# Chapter 4: LR Parsing

# 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

$$
\begin{array}{c|ccl}
1 & S & \rightarrow & \mathtt{a}AB\mathtt{e} \\
2 & A & \rightarrow & A\mathtt{bc} \\
3 & & | & \mathtt{b} \\
4 & B & \rightarrow & \mathtt{d}
\end{array}
$$

and the input string `abbcde`

| Prod'n. | Sentential Form |
|:---:|:---|
| 3 | a b bcde |
| 2 | a $A$bc de |
| 4 | a$A$ d e |
| 1 | a$AB$e |
| – | $S$ |

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 $A \to \alpha$ 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 $A \to \beta$ 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 $S \Rightarrow_{rm}^{*} \alpha A w \Rightarrow_{rm} \alpha\beta w$ then $A \to \beta$ 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



The handle $A \rightarrow \beta$ in the parse tree for $\alpha\beta w$

# Handles

*Theorem*:

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

*Proof: (by defi nition)*

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

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

# 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) *fi nd 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 $x - 2 * y$

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

| Stack | Input | Action |
|---|---|---|
| $\$$ | $\texttt{id} - \texttt{num} * \texttt{id}$ | shift |
| $\$\underline{\texttt{id}}$ | $- \texttt{num} * \texttt{id}$ | reduce 9 |
| $\$\underline{\langle\text{factor}\rangle}$ | $- \texttt{num} * \texttt{id}$ | reduce 7 |
| $\$\underline{\langle\text{term}\rangle}$ | $- \texttt{num} * \texttt{id}$ | reduce 4 |
| $\$\langle\text{expr}\rangle$ | $- \texttt{num} * \texttt{id}$ | shift |
| $\$\langle\text{expr}\rangle -$ | $\texttt{num} * \texttt{id}$ | shift |
| $\$\langle\text{expr}\rangle - \underline{\texttt{num}}$ | $* \texttt{id}$ | reduce 8 |
| $\$\langle\text{expr}\rangle - \underline{\langle\text{factor}\rangle}$ | $* \texttt{id}$ | reduce 7 |
| $\$\langle\text{expr}\rangle - \langle\text{term}\rangle$ | $* \texttt{id}$ | shift |
| $\$\langle\text{expr}\rangle - \langle\text{term}\rangle *$ | $\texttt{id}$ | shift |
| $\$\langle\text{expr}\rangle - \langle\text{term}\rangle * \underline{\texttt{id}}$ | | reduce 9 |
| $\$\langle\text{expr}\rangle - \underline{\langle\text{term}\rangle * \langle\text{factor}\rangle}$ | | reduce 5 |
| $\$\underline{\langle\text{expr}\rangle - \langle\text{term}\rangle}$ | | reduce 3 |
| $\$\underline{\langle\text{expr}\rangle}$ | | reduce 1 |
| $\$\langle\text{goal}\rangle$ | | 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

The key problem: to recognize handles (not covered in this course).

# 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 $\gamma_i$.

# LR($k$) grammars

Formally, a grammar $G$ is LR($k$) iff.:

1. $S \Rightarrow^*_{\text{rm}} \alpha A w \Rightarrow_{\text{rm}} \alpha \beta w$, and
2. $S \Rightarrow^*_{\text{rm}} \gamma B x \Rightarrow_{\text{rm}} \alpha \beta y$, and
3. $\text{FIRST}_k(w) = \text{FIRST}_k(y)$

$\Rightarrow \alpha A y = \gamma B x$

i.e., Assume sentential forms $\alpha\beta w$ and $\alpha\beta y$, with common prefix $\alpha\beta$ and common k-symbol lookahead $\text{FIRST}_k(y) = \text{FIRST}_k(w)$, such that $\alpha\beta w$ reduces to $\alpha A w$ and $\alpha\beta y$ reduces to $\gamma B x$.

But, the common prefix means $\alpha\beta y$ also reduces to $\alpha A y$, for the same result.

Thus $\alpha A y = \gamma B x$.

# 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

  **LL**$(k)$**:** recognize use of a production $A \rightarrow \beta$ seeing first $k$ symbols of $\beta$

  **LR**$(k)$**:** recognize occurrence of $\beta$ (the handle) having seen all of what is derived from $\beta$ plus $k$ symbols of lookahead

# Left versus right recursion

Right Recursion:

- needed for termination in predictive parsers

- requires more stack space

- right associative operators

Left Recursion:

- works fine in bottom-up parsers

- limits required stack space

- left associative operators

Rule of thumb:

- right recursion for top-down parsers

- left recursion for bottom-up 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.

The dilemmas:

- LL dilemma: pick $A \rightarrow b$ or $A \rightarrow c$ ?

- LR dilemma: pick $A \rightarrow b$ or $B \rightarrow b$ ?