

# COM SCI 132 Week 7

Aidan Jan

May 15, 2024

## LR Parsing

### Review

Recall that

- 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

The 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:

```
1 | S -> aABe
2 | A -> Abc
3 |   | b
4 | B -> d
```

and the input string `abbcd`.

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

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

# Handles

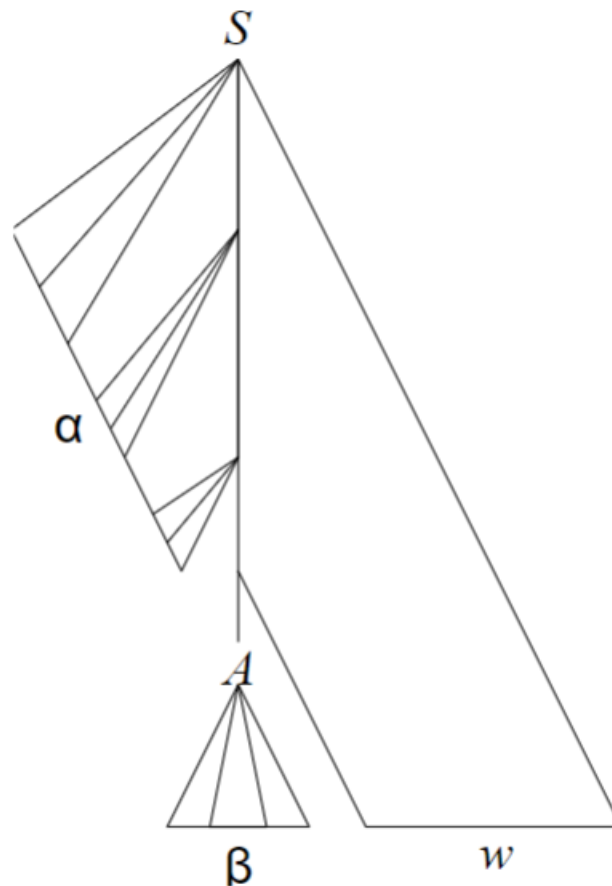
What are we trying to find?

- A substring  $\alpha$  of the tree's upper frontier that matches some production  $A \rightarrow \alpha$  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  $S \Rightarrow_{rm}^* \alpha A w \Rightarrow_{rm} \alpha \beta w$ , 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.



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

## 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

			Prod'n.	Sentential Form
1	$\langle \text{goal} \rangle$	$::=$	$\langle \text{expr} \rangle$	$\langle \text{goal} \rangle$
2	$\langle \text{expr} \rangle$	$::=$	$\langle \text{expr} \rangle + \langle \text{term} \rangle$	$\langle \text{expr} \rangle$
3		$ $	$\langle \text{expr} \rangle - \langle \text{term} \rangle$	$\langle \text{expr} \rangle - \langle \text{term} \rangle$
4		$ $	$\langle \text{term} \rangle$	$\langle \text{expr} \rangle - \langle \text{term} \rangle * \langle \text{factor} \rangle$
5	$\langle \text{term} \rangle$	$::=$	$\langle \text{term} \rangle * \langle \text{factor} \rangle$	$\langle \text{expr} \rangle - \langle \text{term} \rangle * \underline{\text{id}}$
6		$ $	$\langle \text{term} \rangle / \langle \text{factor} \rangle$	$\langle \text{expr} \rangle - \langle \text{factor} \rangle * \text{id}$
7		$ $	$\langle \text{factor} \rangle$	$\langle \text{expr} \rangle - \underline{\text{num}} * \text{id}$
8	$\langle \text{factor} \rangle$	$::=$	$\text{num}$	$\langle \text{term} \rangle - \text{num} * \text{id}$
9		$ $	$\text{id}$	$\langle \text{factor} \rangle - \text{num} * \text{id}$
				$\underline{\text{id}} - \text{num} * \text{id}$

## 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
    find the handle  $A_i \rightarrow \beta_i$  in  $\gamma_i$ 
    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

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

1. Shift until top of stack is the right end of a handle
2. Find the left end of the handle and reduce

For this example: 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 the 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

$$S = \gamma_0 \Rightarrow \gamma_1 \Rightarrow \gamma_2 \Rightarrow \dots \Rightarrow \gamma_n = w$$

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

1. isolate the handle of each right-sequential 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  $\gamma_i$ .  
Formally, a grammar  $G$  is LR( $k$ ) if and only if:

1.  $S \Rightarrow_{rm}^* \alpha Aw \Rightarrow_{rm} \alpha \beta w$ , and
2.  $S \Rightarrow_{rm}^* \gamma Bx \Rightarrow_{rm} \alpha \beta y$ , and
3.  $\text{FIRST}_k(w) = \text{FIRST}_k(y) \Rightarrow \alpha Ay = \gamma Bx$

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  $\alpha Aw$  and  $\alpha \beta y$  reduces to  $\gamma Bx$ .

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

Thus  $\alpha Ay = \gamma Bx$

## 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 operatives

Rule of thumb:

- right recursion for top-down parsers
- left recursion for bottom-up parsers

## Parsing Review

- R. 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.
- Dilemmas:
  - LL dilemma: pick  $A \rightarrow b$  or  $A \rightarrow c$ ?
  - LR dilemma: pick  $A \rightarrow b$  or  $B \rightarrow b$ ?

## Type Checking Generics

A *generic* is a type that has to be parameterized. For example, the compiler needs to be told which object is stored in a given List.

```
class List<A> extends Object {
  public <R> R accept(ListVisitor<R, A> f) {
    // ... do stuff
    return this.accept(f);
  }
}

class ListNull<A> extends List<A> {
  public<R> R accept(ListVisitor<R, A> f) {
    return f.visit(this)
  }
}
```

The above program above will in fact, type check.

## Example

Grammar:

(Class)  $L ::= \text{class } C\langle\bar{x}\rangle \text{ extends } N \{ \bar{s}\bar{f}; K \bar{M} \}$   
(Constructor)  $K ::= C(\bar{T}\bar{f}) \{ \text{super}(\bar{f}); \text{this}.\bar{f} = \bar{f} \}$   
(Method)  $M ::= \langle\bar{Y}\rangle U \ m(\bar{U}\bar{x}) \{ \text{return } e; \}$   
(Expression)  $e ::= x \mid e.f \mid e.\langle\bar{v}\rangle m(\bar{e}) \mid \text{new } N(\bar{e})$   
(Type)  $S, T, U, V ::= X \mid N$   
 $N ::= C\langle\bar{T}\rangle$

Type Rules:

$$\Gamma + e : T$$

where

- $\Gamma$  represents the type environment

- $e$  represents an expression
- $T$  represents a type

$$\frac{\Gamma \vdash e : T \quad \mathbf{fields}(T) = \bar{T}\bar{f}}{\Gamma \vdash e.f_i : T_i} \quad (1)$$

$$\frac{\mathbf{fields}(N) = \bar{T}\bar{f} \quad \Gamma + \bar{e} : \bar{S} \quad \bar{S} \leq \bar{T}}{\Gamma \vdash \mathbf{new} \ N(\bar{e}) : N} \quad (2)$$

$$\frac{\Gamma + e_0 : T_0 \quad \mathbf{mtype}(m, T_0) = \langle \bar{Y} \rangle \bar{U} \rightarrow U \quad \bar{S} \leq [\bar{V}/\bar{Y}]\bar{U}}{\Gamma \vdash e_0. \langle \bar{V} \rangle m(\bar{e}) : [\bar{V}/\bar{Y}]U} \quad (3)$$

- **mtype** refers to method type
- $\langle \bar{Y} \rangle \bar{U}$  are the types of formal parameters and formal parameters, respectively.
- Type expression (3) is created via substitution with (1) and (2).
- $\bar{S} \leq \bar{T}$  is to check class extends; S and T are in a class-superclass relationship.
- Substitution Notation:  $[\bar{V}/\bar{X}]$  means the same thing as  $[\bar{X} := \bar{V}]$
- **mtype** is a placeholder function for method type checking

At this point, we have type checked methods and classes. However, we still have to do the following:

- Override methods
- Fields
- Subtyping
- **mtype**

## Subtyping

Recall that our types:

$$(Type) \ S, T, U, V ::= X \mid N \\ N ::= C \langle \bar{T} \rangle$$

Subtyping rules:

$$T \leq T$$

(one type is a subtype of another type)

$$\frac{S \leq T \quad T \leq U}{S \leq U} \quad (4)$$

$$\frac{\mathbf{class} \ C \langle \bar{x} \rangle \mathbf{extends} \ N\{\dots\}}{C \langle \bar{T} \rangle \leq [\bar{T}/\bar{X}]N} \quad (5)$$

- (5) states that if  $C$  is a use of  $N$ , as in,  $N$  is a parent class to  $C$ , then a use of  $T$  in  $C$  (as a parameter of  $C$ ) is a subtype of the use of  $T$  in  $N$ .
  - In simpler terms,  $C$  parameterized with  $T$  would be a subtype of  $N$  parameterized with  $T$ , because  $C$  is a subtype of  $N$ .

## Fields

$$\frac{\text{class } C < \bar{x} > \text{ extends } N\{\bar{S}\bar{f}; \dots\} \quad \text{fields } ([\bar{T}/\bar{x}]N) = \bar{V}\bar{g}}{\text{fields}(C < \bar{T} >) = \bar{U}\bar{g}, [\bar{T}/\bar{X}]\bar{f}} \quad (6)$$

- $\bar{S}\bar{f}$  is one parameter, where  $S$  is the type,  $f$  is the field.

## mtype (Method Typing)

$$\frac{\text{class } C < \bar{X} > \text{ extends } N\{\dots \bar{M}\} \quad (< \bar{Y} > U \ m(\bar{U} \ \bar{x})\{\text{return } e\}) \in \bar{M}}{\text{mtype}(m, C < \bar{T} >) = [\bar{T}/\bar{X}](< \bar{Y} > \bar{U} \rightarrow U)} \quad (7)$$

$$\frac{\text{class } C < \bar{x} > \text{ extends } N\{\dots \bar{M}\} \quad m \notin \bar{M}}{\text{mtype}(m, C < \bar{T} >) = \text{mtype}(m, [\bar{T}/\bar{X}]\bar{N})} \quad (8)$$

- $\bar{M}$  are the methods of the class  $C$ .

## Method Overriding

$$\frac{\text{If } \{\text{mtype}(m, N) = < \bar{z} > \bar{U} \rightarrow U\} \text{ Then } \{\bar{T} = [Y/\bar{Z}]\bar{U} \quad T = [\bar{Y}/\bar{Z}]U\}}{\text{override}(m, N, < \bar{Y} > \bar{T} \rightarrow T)} \quad (9)$$

In the override statement...

- $m$  = method name
- $N$  = superclass of the current class
- $< \bar{Y} > \bar{T} \rightarrow \bar{T}$  = method header