

## Bottom-up Parsing

Bottom-up parsing is

- more general than (deterministic) top-down parsing (but just as efficient)
- builds on ideas in top-down parsing
- hence, is the preferred method in parser generation tools

Revert to the (unambiguous) natural grammar for our example

```
E --> T + E | T
T --> int * T | int | (E)
```

Consider the string `int * int + int`

Bottom-up parsing *reduces* a string to the start symbol by inverting productions

<code>int * int + int</code>	<code>int * T + int</code>	<code>T --&gt; int</code>
<code>int * T + int</code>	<code>T + int</code>	<code>T --&gt; int * T</code>
<code>T + int</code>	<code>T + T</code>	<code>T --&gt; int</code>
<code>T + T</code>	<code>T + E</code>	<code>E --&gt; T</code>
<code>T + E</code>	<code>E</code>	<code>E --&gt; T + E</code>
<code>E</code>		

Reading the productions bottom-to-top, these are the productions.

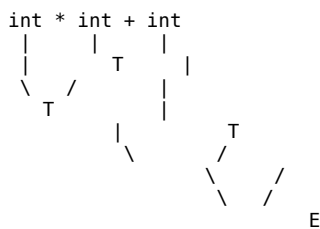
Reading them top-to-bottom, these are reductions

What's the point? The productions read backwards, trace a rightmost derivation, i.e., we are producing the right-most non-terminal at each step

**Important fact #1** about bottom-up parsing

*Bottom-up parsing traces a right-most derivation in reverse*

Picture



## Shift-reduce parsing : strategy used by bottom-up parsers

Important fact #1 has an interesting consequence:

- Let  $abw$  be a step of the bottom-up parse
- Assume the next reduction is by  $X \rightarrow b$
- Then  $w$  is a string of terminals
  - because  $aXw \rightarrow abw$  is a step in a right-most derivation

Idea: split string into two substrings

- Right substring is as yet unexamined by parsing
- Left substring has terminals and non-terminals
- The dividing point is marked by  $|$ 
  - $aX|w$
- Need two kind of moves: shift-move (discussed next) and reduce-move (already discussed)
- The shift-move moves  $|$  to the right by one character (signifying that the parser has seen this character)
- The reduce-move applies an inverse reduction to the right-hand of the left-hand string.
  - If  $A \rightarrow xy$  is a production, then

$$Cbxy|ijk \Rightarrow CbA|ijk$$

Example : `int * int + int`

Assume an oracle tells us whether to shift or reduce

<code> int * int + int</code>	shift
<code>int  * int + int</code>	shift
<code>int *  int + int</code>	shift
<code>int * int  + int</code>	shift
<code>int * T  + int</code>	reduce
<code>T   + int</code>	reduce
<code>T +  int</code>	shift
<code>T + int </code>	shift
<code>T + T </code>	reduce
<code>T + E </code>	reduce
<code>E </code>	reduce

Left string can be implemented by a stack, because we only reduce the suffix of the left string. Top of the stack is  $|$

shift pushes a terminal on the stack

reduce pops off of the stack (production rhs) and pushes a non-terminal on the stack (production lhs)

In a given state, more than one action (shift or reduce) may lead to a valid parse.

If it is legal to shift or reduce, there is a *shift-reduce* conflict. Need some techniques to remove them

If it is legal to reduce using two productions, there is a *reduce-reduce* conflict. Are always bad and usually indicate a serious problem with the grammar.

In either case, the parser does not know what to do, and we either need to rewrite the grammar, or need to give the parser a hint on what to do in this situation.

## Deciding when to shift and when to reduce

Example: `int * int + int`

Assume an oracle tells us whether to shift or reduce

```
int * int + int      shift
int | * int + int    shift
```

At this point, we could reduce.

```
T | * int + int      reduce
```

But this would be a fatal mistake because there is no way we can reduce to E (there is no production that begins with T \*).

Intuition: want to reduce only if the result can still be reduced to the start symbol (E)

Assume a right-most derivation:

$S \rightarrow^* \alpha X \omega \rightarrow \alpha \beta \omega$

Then  $\alpha \beta$  is a handle of  $\alpha \beta \omega$

Handles formalize the intuition: A handle is a reduction that allows further reductions back to the start symbol.

We only want to reduce at handles.

Note: we have said what a handle is, not how to find handles.

**Important fact #2:** In shift-reduce parsing, handles appear only at the top of the stack, never inside. Proof by induction

1. True initially, stack is empty
2. Immediately after reducing a handle
  - rightmost non-terminal on top of stack
  - next handle must be to right of rightmost non-terminal, because this is a right-most derivation.
  - i.e., sequence of shift moves reaches the next handle

Handles are never to the left of the rightmost non-terminal, and are always at the top of the stack.

Hence shift-reduce parsers can only shift or reduce at the top of the stack. However, how to recognize handles, i.e., when to shift and when to reduce?

## Recognizing handles

For an arbitrary grammar, no efficient algorithms known to recognize handles.

Heuristics to identify handles. On certain types of CFGs, heuristics are always correct.

Venn diagram: All CFGs  $\supset$  Unambiguous CFGs  $\supset$  LR(k) CFGs  $\supset$  LALR(k) CFGs  $\supset$  SLR(k)  $\supset$  LR(0)

LR(k) are quite general, but most practical grammars are LALR(k). SLR(k) are simplifications over LALR(k) [simple LR grammars]. There are grammars that are SLR(k) but not LALR(k) and so on.

$\alpha$  is a *viable prefix* if there is a valid right-most derivation of the string:  $S' \rightarrow \dots \rightarrow \alpha \omega \rightarrow \dots \rightarrow \text{string}$ . e.g.,  $\epsilon$  is a viable prefix

In other words,  $\alpha$  is a *viable prefix* if there is an  $\omega$  such that  $\alpha \omega$  is a state of a shift-reduce parser. Here  $\alpha$  is the stack and  $\omega$  is the rest of the input. It can lookahead at  $\omega$ , but the parser does not know the whole thing. The parser knows the whole stack.

What does this mean? A viable prefix does not extend past the right end of the handle. It's a viable prefix because it is a prefix of the handle. As long as a parser has viable prefixes on the stack, no parsing error has been detected.

An item is a production with a "." somewhere on the RHS in the production: e.g., the items for  $T \rightarrow (E)$  are:  $T \rightarrow \cdot (E)$ ,  $T \rightarrow (\cdot E)$ ,  $T \rightarrow (E \cdot)$ ,  $T \rightarrow (E) \cdot$ .

The only item for an  $\epsilon$  production,  $X \rightarrow \epsilon$  is  $X \rightarrow \cdot$ . Items are often called "LR(0) items".

The problem in recognizing viable prefixes is that the stack has only bits and pieces of the RHS of productions

- If it had a complete RHS, we could reduce

These bits and pieces are always prefixes of RHS of some production(s).

Consider the input  $(int)$  for the grammar:

```
E → T + E | T
T → int * T | int | (E)
```

- Then  $(E|)$  is a state of a shift-reduce parse
- $(E$  is a prefix of the RHS of  $T \rightarrow (E)$ 
  - Will be reduced after the next shift
- Item  $T \rightarrow (E \cdot)$  says that so far we have seen  $(E$  of this production and hope to see  $)$

The stack may have many prefixes on RHS's:

Prefix1 Prefix2 Prefix3 ... Prefix(n-1) Prefix(n)

Let Prefix(i) be a prefix of RHS of  $X_i \rightarrow \alpha_i$

- Prefix(i) will eventually reduce to  $X(i)$
- The missing part of  $\alpha_{(i-1)}$  starts with  $X(i)$
- i.e., there is a  $X(i-1) \rightarrow \text{Prefix}(i-1)X(i)\beta$  production
- and so on.. (e.g., there is a  $X(i-2) \rightarrow \text{Prefix}(i-2)X(i-1)\gamma$  production.

Recursively, Prefix(k+1) ... Prefix(n-1) Prefix(n) eventually reduces to the missing part of  $\alpha(k)$

Important fact #3 about bottom-up parsing: *For any grammar, the set of viable prefixes is a regular language.* The regular language represents the language formed by concatenating 0 or more prefixes of the productions (items).

For example, the language of viable prefixes for the example grammar:

```
S → ε | [S]
```

is

```
ε | "[" | "["* | "["*S | "["*S"]"
```

The language of viable prefixes for the example grammar:

```
S → ε | [S] | S.S
```

is

```
X = ε | "[" | "["* | "["*S | "["*S"]"
Y = "["*S.(X | Y)*
Z = X + Y
```

Conversely, In a string is parse-able through the bottom-up parser, then every potential state of the stack should always be a viable prefix.

Consider the string  $(int * int)$ :

- $int * | int)$  is a state of a shift-reduce parse
- $($  is a prefix of the RHS of  $T \rightarrow (E)$
- $\epsilon$  is a prefix of the RHS of  $E \rightarrow T$ . This is an interesting case as we are considering  $\epsilon$  as a prefix too!
- $int *$  is a prefix of the RHS of  $T \rightarrow int * T$

Alternatively, we can represent this as a "stack of items"

```
T → ( . E
E → . T
T → int * . T
```

says that

- We have seen  $($  of  $T \rightarrow (E)$
- We have seen  $\epsilon$  of  $E \rightarrow T$
- We have seen  $int *$  of  $T \rightarrow int * T$

Reading backwards, the LHS of every item becomes the RHS of the predecessor production.

In other words, every viable prefix can be represented as a stack of items, where the  $(n+1)$ th item is a production for a non-terminal that follows the "." in the  $n$ th item.

To recognize viable prefixes, we must

- Recognize a sequence of partial RHS's of productions, where
- Each partial RHS can eventually reduce to part of the missing suffix of its predecessor

## Recognizing Viable Prefixes

1. Add a dummy production  $S' \rightarrow S$  to  $G$   
This ensures that the start symbol ( $S'$ ) is used only in one place which is the LHS of this production
2. We will construct an NFA that will behave as follows:  
$$\text{NFA}(\text{stack}) = \begin{cases} \text{yes} & \text{if stack is a viable prefix} \\ \text{no} & \text{otherwise} \end{cases}$$
3. The NFA will read the input (stack) bottom-to-top
4. The NFA states are the items of  $G$ 
  - Including the extra production  $S' \rightarrow S$
5. For item  $E \rightarrow \alpha.X\beta$ , add transition from  $(E \rightarrow \alpha.X\beta) \xrightarrow{X} (E \rightarrow \alpha X\beta)$ 
  - i.e., if we see  $X$  in the left state, then we go to the right state.
6. For item  $E \rightarrow \alpha.X\beta$  and production  $X \rightarrow \gamma$ , add  $(E \rightarrow \alpha.X\beta) \xrightarrow{\epsilon} (X \rightarrow \gamma)$
7. Every state is an accepting state (i.e., if the entire stack is consumed, the stack is a viable prefix)
8. Start state is  $(S' \rightarrow S)$

## Valid Items

- Can construct a DFA from NFA using the standard subset-of-states construction
- The states of the DFA are "canonical collections of items" or "canonical collections of LR(0) items"
  - The dragon book gives another way of constructing the LR(0) items

Item  $X \rightarrow \beta.\gamma$  is *valid* for a viable prefix  $\alpha\beta$  if

$S' \rightarrow^* \alpha X \omega \rightarrow \alpha \beta \gamma \omega$

by a right-most derivation.

After parsing  $\alpha\beta$ , the valid items are the possible tops of the stack of items

An item  $I$  is valid for a viable prefix  $\alpha$  if the DFA recognizing viable prefixes terminates on input  $\alpha$  in a state  $s$  containing  $I$ .

The items in  $s$  describe what the top of the item stack might be after reading input  $\alpha$

An item is often valid for many prefixes. e.g., The item  $T \rightarrow (.E)$  is valid for prefixes

```
(
((
(((
((((
...
```

We can see this by looking at the DFA, which will keep looping into the same DFA state for each open paren. Need to show the NFA and DFA construction for our example grammar, and the valid items for these string prefixes. Will need a laptop and a projector!

## Simple LR (SLR) parsing

- LR(0) Parsing: Assume
  - stack contains  $\alpha$
  - next input is  $t$
  - DFA on input  $\alpha$  terminates in state  $s$
- Reduce by  $X \rightarrow \beta$  if
  - $s$  contains item  $X \rightarrow \beta$
- Shift if
  - $s$  contains item  $X \rightarrow \beta.t\omega$
  - equivalent to saying that  $s$  has a transition labeled  $t$
- LR(0) has a reduce/reduce conflict if:
  - Any state has two reduce items:  $X \rightarrow \beta$  and  $Y \rightarrow \omega$ .
- LR(0) has a shift/reduce conflict if:
  - Any state has a reduce item and a shift item:  $X \rightarrow \beta$  and  $Y \rightarrow \omega.t\delta$

SLR = "Simple LR": improves on LR(0) shift/reduce heuristics so fewer states have conflicts

Idea: Assume

- stack contains  $\alpha$
- next input is  $t$
- DFA on input  $\alpha$  terminates in state  $s$
- Reduce by  $X \rightarrow \beta$  if
  - $s$  contains item  $X \rightarrow \beta$
  - $t \in \text{Follow}(X)$  [only change to LR(0)]
- Shift if
  - $s$  contains item  $X \rightarrow \beta.t\omega$

If there are conflicts under these rules, the grammar is not SLR

- In other words, SLR grammars are those where the heuristics detect exactly the handles
- SLR(1) grammars work with the SLR algorithm and one lookahead
- Lots of grammars are not SLR
  - including all ambiguous grammars
- We can parse more grammars by using precedence declarations
  - Instructions for resolving conflicts
- Consider the ambiguous grammar:
  - $E \rightarrow E+E|E^*E|(E)|int$
- The DFA for this grammar contains a state with the following items:
  - $E \rightarrow E^*E$ . and  $E \rightarrow E.E$
  - shift/reduce conflict!
- Declaring "\*" has higher precedence than "+" resolves this conflict in favor of reducing
- The term "precedence declaration" is misleading
- These declarations do not define precedence; they define conflict resolutions
  - Not quite the same thing!
  - Tools allow you to print out the parsing automaton, and you will be able to see the conflict resolution in the automaton (recommended)

#### SLR Parsing algorithm

1. Let M be DFA for viable prefixes of G
2. Let  $|x_1 \dots x_n\$$  be initial configuration
3. Repeat until configuration is  $S|\$$ 
  - Let  $\alpha|\omega$  be current configuration
  - Run M on current stack  $\alpha$
  - If M rejects  $\alpha$ , report parsing error
    - Stack  $\alpha$  is not a viable prefix
  - If M accepts  $\alpha$  with items I, let a be next input
  - Shift if  $X \rightarrow \alpha.a\gamma$  in I
  - Reduce if  $X \rightarrow \alpha$  in I and a in Follow(X)
  - Report parsing error if neither applies

If there is a conflict in the last step, grammar is not SLR(k). k is the amount of lookahead (in practice,  $k = 1$ )

#### SLR Improvements

- Rerunning the viable prefixes automaton on the stack at each step is wasteful
- Most of the work is repeated
- Remember the state of the automaton on each prefix of the stack
- Change stack to contain pairs  $\langle \text{Symbol}, \text{DFA state} \rangle$
- For a stack:  $\langle \text{sym}_1, \text{state}_1 \rangle \dots \langle \text{sym}_n, \text{state}_n \rangle$ 
  - $\text{state}_n$  is the final state of the DFA on  $\text{sym}_1 \dots \text{sym}_n$
- Detail: the bottom of the stack is  $\langle \text{any}, \text{start} \rangle$  where
  - any is any dummy symbol
  - start is the start state of the DFA
- Define  $\text{goto}[i, A] = j$  if  $\text{state}_i \rightarrow A \text{ state}_j$
- $\text{goto}$  is just the transition function of the DFA
- Shift x
  - Push  $\langle a, x \rangle$  on the stack
  - a is current input
  - x is a DFA state
- Reduce  $X \rightarrow \alpha$ 
  - As before
- Accept
- Error

Action table: For each state  $s_i$  and terminal a

- If  $s_i$  has item  $X \rightarrow \alpha.a\beta$  and  $\text{goto}[i, a] = j$ , then  $\text{action}[i, a] = \text{shift } j$
- If  $s_i$  has item  $X \rightarrow \alpha.$  and a in Follow(X) and  $X \neq S'$ , then  $\text{action}[i, a] = \text{reduce } X \rightarrow \alpha$
- If  $s_i$  has item  $S' \rightarrow S.$ , then  $\text{action}[i, \$] = \text{accept}$
- Otherwise,  $\text{action}[i, a] = \text{error}$

#### SLR Improvements

- Let  $I = w\$$  be initial input
- Let  $j = 0$
- Let DFA state 1 have item  $S' \rightarrow .S$
- Let stack =  $\langle \text{dummy}, 1 \rangle$
- repeat
  - case  $\text{action}[\text{top\_state}(\text{stack}), I[j]]$  of
    - shift k: push  $\langle I[j++], k \rangle$
    - reduce  $X \rightarrow A$ :
      - pop  $|A|$  pairs
      - push  $\langle X, \text{goto}[\text{top\_state}(\text{stack}), X] \rangle$
    - accept: halt normally
    - error: halt and report error
- Note that the algorithm uses only the DFA states and the input
  - The stack symbols are never used!
- However, we still need the symbols for semantic actions (e.g., code generation)
- Some common constructs are not SLR(1)
- LR(1) is more powerful
  - Build lookahead into the items (fine-grained disambiguation for each item, rather than just checking the follow set)
  - An LR(1) item is a pair: LR(0) item, lookahead
  - $[T \rightarrow \dots .int * T, \$]$  means
    - After seeing  $T \rightarrow \dots int * T$ , reduce if lookahead is  $\$$
  - More accurate than just using follow sets
  - Take a look at the LALR(1) automaton for your parser! (uses these items, but has a slight optimization over it)

- LALR automaton is composed of states containing LR(1) items, with the added optimization that if two states differ only in lookahead then we combine those states. If the resulting states (after this minimization), have no conflict in the grammar, then that grammar is LALR also.

## SLR Examples

```
S' --> S
S --> Sa
S --> b
```

SLR parsers do not mind left-recursive grammars

The first state in the corresponding DFA will look like ( $\epsilon$ -closure of the NFA state): (state 1)

```
S' --> .S
S --> .Sa
S --> .b
```

If we see a b in this state, we get another DFA state: (state 2)

```
S --> b.
```

Alternatively, if we see a S in this state, we get another DFA state: (state 3)

```
S' --> S.
S --> S.a
```

From this state, if we see a, we get (state 4)

```
S --> Sa.
```

The only state with a shift-reduce conflict is state3. Here, if we look at follow of S', we have only "\$", and hence we can resolve this conflict by one lookahead. Hence this is an SLR grammar

Another example grammar

```
S' --> S
S --> SaS
S --> b
```

Looking at the corresponding DFA: (state 1)

```
S' --> .S
S' --> .SaS
S --> .b
```

One possibility is that we see b in this state to get: (state 2)

```
S --> b.
```

Another possibility is that we see S in this state to get: (state 3)

```
S' --> S.
S' --> S.aS
```

If we get a in this state, we get (state 4)

```
S' --> Sa.S
S --> .SaS
S --> .b
```

(notice that we formed an  $\epsilon$ -closure of the first item to add more items

From here, if we get S, we get the following state: (state 5)

```
S --> SaS.
S --> S.aS
```

From here, if we get a again, we go back to state 4! If we get b, we go to state 2

The only states that have conflicts are: state3 (resolved by follow as follow(S') = \$) and state5 (has a shift/reduce conflict because a  $\in$  follow(S))

Thus this is not an SLR(1) grammar