L only after seeing * or =

I₀: S' -> •S
 S -> •L = R
 S -> •R
 L -> •*R
 L -> •id
 R -> •L

I₅: L -> id•
 I₆: S -> L = •R
 R -> •L
 L -> •*R
 L -> •id

I₁: S' -> S•

I₇: L -> *R•

I₂: S -> L• = R
 R -> L•

I₈: R -> L•

I₉: S -> L = R•

I₃: S -> R•

I₄: L -> *•R
 R -> •L
 L -> •*R
 L -> •id

For any item $A \rightarrow u \cdot x v$ in the set, with terminal x , there is no complete item $B \rightarrow w \cdot$ in that set with x in FOLLOW(B)

SLR(1) Improvement[改进]

- We don't need to see additional symbols beyond the first token in the input, we have already seen the info that allows us to determine the correct choice[展望信息已足够]
- Retain a little more of the left **context** that brought us here[历史路径]
 - Divide an SLR(1) state into separate states to differentiate the possible means by which that sequence has appeared on the stack
- Just using the entire FOLLOW set is not discriminating enough as the guide for when to reduce[FOLLOW集不够]
 - For the example, the FOLLOW set contains symbols that can follow R in any position within a valid sentence
 - But it does not precisely indicate which symbols follow R at this particular point in a derivation

LR(1) Parsing

- LR parsing adds the required extra info into the state
 - By redefining items to include a **terminal symbol** as an added component[让项目中包含终结符]
- General form of **LR(1) items**[项目]
 - $A \rightarrow X_1 \dots X_i \bullet X_{i+1} \dots X_j, a$
 - We have states $X_1 \dots X_i$ on the stack and are looking to put states $X_{i+1} \dots X_j$ on the stack and then reduce
 - But **only if** the token following X_j is the terminal a
 - a is called the lookahead of the configuration
- The lookahead **only** works with completed items[完成项]
 - $A \rightarrow X_1 \dots X_j \bullet, a$
 - All states are now on the stack, but only reduce when next symbol is a (a is either a terminal or \$)
 - Multi lookahead symbols: $A \rightarrow u \bullet, a/b/c$

LR(1) Parsing (cont.)

- When to reduce?
 - LR(0): if the configuration set has a **completed item** (i.e., dot at the end)
 - SLR(1): only if the next input token is in the **FOLLOW** set
 - LR(1): only if the next input token is exactly ***a*** (terminal or \$)
 - Trend: **more and more precise**
- **LR(1) items**: LR(0) item + lookahead terminals
 - Many differ only in their lookahead components[仅展望不同]
 - The extra lookahead terminals allow to make parsing decisions beyond the SLR(1) capability, but with **a big price**[代价]
 - More distinguished items and thus more sets
 - Greatly increased GOTO and ACTION table sizes

$S' \rightarrow \cdot S$

LR(0)

$S' \rightarrow \cdot S, \$$

LR(1)

LR(1) Construction

- Configuration sets

- Sets construction are essentially the same with SLR, but differing on Closure() and Goto()
 - Because we must respect the lookahead

- Closure()

- For each item $[A \rightarrow u \cdot Bv, a]$ in I , for each production rule $B \rightarrow w$ in G' , add $[B \rightarrow \cdot w, b]$ to I , if
 - $b \in \text{FIRST}(va)$ and $[B \rightarrow \cdot w, b]$ is not already in I
- Lookahead is the **FIRST(va)**, which are what can follow B
 - v can be nullable

(0) $S' \rightarrow S$

(1) $S \rightarrow XX$

(2) $X \rightarrow aX$

(3) $X \rightarrow b$

$S' \rightarrow \cdot S, \$$

I_0 :
 $S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot XX, \text{First}(\epsilon \$)$
 $X \rightarrow \cdot aX, \text{First}(X \$)$
 $X \rightarrow \cdot b, \text{First}(X \$)$

I_0 :
 $S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot XX, \$$
 $X \rightarrow \cdot aX, a/b$
 $X \rightarrow \cdot b, a/b$

LR(1) Construction (cont.)

- **Goto(I, X)**

- For item $[A \rightarrow u \cdot Xv, a]$ in I , $\text{Goto}(I, X) = \text{Closure}([A \rightarrow uX \cdot v, a])$
- Basically the same Goto function as defined for LR(0)
 - But have to **propagate the lookahead**[传递] when computing the transitions

- Overall steps

- Start from the initial set $\text{Closure}([S' \rightarrow \cdot S, \$])$
- Construct configuration sets following $\text{Goto}(I, X)$
- Repeat until no new sets can be added

I_0 :
 $S' \rightarrow \cdot S, \$$
 $S \rightarrow \cdot XX, \$$
 $X \rightarrow \cdot aX, a/b$
 $X \rightarrow \cdot b, a/b$

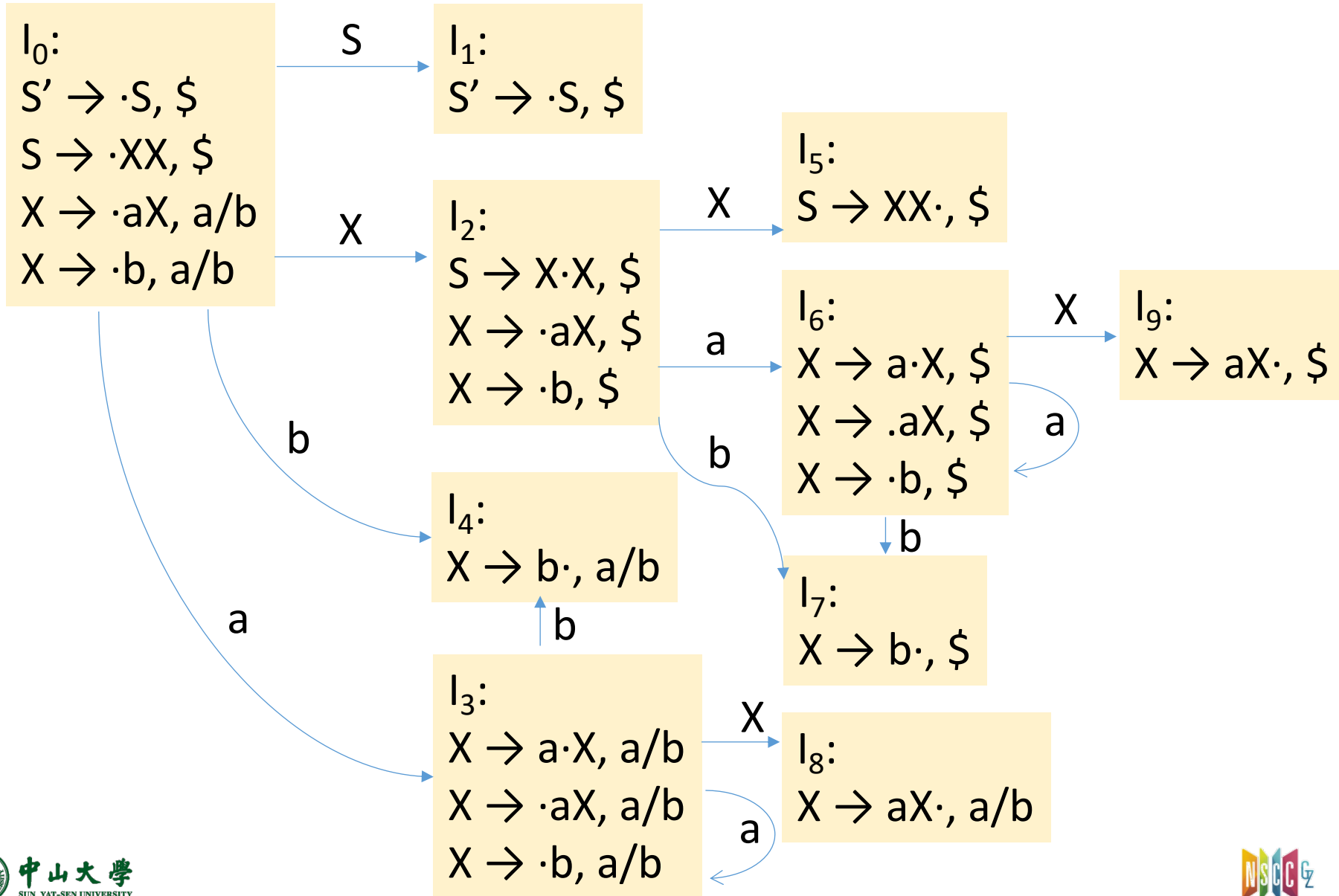


I_2 :
 $S \rightarrow X \cdot X, \$$
 $X \rightarrow \cdot aX, \text{First}(\epsilon \$)$
 $X \rightarrow \cdot b, \text{First}(\epsilon \$)$

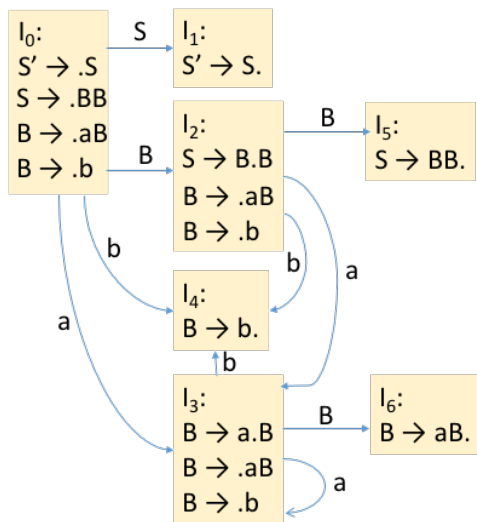


I_2 :
 $S \rightarrow X \cdot X, \$$
 $X \rightarrow \cdot aX, \$$
 $X \rightarrow \cdot b, \$$

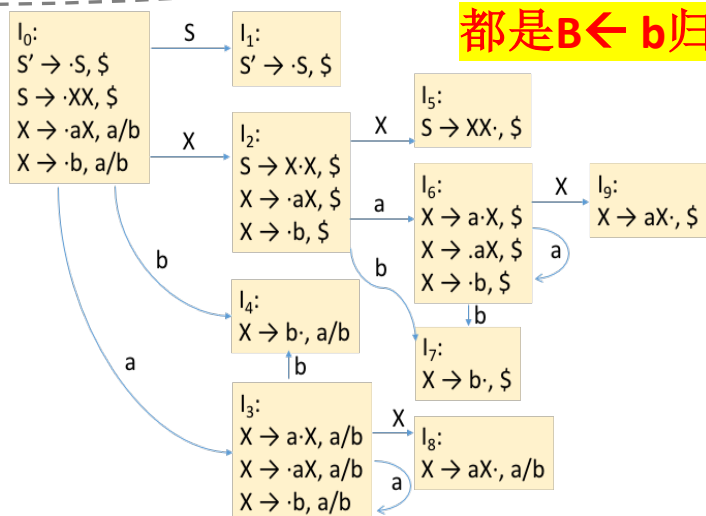
Example



Example (cont.)



LR(0)



LR(1)

String: bab

$\#bab\$ \Rightarrow b\#ab\$ \Rightarrow B\#ab\$ \Rightarrow Ba\#b\$$
 $\begin{matrix} 0 & & 04 & & 02 & & 023 \end{matrix}$
 $\Rightarrow Bab\#\$ \Rightarrow BaB\#\$ \Rightarrow BB\#\$ \Rightarrow S\#\$$
 $\begin{matrix} 0234 & 0236 & 025 & 01 \end{matrix}$

String: bab

$\#bab\$ \Rightarrow b\#ab\$ \Rightarrow B\#ab\$ \Rightarrow Ba\#b\$$
 $\begin{matrix} 0 & & 04 & & 02 & & 026 \end{matrix}$
 $\Rightarrow Bab\#\$ \Rightarrow BaB\#\$ \Rightarrow BB\#\$ \Rightarrow S\#\$$
 $\begin{matrix} 0237 & 0239 & 025 & 01 \end{matrix}$

String: abb

$\#abb\$ \Rightarrow a\#bb\$ \Rightarrow ab\#b\$ \Rightarrow aB\#b\$$
 $\begin{matrix} 0 & & 03 & & 034 & & 038 \end{matrix}$
 $\Rightarrow B\#b\$ \Rightarrow Bb\#\$ \Rightarrow BB\#\$ \Rightarrow S\#\$$
 $\begin{matrix} 02 & 027 & 025 & 01 \end{matrix}$

都是 $B \leftarrow b$ 归约

都是 $B \leftarrow aB$ 归约