

PARSING

Baishakhi Ray
Fall 2018

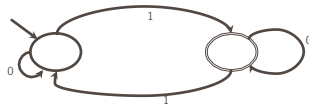
These slides are motivated from Prof. Alex Aiken: Compilers (Stanford)



Languages and Automata

- Formal languages are very important in CS
 - Especially in programming languages
- Regular Languages
 - Weakest formal languages that are widely used
 - Many applications
- Many Languages are not regular

Automata that accepts odd numbers of 1



How many 1s it has accepted?

- Only solution is duplicate state

Automata does not have any memory

Intro to Parsing

- Regular Languages
 - Weakest formal languages that are widely used
 - Many applications
- Consider the language $\{(^i)^i \mid i \geq 0\}$
 - $()()$
 - $(())$
 - $(((+ 2) * 3))$
- Nesting structures
 - if .. if.. else.. else..

Regular languages cannot handle well

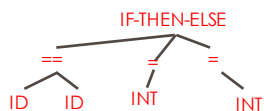
Intro to Parsing

- Input: `if(x==y) 1 else 2;`

- Parser Input (Lexical Input):

`KEY(IF) ' (' ID(x) OP('==') ')' INT(1) KEY(ELSE) INT(2) ';'`

- Parser Output



Intro to Parsing

Character stream

Lexical Analysis

Token stream

Parser

Syntax trees

Semantic Analysis

Syntax trees

Code Generation

- Nor every strings of tokens are valid

- Parser must distinguish between valid and invalid token strings.

- We need

- A Language: to describe valid string
- A method: to distinguish valid from invalid.

Context Free Grammar

- A CFG consists of
 - A set of terminal T
 - A set of non-terminal N
 - A start symbol S ($S \in N$)
 - A set of production rules
 - $X \rightarrow Y_1 \dots Y_n$
 - $X \in N$
 - $Y_i \in \{N, T, \epsilon\}$
- Ex: $S \rightarrow (S) \mid \epsilon$
 - $N = \{S\}$
 - $T = \{ (,), \epsilon \}$

Context Free Grammar

- Begin with a string with only the start symbol S
- Replace a non-terminal X with in the string by the RHS of some production rule: $X \rightarrow Y_1 \dots Y_n$
- Repeat 2 again and again until there are no non-terminals

$X_1 \dots X_i \underline{X} X_{i+1} \dots X_n \rightarrow X_1 \dots X_i Y_1 \dots Y_k X_{i+1} \dots X_n$

For the production rule $X \rightarrow Y_1 \dots Y_k$

$$\alpha_0 \rightarrow \alpha_1 \rightarrow \dots \rightarrow \alpha_n$$

$$\alpha_0 \xrightarrow{*} \alpha_n, n \geq 0$$

Context Free Grammar

- Let G be a CFG with start symbol S . Then the language $L(G)$ of G is:

$$\{a_1 \dots a_n \mid \forall i, a_i \in T \wedge S \xrightarrow{*} a_1 a_2 \dots a_n\}$$

Context Free Grammar

- There are no rules to replace terminals.
- Once generated, terminals are permanent
- Terminals ought to be tokens of programming languages
- Context-free grammars are a natural notation for this recursive structure

CFG: Simple Arithmetic expression

$E \rightarrow E + E$
 $\mid E * E$
 $\mid (E)$
 $\mid id$

Languages can be generated: $id, (id), (id + id)^* id, \dots$

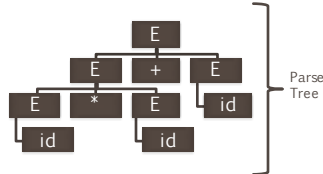
Derivation

- A derivation is a sequence of production
 - $S \rightarrow \dots \rightarrow \dots \rightarrow$
- A derivation can be drawn as a tree
 - Start symbol is tree's root
 - For a production $X \rightarrow Y_1 \dots Y_n$, add children $Y_1 \dots Y_n$ to node X

- Grammar
 - $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- String
 - $id * id + id$

Derivation

$E \rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$
 $\rightarrow id * id + id$



Parse Tree

- A parse tree has
 - Terminals at the leaves
 - Nonterminals at the interior nodes
- An in-order traversal of the leaves is the original input
- The parse tree shows the association of operations, the input string does not

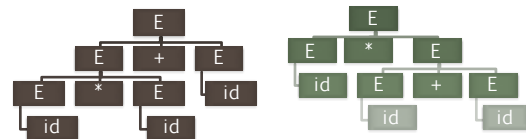
Parse Tree

- Left-most derivation
 - At each step, replace the leftmost non-terminal
 - Right-most derivation
 - At each step, replace the rightmost non-terminal
- $E \rightarrow E + E$
 $\rightarrow E * E + E$
 $\rightarrow id * E + E$
 $\rightarrow id * id + E$
 $\rightarrow id * id + id$
- $E \rightarrow E + E$
 $\rightarrow E + id$
 $\rightarrow E * E + id$
 $\rightarrow E * id + id$
 $\rightarrow id * id + id$

Note that, right-most and left-most derivations have the same parse tree

Ambiguity

- Grammar
 - $E \rightarrow E + E \mid E * E \mid (E) \mid id$
- String
 - $id * id + id$



Ambiguity

- A grammar is ambiguous if it has more than one parse tree for a string
 - There are more than one right-most or left-most derivation for some string
- Ambiguity is bad
 - Leaves meaning for some programs ill-defined

Error Handling

- Purpose of the compiler is
 - To detect nonvalid programs
 - To translate the valid ones
- Many kinds of possible errors (e.g., in C)

Error Kind	Example	Detected by
Lexical	... \$...	Lexer
Syntax	... x*%...	Parser
Semantic	... int x; y = x(3);...	Type Checker
Correctness	your program	tester/user

Error Handling

- Error Handler should
 - Recover errors accurately and quickly
 - Recover from an error quickly
 - Not slow down compilation of valid code
- Types of Error Handling
 - Panic mode
 - Error productions
 - Automatic local or global correction

Panic Mode Error Handling

- Panic mode is simplest and most popular method
- When an error is detected
 - Discard tokens until one with a clear role is found
 - Continue from there
- Typically looks for "synchronizing" tokens
 - Typically the statement or expression terminators

Panic Mode Error Handling

- Example:
 - $(1 + + 2) + 3$
- Panic-mode recovery:
 - Skip ahead to the next integer and then continue
- Bison: use the special terminal **error** to describe how much input to skip
 - $E \rightarrow \text{int} \mid E + E \mid (E) \mid \text{error int} \mid (\text{error})$

Error Productions

- Specify known common mistakes in the grammar
- Example:
 - Write 5x instead of $5 * x$
 - Add production rule $E \rightarrow \dots \mid E E$
- Disadvantages
 - complicates the grammar

Error Corrections

- Idea: find a correct "nearby" program
 - Try token insertions and deletions (goal: minimize edit distance)
 - Exhaustive search
- Disadvantages
 - Hard to implement
 - Slows down parsing of correct programs
 - "Nearby" is not necessarily "the intended" program

Error Corrections

- Past
 - Slow recompilation cycle (even once a day)
 - Find as many errors in once cycle as possible
- Disadvantages
 - Quick recompilation cycle
 - Users tend to correct one error/cycle
 - Complex error recovery is less compelling

Abstract Syntax Trees

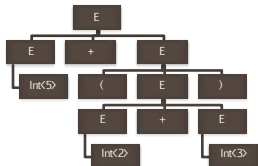
- A parser traces the derivation of a sequence of tokens
- But the rest of the compiler needs a structural representation of the program
- Abstract Syntax Trees
 - Like parse trees but ignore some details
 - Abbreviated as AST

Abstract Syntax Trees

- Grammar
 - $E \rightarrow \text{int} \mid (E) \mid E + E$
- String
 - $5 + (2 + 3)$
- After lexical analysis
 - $\text{Int}\langle 5 \rangle \text{'+' Int}\langle 2 \rangle \text{'+' Int}\langle 3 \rangle \text{'}'}$

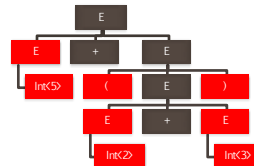
Abstract Syntax Trees: $5 + (2 + 3)$

Parse Trees



Abstract Syntax Trees: $5 + (2 + 3)$

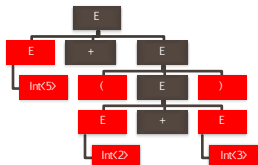
Parse Trees



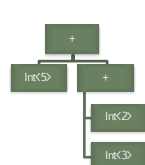
- Have too much information
 - Parentheses
 - Single-successor nodes

Abstract Syntax Trees: $5 + (2 + 3)$

Parse Trees



AST



- Have too much information
 - Parentheses
 - Single-successor nodes
- ASTs capture the nesting structure
 - But abstracts from the concrete syntax
 - More compact and easier to use

Disadvantages of ASTs

- AST has many similar forms
 - E.g., for, while, repeat...until
 - E.g., if, ?:, switch
- Expressions in AST may be complex, nested
 - $(x * y) + (z > 5 ? 12 * z : z + 20)$
- Want simpler representation for analysis
 - ...at least, for dataflow analysis

Parsing algorithm: Recursive Descent Parsing

- The parse tree is constructed
 - From the top
 - From left to right
- Terminals are seen in order of appearance in the token stream

Parsing algorithm: Recursive Descent Parsing

- Grammar:
 - $E \rightarrow T \mid T + E$
 - $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$
- Token Stream: $(\text{int} \langle 5 \rangle)$
- Start with top level non-terminal E
 - Try the rules for E in order

Recursive Descent Parsing Example

$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

E
 $|$
 T
 $|$
 int

mismatch: int does not match arrowhead (
 backtrack

$(\text{int} \langle 5 \rangle)$



Recursive Descent Parsing Example

$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

E
 $|$
 T
 $/ \quad | \quad \backslash$
 $\text{int} \quad * \quad T$

backtrack

$(\text{int} \langle 5 \rangle)$



Recursive Descent Parsing Example

$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

E
 $|$
 T
 $/ \quad | \quad \backslash$
 $(\quad E \quad)$

Match! Advance input

$(\text{int} \langle 5 \rangle)$



Recursive Descent Parsing Example

$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

E
 $|$
 T
 $/ \quad | \quad \backslash$
 $(\quad E \quad)$
 $|$
 T
 $|$
 int

Match! Advance input

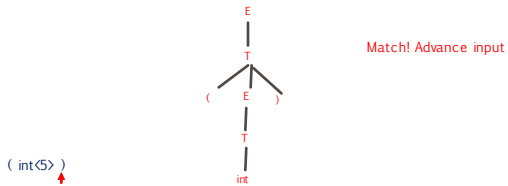
$(\text{int} \langle 5 \rangle)$



Recursive Descent Parsing Example

$E \rightarrow T \mid T + E$

$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$



A Recursive Descent Parser. Preliminaries

- Let **TOKEN** be the type of tokens
 - Special tokens **INT**, **OPEN**, **CLOSE**, **PLUS**, **TIMES**.
- Let the global **next** point to the next token

A (Limited) Recursive Descent Parser

- Define boolean functions that check the token string for a match of
 - A given token terminal


```
bool term (TOKEN tok) { return *next++ == tok; }
```
 - The n^{th} production of S :


```
bool Sn() { ... }
```
 - Try all productions of S :


```
bool S() { ... }
```

A (Limited) Recursive Descent Parser

- For production $E \rightarrow T$

```
bool E1() { return T0(); }
```
- For production $E \rightarrow T + E$

```
bool E2() { return T0() && term(PLUS) && E0(); }
```
- For all productions of E (with backtracking)

```
bool E() {
    TOKEN *save = next;
    return (next = save, E1()) || (next = save, E2());
}
```

A (Limited) Recursive Descent Parser (4)

- Functions for non-terminal T

```
bool T1() { return term(INT); }
bool T2() { return term(INT) && term(TIMES) && T0(); }
bool T3() { return term(OPEN) && E0 && term(CLOSE); }

bool T0() {
    TOKEN *save = next;
    return (next = save, T1())
        || (next = save, T2())
        || (next = save, T3());
}
```

Recursive Descent Parsing

- To start the parser
 - Initialize **next** to point to first token
 - Invoke **E()**. Notice how this simulates the example parse.

Example

Grammar

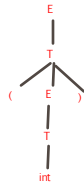
$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

Input: `(int)`

Code: `bool term(TOKEN tok) { return *next++ == tok; }`

```
bool E() { return T(); }
bool E() { return T() && term(PLUS) && E(); }
bool E() { TOKEN *save = next;
  return (next = save, E()) || (next = save, E()); }

bool T() { return term(INT); }
bool T() { return term(INT) && term(TIMES) && T(); }
bool T() { return term(OPEN) && E() && term(CLOSE); }
bool T() { TOKEN *save = next;
  return (next = save, T())
    || (next = save, T())
    || (next = save, T()); }
```



When Recursive Descent Does Not Work

Grammar

$E \rightarrow T \mid T + E$
 $T \rightarrow \text{int} \mid \text{int} * T \mid (E)$

Input: `int * int`

Code: `bool term(TOKEN tok) { return *next++ == tok; }`

```
bool E() { return T(); }
bool E() { return T() && term(PLUS) && E(); }
bool E() { TOKEN *save = next;
  return (next = save, E()) || (next = save, E()); }

bool T() { return term(INT); }
bool T() { return term(INT) && term(TIMES) && T(); }
bool T() { return term(OPEN) && E() && term(CLOSE); }
bool T() { TOKEN *save = next;
  return (next = save, T())
    || (next = save, T())
    || (next = save, T()); }
```

Recursive Descent Parsing: Limitation

- If production for non-terminal X **succeeds**
 - Cannot backtrack to try different production for X later
- General recursive descent algorithms support such full backtracking
 - Can implement any grammar
- Presented RDA is not general
 - But easy to implement
- Sufficient for grammars where for any non-terminal at most one production can succeed
- The grammar can be rewritten to work with the presented algorithm
 - By left factoring

Left Factoring

$A \rightarrow \alpha\beta_1 \mid \alpha\beta_2$

- The input begins with a nonempty string derived from α , we do not know whether to expand A to $\alpha\beta_1$ or $\alpha\beta_2$.
- We can defer the decision by expanding A to $\alpha A'$.
- Then, after seeing the input derived from α , we expand A' to β_1 or β_2 (left-factored)
- The original productions become:

$A \rightarrow \alpha A', A' \rightarrow \beta_1 \mid \beta_2$

When Recursive Descent Does Not Work

- Consider a production $S \rightarrow S a$

```
bool S() { return S() && term(a); }
bool S() { return S(); }
```
- S() goes into an infinite loop
- A **left-recursive grammar** has a non-terminal S

$S \rightarrow \alpha S a$ for some α
- Recursive descent does not work for left recursive grammar

Elimination of Left Recursion

- Consider the left-recursive grammar

$$S \rightarrow S \alpha \mid \beta$$
- S generates all strings starting with a β and followed by a number of α
- Can rewrite using right-recursion

$$S \rightarrow \beta S'$$

$$S' \rightarrow \alpha S' \mid \epsilon$$

More Elimination of Left-Recursion

- In general

$$S \rightarrow S \alpha_1 \mid \dots \mid S \alpha_n \mid \beta_1 \mid \dots \mid \beta_m$$
- All strings derived from S start with one of β_1, \dots, β_m and continue with several instances of $\alpha_1, \dots, \alpha_n$
- Rewrite as

$$S \rightarrow \beta_1 S' \mid \dots \mid \beta_m S'$$

$$S' \rightarrow \alpha_1 S' \mid \dots \mid \alpha_n S' \mid \epsilon$$

General Left Recursion

- The grammar

$$S \rightarrow A \alpha \mid \delta$$

$$A \rightarrow S \beta$$
 is also left-recursive because

$$S \rightarrow^+ S \beta \alpha$$
- This left-recursion can also be eliminated

Summary of Recursive Descent

- Simple and general parsing strategy
 - Left-recursion must be eliminated first
 - ... but that can be done automatically
- Unpopular because of backtracking
 - Thought to be too inefficient
- In practice, backtracking is eliminated by restricting the grammar

Predictive Parsers

- Like recursive-descent but parser can “predict” which production to use
 - By looking at the next few tokens
 - No backtracking
- Predictive parsers accept LL(k) grammars
 - L means “left-to-right” scan of input
 - L means “leftmost derivation”
 - k means “predict based on k tokens of lookahead”
 - In practice, LL(1) is used

LL(1) vs. Recursive Descent

- In recursive-descent
 - At each step, many choices of production to use
 - Backtracking used to undo bad choices
- In LL(1)
 - At each step, only one choice of production
 - That is
 - When a non-terminal A is leftmost in a derivation
 - The next input symbol is t
 - There is a unique production $A \rightarrow \alpha$ to use
 - Or no production to use (an error state)
- LL(1) is a recursive descent variant without backtracking

Predictive Parsing and Left Factoring

- Recall the grammar

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E) \cdot$$
- Hard to predict because
 - For T two productions start with int
 - For E it is not clear how to predict
- We need to left-factor the grammar

Left-Factoring Example

- Grammar

$$E \rightarrow T + E \mid T$$

$$T \rightarrow \text{int} \mid \text{int} * T \mid (E)$$
- Factor out common prefixes of productions

$$E \rightarrow T X$$

$$X \rightarrow + E \mid \epsilon$$

$$T \rightarrow (E) \mid \text{int } Y$$

$$Y \rightarrow * T \mid \epsilon$$

LL(1) Parsing Table Example

- Left-factored grammar

$$E \rightarrow T X$$

$$X \rightarrow + E \mid \epsilon$$

$$T \rightarrow (E) \mid \text{int } Y$$

$$Y \rightarrow * T \mid \epsilon$$

- The LL(1) parsing table:

		next input tokens					
Left-most non- terminals		int	*	+	()	\$
	E	TX			TX		
	X			+E		ϵ	ϵ
	T	int Y			(E)		
	Y		*T	ϵ		ϵ	ϵ

LL(1) Parsing Table Example (Cont.)

- Consider the [E, int] entry
 - "When current non-terminal is E and next input is int, use production $E \rightarrow T X$ "
 - This can generate an int in the first position
- Consider the [Y,+] entry
 - "When current non-terminal is Y and current token is +, get rid of Y"
 - Y can be followed by + only if $Y \rightarrow \epsilon$

LL(1) Parsing Tables. Errors

- Blank entries indicate error situations
- Consider the [E,*) entry
 - "There is no way to derive a string starting with * from non-terminal E"

Using Parsing Tables

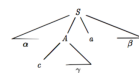
- Method similar to recursive descent, except
 - For the leftmost non-terminal S
 - We look at the next input token a
 - And choose the production shown at [S,a]
- A stack records frontier of parse tree
 - Nonterminals that have yet to be expanded
 - Terminals that have yet to match against the input
 - Top of stack = leftmost pending terminal or non-terminal
- Reject on reaching error state
- Accept on end of input & empty stack

First & Follow

- During top down parsing, FIRST and FOLLOW allow us to choose which production to apply, based on the next input symbol.

- FIRST(α), α is any string of grammar symbols

- A set of terminals that begin strings derived from α .
- If $\alpha \rightarrow \epsilon$, then ϵ is in FIRST(α).
- If $\alpha \rightarrow cY$, the c is in FIRST(α).



- FOLLOW(A), A is a nonterminal

- the set of terminals that can appear immediately to the right of A
- A set of terminals "a" such that $S \rightarrow \alpha A \beta$ for some α and β .

Constructing Parsing Tables: The Intuition

- Consider non-terminal A , production $A \rightarrow \alpha$, & token t
- $T[A, t] = \alpha$ in two cases:
 - If $\alpha \rightarrow^* t \beta$
 - α can derive a t in the first position
 - We say that $t \in \text{First}(\alpha)$
 - If $A \rightarrow \alpha$ and $\alpha \rightarrow^* \epsilon$ and $S \rightarrow^* \beta A t \delta$
 - Useful if stack has A , input is t , and A cannot derive t
 - In this case only option is to get rid of A (by deriving ϵ)
 - We say $t \in \text{Follow}(A)$

Computing First Sets

Definition

$\text{First}(X) = \{ t \mid X \rightarrow^* t\alpha \} \cup \{ \epsilon \mid X \rightarrow^* \epsilon \}$, X can be single terminal, single nonterminal, or string including both

Algorithm sketch:

- $\text{First}(t) = \{ t \}$, t is terminal
- $\epsilon \in \text{First}(X)$
 - if $X \rightarrow \epsilon$
 - if $X \rightarrow A_1 \dots A_n$ and $\epsilon \in \text{First}(A_i)$ for $1 \leq i \leq n$
- $\text{First}(\alpha) \subseteq \text{First}(X)$ if $X \rightarrow A_1 \dots A_n \alpha$
 - $\epsilon \in \text{First}(A_i)$ for $1 \leq i \leq n$

First Sets. Example

grammar

$E \rightarrow T X$
 $X \rightarrow + E \mid \epsilon$
 $T \rightarrow (E) \mid \text{int } Y$
 $Y \rightarrow * T \mid \epsilon$

First sets

$\text{First}(\epsilon) = \{ \epsilon \}$
 $\text{First}(\epsilon) = \{ \epsilon \}$
 $\text{First}(\text{int}) = \{ \text{int} \}$
 $\text{First}(+) = \{ + \}$
 $\text{First}(\epsilon) = \{ \epsilon \}$

$\text{First}(E) \supseteq \text{First}(T) = \{ \text{int}, (\}$
 $\text{First}(X) = \{ +, \epsilon \}$
 $\text{First}(Y) = \{ *, \epsilon \}$

Computing Follow Sets

Definition:

$\text{Follow}(X) = \{ t \mid S \rightarrow^* \beta X t \delta \}$

Intuition:

- If $X \rightarrow A B$ then $\text{First}(B) \subseteq \text{Follow}(A)$ and $\text{Follow}(X) \subseteq \text{Follow}(B)$
- If $B \rightarrow^* \epsilon$ then $\text{Follow}(X) \subseteq \text{Follow}(A)$
- If S is the start symbol then $\$ \in \text{Follow}(S)$

Computing Follow Sets (Cont.)

Algorithm sketch:

- $\$ \in \text{Follow}(S)$
- $\text{First}(\beta) - \{ \epsilon \} \subseteq \text{Follow}(X)$
 - For each production $A \rightarrow \alpha X \beta$
- $\text{Follow}(A) \subseteq \text{Follow}(X)$
 - For each production $A \rightarrow \alpha X \beta$ where $\epsilon \in \text{First}(\beta)$

Follow Sets. Example

Recall the grammar

$E \rightarrow T X$ $X \rightarrow + E \mid \epsilon$
 $T \rightarrow (E) \mid \text{int } Y$ $Y \rightarrow * T \mid \epsilon$

Follow sets

$\text{Follow}(+) = \{ \text{int}, (\}$
 $\text{Follow}(\epsilon) = \{ \text{int}, (\}$
 $\text{Follow}(\epsilon) = \{ \text{int}, (\}$
 $\text{Follow}(\epsilon) = \{ +,) , \$ \}$
 $\text{Follow}(\text{int}) = \{ +, +,) , \$ \}$

$\text{Follow}(E) = \{ \epsilon, \$ \}$
 $\text{Follow}(T) = \{ +,) , \$ \}$
 $\text{Follow}(Y) = \{ +,) , \$ \}$
 $\text{Follow}(X) = \{ \$,) \}$

Constructing LL(1) Parsing Tables

- Construct a parsing table T for CFG G
- For each production $A \rightarrow \alpha$ in G do:
 - For each terminal $t \in \text{First}(\alpha)$ do
 - $T[A, t] = \alpha$
 - If $\epsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do
 - $T[A, t] = \alpha$
 - If $\epsilon \in \text{First}(\alpha)$ and $\$ \in \text{Follow}(A)$ do
 - $T[A, \$] = \alpha$

LL(1) Parsing Table Example

- Left-factored grammar

Rules:
 $E \rightarrow T X$
 $X \rightarrow + E \mid \epsilon$
 $T \rightarrow (E) \mid \text{int } Y$
 $Y \rightarrow * T \mid \epsilon$

For each production $A \rightarrow \alpha$ in G do:
 For each terminal $t \in \text{First}(\alpha)$ do
 $T[A, t] = \alpha$
 If $\epsilon \in \text{First}(\alpha)$, for each $t \in \text{Follow}(A)$ do
 $T[A, t] = \alpha$
 If $\epsilon \in \text{First}(\alpha)$ and $\$ \in \text{Follow}(A)$ do
 $T[A, \$] = \alpha$

- The LL(1) parsing table:

Left-most non-terminals		next input tokens					
	E	int	*	+	()	\$
	X	TX			TX		
	T			+E		ε	ε
	Y	int Y			(E)		
			*T	ε		ε	ε

Notes on LL(1) Parsing Tables

- If any entry is multiply defined then G is not LL(1) [Eg: $S \rightarrow Sa|b$]
 - If G is ambiguous
 - If G is left recursive
 - If G is not leftfactored
 - other e.g., LL(2)
- Most programming language CFGs are not LL(1)
 - too weak
 - However they build on these basic ideas

Bottom-Up Parsing

- Bottom-up parsing is more general than (deterministic) top-down parsing
 - just as efficient
 - Builds on ideas in top-down parsing
- Bottom-up parsers don't need left-factored grammars
- Revert to the "natural" grammar for our example:
 - $E \rightarrow T + E \mid T$
 - $T \rightarrow \text{int } * T \mid \text{int } \mid (E) \cdot$
- Consider the string: $\text{int } * \text{int } + \text{int}$

Bottom-Up Parsing

- Revert to the "natural" grammar for our example:
 - $E \rightarrow T + E \mid T$
 - $T \rightarrow \text{int } * T \mid \text{int } \mid (E) \cdot$
- Consider the string: $\text{int } * \text{int } + \text{int}$
- Bottom-up parsing **reduces** a string to the start symbol by **inverting** productions:

$\text{int } * \text{int } + \text{int}$	$T \rightarrow \text{int}$
$\text{int } * T + \text{int}$	$T \rightarrow \text{int } * T$
$T + \text{int}$	$T \rightarrow \text{int}$
$T + T$	$E \rightarrow T$
$T + E$	$E \rightarrow T + E$
E	

Observation

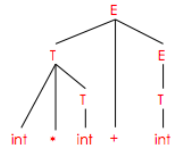
- Read the productions in reverse (from bottom to top)
- This is a rightmost derivation!

$\text{int } * \text{int } + \text{int}$	$T \rightarrow \text{int}$
$\text{int } * T + \text{int}$	$T \rightarrow \text{int } * T$
$T + \text{int}$	$T \rightarrow \text{int}$
$T + T$	$E \rightarrow T$
$T + E$	$E \rightarrow T + E$
E	

Bottom-Up Parsing

- A bottom-up parser traces a **rightmost derivation in reverse**

int * int + int	$T \rightarrow \text{int}$
int * T + int	$T \rightarrow \text{int} * T$
T + int	$T \rightarrow \text{int}$
T + T	$E \rightarrow T$
T + E	$E \rightarrow T + E$
E	



A trivial Bottom-Up Parsing Algorithm

Let I = input string
 repeat
 pick a non-empty substring β of I
 where $X \rightarrow \beta$ is a production
 if no such β , backtrack
 replace one β by X in I
 until I = "S" (the start symbol) or all possibilities are exhausted

Where Do Reductions Happen?

- Right-most derivation has an interesting consequence:
 - Let $\alpha\beta\omega$ be a step of a bottom-up parse
 - Assume the next reduction is by $X \rightarrow \beta$
 - Then ω is a string of terminals
- Why? Because $\alpha\beta\omega \rightarrow \alpha\omega$ is a step in a rightmost derivation

Notation

- Idea: Split string into two substrings
 - Right substring is as yet unexamined by parsing (a string of terminals)
 - Left substring has terminals and non-terminals
- The dividing point is marked by a $|$
 - The $|$ is not part of the string
- Initially, all input is unexamined $| x_1x_2 \dots x_n$

Shift-Reduce Parsing

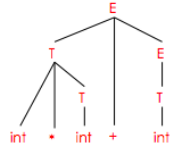
- Bottom-up parsing uses only two kinds of actions:
 - Shift
 - Reduce
- Shift: Move $|$ one place to the right
 - Shifts a terminal to the left string $ABC|xyz \Rightarrow ABCx|yz$
- Reduce: Apply an inverse production at the right end of the left string
 - If $A \rightarrow xy$ is a production, then $Cbxy|ijk \Rightarrow CbA|ijk$

The Example with Reductions Only

int * int + int	reduce $T \rightarrow \text{int}$
int * T + int	reduce $T \rightarrow \text{int} * T$
T + int	reduce $T \rightarrow \text{int}$
T + T	reduce $E \rightarrow T$
T + E	reduce $E \rightarrow T + E$
E	

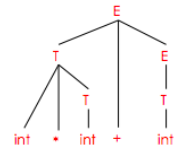
The Example with Shift-Reduce Parsing

int * int + int	shift
int * int + int	shift
int * int + int	shift
int * int + int	reduce $T \rightarrow \text{int}$
int * T + int	reduce $T \rightarrow \text{int} * T$
T + int	shift
T + int	shift
T + int	reduce $T \rightarrow \text{int}$
T + T	reduce $E \rightarrow T$
T + E	reduce $E \rightarrow T + E$
E	



A shift-reduce parsing

int * int + int
int * int + int
int * int + int
int * int + int
int * T + int
T + int
T + int
T + int
T + T
T + E
E



The Stack

- Left string can be implemented by a stack
 - Top of the stack is the |
- Shift pushes a terminal on the stack
- Reduce pops 0 or more symbols off of the stack (production rhs) and pushes a nonterminal on the stack (production lhs)

Conflicts

- 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
- If it is legal to reduce by two different productions, there is a reduce-reduce conflict.