

# Some definitions

---

## *Recall*

For a grammar  $G$ , with start symbol  $S$ , any string  $\alpha$  such that  $S \Rightarrow^* \alpha$  is called a *sentential form*

- If  $\alpha \in V_t^*$ , then  $\alpha$  is called a *sentence* in  $L(G)$
- Otherwise it is just a sentential form (not a sentence in  $L(G)$ )

A *left-sentential form* is a sentential form that occurs in the leftmost derivation of some sentence.

A *right-sentential form* is a sentential form that occurs in the rightmost derivation of some sentence.

# Bottom-up parsing

---

Goal:

*Given an input string  $w$  and a grammar  $G$ , construct a parse tree by starting at the leaves and working to the root.*

The parser repeatedly matches a *right-sentential* form from the language against the tree's upper frontier.

At each match, it applies a *reduction* to build on the frontier:

- each reduction matches an upper frontier of the partially built tree to the RHS of some production
- each reduction adds a node on top of the frontier

The final result is a rightmost derivation, in reverse.

# Example

---

Consider the grammar

$$\begin{array}{l|l} 1 & S \rightarrow aABe \\ 2 & A \rightarrow Abc \\ 3 & \quad | b \\ 4 & B \rightarrow d \end{array}$$

and the input string `abbcd`

| Prod'n. | Sentential Form  |
|---------|--|
| 3       | a <span style="border: 1px solid black; padding: 0 2px;">b</span> bcde |
| 2       | a <span style="border: 1px solid black; padding: 0 2px;">Abc</span> de |
| 4       | aA <span style="border: 1px solid black; padding: 0 2px;">d</span> e   |
| 1       | <span style="border: 1px solid black; padding: 2px;">aABe</span>       |
| –       | <i>S</i>   |

The trick appears to be scanning the input and finding valid sentential forms.

# Handles

---

*What are we trying to find?*

A substring  $\alpha$  of the tree's upper frontier that

matches some production  $A \rightarrow \alpha$  where reducing  $\alpha$  to  $A$  is one step in the reverse of a rightmost derivation

We call such a string a *handle*.

Formally:

a *handle* of a right-sentential form  $\gamma$  is a production  $A \rightarrow \beta$  and a position in  $\gamma$  where  $\beta$  may be found and replaced by  $A$  to produce the previous right-sentential form in a rightmost derivation of  $\gamma$ .

i.e., if  $S \Rightarrow_{\text{rm}}^* \alpha A w \Rightarrow_{\text{rm}} \alpha \beta w$  then  $A \rightarrow \beta$  in the position following  $\alpha$  is a handle of  $\alpha \beta w$

Because  $\gamma$  is a right-sentential form, the substring to the right of a handle contains only terminal symbols.

# Handles

---

*Theorem:*

If  $G$  is unambiguous then every right-sentential form has a unique handle.

*Proof: (by definition)*

1.  $G$  is unambiguous  $\Rightarrow$  rightmost derivation is unique
2.  $\Rightarrow$  a unique production  $A \rightarrow \beta$  applied to take  $\gamma_{i-1}$  to  $\gamma_i$
3.  $\Rightarrow$  a unique position  $k$  at which  $A \rightarrow \beta$  is applied
4.  $\Rightarrow$  a unique handle  $A \rightarrow \beta$

# Example

---

The left-recursive expression grammar

(original form)

|   |  |   |
|---|--|---|
| 1 |  | $\langle \text{goal} \rangle ::= \langle \text{expr} \rangle$                                 |
| 2 |  | $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$   |
| 3 |  | $\langle \text{expr} \rangle - \langle \text{term} \rangle$                                   |
| 4 |  | $\langle \text{term} \rangle$   |
| 5 |  | $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$ |
| 6 |  | $\langle \text{term} \rangle / \langle \text{factor} \rangle$                                 |
| 7 |  | $\langle \text{factor} \rangle$   |
| 8 |  | $\langle \text{factor} \rangle ::= \text{num}$  |
| 9 |  | $\langle \text{factor} \rangle \text{id}$   |

| Prod'n. | Sentential Form   |
|---------|---|
| —       | $\langle \text{goal} \rangle$   |
| 1       | $\langle \text{expr} \rangle$   |
| 3       | $\langle \text{expr} \rangle - \langle \text{term} \rangle$                                 |
| 5       | $\langle \text{expr} \rangle - \langle \text{term} \rangle * \langle \text{factor} \rangle$ |
| 9       | $\langle \text{expr} \rangle - \langle \text{term} \rangle * \underline{\text{id}}$         |
| 7       | $\langle \text{expr} \rangle - \underline{\langle \text{factor} \rangle} * \text{id}$       |
| 8       | $\langle \text{expr} \rangle - \underline{\text{num}} * \text{id}$                          |
| 4       | $\underline{\langle \text{term} \rangle} - \text{num} * \text{id}$                          |
| 7       | $\underline{\langle \text{factor} \rangle} - \text{num} * \text{id}$                        |
| 9       | $\underline{\text{id}} - \text{num} * \text{id}$  |

# Handle-pruning

---

The process to construct a bottom-up parse is called *handle-pruning*.

To construct a rightmost derivation

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \cdots \Rightarrow \gamma_{n-1} \Rightarrow \gamma_n = w$$

we set  $i$  to  $n$  and apply the following simple algorithm

for  $i = n$  downto 1

1. find the handle  $A_i \rightarrow \beta_i$  in  $\gamma_i$
2. replace  $\beta_i$  with  $A_i$  to generate  $\gamma_{i-1}$

*This takes  $2n$  steps, where  $n$  is the length of the derivation*

# Stack implementation

---

One scheme to implement a handle-pruning, bottom-up parser is called a *shift-reduce* parser.

Shift-reduce parsers use a *stack* and an *input buffer*

1. initialize stack with \$
2. Repeat until the top of the stack is the goal symbol and the input token is \$
  - a) *find the handle*  
if we don't have a handle on top of the stack, *shift* an input symbol onto the stack
  - b) *prune the handle*  
if we have a handle  $A \rightarrow \beta$  on the stack, *reduce*
    - i) pop  $|\beta|$  symbols off the stack
    - ii) push  $A$  onto the stack

## Example: back to $x - 2 * y$

|   | Stack  | Input         | Action   |
|---|--|---------------|----------|
|   | \$   | id - num * id | shift    |
| 1 $\langle \text{goal} \rangle ::= \langle \text{expr} \rangle$                                 | \$ <u>id</u>   | - num * id    | reduce 9 |
| 2 $\langle \text{expr} \rangle ::= \langle \text{expr} \rangle + \langle \text{term} \rangle$   | \$ <u><math>\langle \text{factor} \rangle</math></u>   | - num * id    | reduce 7 |
| 3   $\langle \text{expr} \rangle - \langle \text{term} \rangle$                                 | \$ <u><math>\langle \text{term} \rangle</math></u>   | - num * id    | reduce 4 |
| 4   $\langle \text{term} \rangle$   | \$ <u><math>\langle \text{expr} \rangle</math></u>   | - num * id    | shift    |
| 5 $\langle \text{term} \rangle ::= \langle \text{term} \rangle * \langle \text{factor} \rangle$ | \$ <u><math>\langle \text{expr} \rangle -</math></u>   | num * id      | shift    |
| 6   $\langle \text{term} \rangle / \langle \text{factor} \rangle$                               | \$ <u><math>\langle \text{expr} \rangle - \text{num}</math></u>  | * id          | reduce 8 |
| 7   $\langle \text{factor} \rangle$   | \$ <u><math>\langle \text{expr} \rangle - \langle \text{factor} \rangle</math></u>                               | * id          | reduce 7 |
| 8 $\langle \text{factor} \rangle ::= \text{num}$  | \$ <u><math>\langle \text{expr} \rangle - \langle \text{term} \rangle</math></u>                                 | * id          | shift    |
| 9   id  | \$ <u><math>\langle \text{expr} \rangle - \langle \text{term} \rangle * \text{id}</math></u>                     | id            | shift    |
|   | \$ <u><math>\langle \text{expr} \rangle - \langle \text{term} \rangle * \langle \text{factor} \rangle</math></u> |               | reduce 9 |
|   | \$ <u><math>\langle \text{expr} \rangle - \langle \text{term} \rangle</math></u>                                 |               | reduce 5 |
|   | \$ <u><math>\langle \text{expr} \rangle</math></u>   |               | reduce 3 |
|   | \$ <u><math>\langle \text{goal} \rangle</math></u>   |               | reduce 1 |
|   |  |               | accept   |

1. Shift until top of stack is the right end of a handle

2. Find the left end of the handle and reduce

5 shifts + 9 reduces + 1 accept

# Shift-reduce parsing

---

*Shift-reduce parsers are simple to understand*

A shift-reduce parser has just four canonical actions:

1. *shift* — next input symbol is shifted onto the top of the stack
2. *reduce* — right end of handle is on top of stack;  
locate left end of handle within the stack;  
pop handle off stack and push appropriate non-terminal LHS
3. *accept* — terminate parsing and signal success
4. *error* — call an error recovery routine

Key insight: recognize handles with a DFA.

# LR parsing

---

The skeleton parser:

```
push  $s_0$ 
token  $\leftarrow$  next_token()
repeat forever
  s  $\leftarrow$  top of stack
  if action[s,token] = "shift  $s_i$ " then
    push  $s_i$ 
    token  $\leftarrow$  next_token()
  else if action[s,token] = "reduce  $A \rightarrow \beta$ "
    then
      pop  $|\beta|$  states
       $s' \leftarrow$  top of stack
      push goto[ $s',A$ ]
  else if action[s, token] = "accept" then
    return
  else error()
```

This takes  $k$  shifts,  $l$  reduces, and 1 accept, where  $k$  is the length of the input string and  $l$  is the length of the reverse rightmost derivation

# Example tables

| state | ACTION |    |    |     | GOTO   |        |          |
|-------|--------|----|----|-----|--------|--------|----------|
|       | id     | +  | *  | \$  | ⟨expr⟩ | ⟨term⟩ | ⟨factor⟩ |
| 0     | s4     | -  | -  | -   | 1      | 2      | 3        |
| 1     | -      | -  | -  | acc | -      | -      | -        |
| 2     | -      | s5 | -  | r3  | -      | -      | -        |
| 3     | -      | r5 | s6 | r5  | -      | -      | -        |
| 4     | -      | r6 | r6 | r6  | -      | -      | -        |
| 5     | s4     | -  | -  | -   | 7      | 2      | 3        |
| 6     | s4     | -  | -  | -   | -      | 8      | 3        |
| 7     | -      | -  | -  | r2  | -      | -      | -        |
| 8     | -      | r4 | -  | r4  | -      | -      | -        |

## The Grammar

|   |          |     |                   |
|---|----------|-----|-------------------|
| 1 | ⟨goal⟩   | ::= | ⟨expr⟩            |
| 2 | ⟨expr⟩   | ::= | ⟨term⟩ + ⟨expr⟩   |
| 3 |          |     | ⟨term⟩            |
| 4 | ⟨term⟩   | ::= | ⟨factor⟩ * ⟨term⟩ |
| 5 |          |     | ⟨factor⟩          |
| 6 | ⟨factor⟩ | ::= | id                |

*Note:* This is a simple little right-recursive grammar; *not* the same as in previous lectures.

# Example using the tables

---

| Stack      | Input        | Action |
|------------|--------------|--------|
| \$ 0       | id* id+ id\$ | s4     |
| \$ 0 4     | * id+ id\$   | r6     |
| \$ 0 3     | * id+ id\$   | s6     |
| \$ 0 3 6   | id+ id\$     | s4     |
| \$ 0 3 6 4 | + id\$       | r6     |
| \$ 0 3 6 3 | + id\$       | r5     |
| \$ 0 3 6 8 | + id\$       | r4     |
| \$ 0 2     | + id\$       | s5     |
| \$ 0 2 5   | id\$         | s4     |
| \$ 0 2 5 4 | \$           | r6     |
| \$ 0 2 5 3 | \$           | r5     |
| \$ 0 2 5 2 | \$           | r3     |
| \$ 0 2 5 7 | \$           | r2     |
| \$ 0 1     | \$           | acc    |

# LR( $k$ ) grammars

---

Informally, we say that a grammar  $G$  is LR( $k$ ) if, given a rightmost derivation

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \cdots \Rightarrow \gamma_n = w,$$

we can, for each right-sentential form in the derivation,

1. *isolate the handle of each right-sentential form, and*
2. *determine the production by which to reduce*

by scanning  $\gamma_i$  from left to right, going at most  $k$  symbols beyond the right end of the handle of  $\gamma_i$ .

## Why study LR grammars?

---

LR(1) grammars are often used to construct parsers.

We call these parsers LR(1) parsers.

- everyone's favorite parser
- virtually all context-free programming language constructs can be expressed in an LR(1) form
- LR grammars are the most general grammars parsable by a deterministic, bottom-up parser
- efficient parsers can be implemented for LR(1) grammars
- LR parsers detect an error as soon as possible in a left-to-right scan of the input
- LR grammars describe a proper superset of the languages recognized by predictive (i.e., LL) parsers

# Parsing review

---

## *Recursive descent*

A hand coded recursive descent parser directly encodes a grammar (typically an LL(1) grammar) into a series of mutually recursive procedures. It has most of the linguistic limitations of LL(1).

## LL( $k$ )

An LL( $k$ ) parser must be able to recognize the use of a production after seeing only the first  $k$  symbols of its right hand side.

## LR( $k$ )

An LR( $k$ ) parser must be able to recognize the occurrence of the right hand side of a production after having seen all that is derived from that right hand side with  $k$  symbols of lookahead.

## Facts to remember

---

LR is more expressive than LL.

LL is more expressive than RE.

A RE can be parsed by a DFA; A CFG can be parsed by DFA+Stack (in LR). Why Stack?

# Applications

---

Machine translation.

Random test generation.

Reverse engineering.