Compilation 0368-3133

Lecture 3a:

Syntax Analysis:

Top-Down parsing

Noam Rinetzky

1

The Real Anatomy of a Compiler



Frontend: Scanning & Parsing





Reminder

- Context-free languages
 - Grammars
 - Pushdown Automata
- Terminology
 - Derivation
 - Sentential form
 - Parse trees
 - Leftmost/rightmost derivation
 - Ambiguous grammars
- Top-down / Botto-up parsers

Parsing

- Construct a structured representation of the input
- Challenges
 - How do you describe the programming language?
 - How do you check validity of an input?
 - Is a sequence of tokens a valid program in the language?
 - How do you construct the structured representation?
 - Where do you report an error?

Top-down parsing



Predictive parsing

- Given a grammar G and a word w attempt to derive w using G
- Idea
 - Apply production to leftmost nonterminal
 - Pick production rule based on next input token
- General grammar
 - More than one option for choosing the next production based on a token
- Restricted grammars (LL)
 - Know exactly which single rule to apply
 - May require some lookahead to decide

Boolean expressions example

 $E \rightarrow LIT \mid (E \text{ OP } E) \mid not E$ LIT \rightarrow true \mid false OP \rightarrow and \mid or \mid xor

not (not true or false)

Boolean expressions example

 $E \rightarrow LIT \mid (E \text{ OP } E) \mid \text{not } E$ LIT \rightarrow true \mid false OP \rightarrow and \mid or \mid xor



Recursive descent parsing

- Define a function for every nonterminal
- Every function work as follows
 - Find applicable production rule
 - Terminal function checks match with next input token
 - Nonterminal function calls (recursively) other functions
- If there are several applicable productions for a nonterminal, use lookahead



Recursive descent

```
void A() {
  choose an A-production, A \rightarrow X_1X_2...X_k;
  for (i=1; i ≤ k; i++) {
    if (X<sub>i</sub> is a nonterminal)
        call procedure X<sub>i</sub>();
    elseif (X<sub>i</sub> == current)
        advance input;
    else
        report error;
  }
}
```

- How do you pick the right A-production?
- Generally try them all and use backtracking
- In our case use lookahead

Problem 1: productions with common prefix

term \rightarrow ID | indexed_elem indexed_elem \rightarrow ID [expr]

- The function for indexed_elem will never be tried...
 - What happens for input of the form ID[expr]

Problem 2: null productions

 $\begin{array}{c} \mathsf{S} \to \mathsf{A} \texttt{ a } \mathsf{ b} \\ \mathsf{A} \to \mathsf{a} \mid \epsilon \end{array}$

```
S() {
   return A() ; match(token('a')) ; match(token('b'))
}
A() {
   match(token('a')) || skip
}
```

- What happens for input "ab"?
- What happens if you flip order of alternatives and try "aab"?



Problem 3: left recursion

```
E \rightarrow E - term | term
```

```
E() {
  return E() ; match(token('-')) ; term()
  ||
  term()
}
```

- What happens with this procedure?
- Recursive descent parsers cannot handle left-recursive grammars

What can we do?

$X \rightarrow YY \mid Z Z \mid Y Z \mid 1 Y$ $Y \rightarrow 4 \mid E$ $Z \rightarrow 2$ $L(Z) = \{2\}$ $L(Y) = \{4, E\}$ $L(X) = \{44, 4, \epsilon, 22, 42, 2, 14, 1\}$

$X \rightarrow YY \mid Z Z \mid Y Z \mid 1 Y$ $Y \rightarrow 4 \mid E$ $Z \rightarrow 2$ $L(Z) = \{2\}$ $L(Y) = \{4, E\}$ $L(X) = \{44, 4, \epsilon, 22, 42, 2, 14, 1\}$

- FIRST(X) = { t | X \rightarrow * t β } \cup { \mathcal{E} | X \rightarrow * \mathcal{E} }
 - FIRST(X) = all terminals that α can appear as first in some derivation for X
 - + E if can be derived from X

- Example:
 - FIRST(LIT) = { true, false }
 - FIRST((E OP E)) = { (}
 - FIRST(not E) = { not }

- No intersection between FIRST sets => can always pick a single rule
- If the FIRST sets intersect, may need longer lookahead
 - LL(k) = class of grammars in which production rule can be determined using a lookahead of k tokens
 - LL(1) is an important and useful class

Computing FIRST sets

• FIRST (t) = { t } // "t" non terminal

• $\mathcal{E} \in FIRST(X)$ if

$$-X \rightarrow \varepsilon$$
 or

 $- X \rightarrow A_1 ... A_k$ and $\mathcal{E} \in FIRST(A_i) i=1...k$

- $FIRST(\alpha) \subseteq FIRST(X)$ if
 - $X \rightarrow A_1 ... A_k \alpha$ and $\mathcal{E} \in FIRST(A_i) i=1...k$

Computing FIRST sets

- Assume no null productions $A \rightarrow \epsilon$
 - 1. Initially, for all nonterminals A, set FIRST(A) = { t | $A \rightarrow t\omega$ for some ω }
 - 2. Repeat the following until no changes occur: for each nonterminal A for each production $A \rightarrow B\omega$ set FIRST(A) = FIRST(A) \cup FIRST(B)
- This is known as fixed-point computation

FIRST sets computation example

```
STMT \rightarrow if EXPR then STMT

| while EXPR do STMT

| EXPR ;

EXPR \rightarrow TERM -> id

| zero? TERM

| not EXPR

| ++ id

| -- id

TERM \rightarrow id

| constant
```

STMT	EXPR	TERM

1. Initialization

```
STMT \rightarrow if EXPR then STMT

| while EXPR do STMT

| EXPR ;

EXPR \rightarrow TERM -> id

| zero? TERM

| not EXPR

| ++ id

| -- id

TERM \rightarrow id

| constant
```

STMT	EXPR	TERM
if while	zero? Not	id constant
	++	

2. Iterate 1

```
\begin{array}{l} \text{STMT} \rightarrow \text{if EXPR then STMT} \\ | \ \text{while EXPR do STMT} \\ | \ \text{EXPR } \text{id} \\ | \ \text{EXPR} ; \\ \text{EXPR} \rightarrow \text{TERM} -> \text{id} \\ | \ \text{zero? TERM} \\ | \ \text{not EXPR} \\ | \ \text{++} \ \text{id} \\ | \ \text{--} \ \text{id} \\ \\ \text{TERM} \rightarrow \text{id} \\ | \ \text{constant} \end{array}
```

STMT	EXPR	TERM
if while	zero? Not ++ 	id constant
zero? Not ++ 		

2. Iterate 2

```
STMT \rightarrow if EXPR then STMT

| while EXPR do STMT

| EXPR ;

EXPR \rightarrow TERM -> id

| zero? TERM

| not EXPR

| ++ id

| -- id

TERM \rightarrow id

| constant
```

STMT	EXPR	TERM
if while	zero? Not ++ 	id constant
zero? Not ++ 	id constant	

2. Iterate 3 – fixed-point

```
STMT → if EXPR then STMT

| while EXPR do STMT

| EXPR ;

EXPR → TERM -> id

| zero? TERM

| not EXPR

| ++ id

| -- id

TERM → id

| constant
```

EXPR	TERM
zero?	id
Not	constant
++	
id	
constant	
	EXPR zero? Not ++ id constant



FOLLOW sets

- What do we do with nullable (ε) productions?
 - $A \rightarrow B C D \quad B \rightarrow \varepsilon C \rightarrow \varepsilon$
 - Use what comes afterwards to predict the right production
- For every production rule $A \to \alpha$
 - FOLLOW(A) = set of tokens that can immediately follow A
- Can predict the alternative A_k for a non-terminal N when the lookahead token is in the set

- FIRST(A_k) \rightarrow (if A_k is nullable then FOLLOW(N))

FOLLOW sets: Constraints

• $\$ \in FOLLOW(S)$

- FIRST(β) { \mathcal{E} } \subseteq FOLLOW(X) – For each A $\rightarrow \alpha \times \beta$
- FOLLOW(A) \subseteq FOLLOW(X) - For each A $\rightarrow \alpha \times \beta$ and $\mathcal{E} \in$ FIRST(β)

Example: FOLLOW sets

- $E \rightarrow TX$ $X \rightarrow + E \mid E$
- $T \rightarrow (E) \mid int Y$ $Y \rightarrow * T \mid E$

Terminal	+	(*)	int
FOLLOW	int, (int, (int, (_,), \$	*,), +, \$

Non. Term.	E	Τ	X	Υ
FOLLOW), \$	+,), \$	\$,)	_,), \$

Prediction Table

• $A \rightarrow \alpha$

- $T[A,t] = \alpha$ if $t \in FIRST(\alpha)$
- $T[A,t] = \alpha$ if $\mathcal{E} \in FIRST(\alpha)$ and $t \in FOLLOW(A)$ - t can also be \$
- T is not well defined → the grammar is not LL(1)

LL(k) grammars

- A grammar is in the class LL(K) when it can be derived via:
 - Top-down derivation
 - Scanning the input from left to right (L)
 - Producing the leftmost derivation (L)
 - With lookahead of k tokens (k)
- A language is said to be LL(k) when it has an LL(k) grammar

LL(1) grammars

- A grammar is in the class LL(1) iff
 - For every two productions $A \rightarrow \alpha$ and $A \rightarrow \beta$ we have
 - FIRST(α) \cap FIRST(β) = {} // including ε
 - If $\varepsilon \in FIRST(\alpha)$ then $FIRST(\beta) \cap FOLLOW(A) = \{\}$
 - If $\varepsilon \in FIRST(\beta)$ then $FIRST(\alpha) \cap FOLLOW(A) = \{\}$

Problem: Non LL Grammars

Problem: Non LL Grammars

```
\begin{array}{c} \mathsf{S} \to \mathsf{A} \texttt{ a } \mathsf{ b} \\ \mathsf{A} \to \mathsf{a} \mid \epsilon \end{array}
```

```
bool S() {
  return A() && match(token('a')) && match(token('b'));
}
```

```
bool A() {
  return match(token('a')) || true;
}
```

- What happens for input "ab"?
- What happens if you flip order of alternatives and try "aab"?
Problem: Non LL Grammars

$$S \rightarrow A a b$$
$$A \rightarrow a \mid \varepsilon$$

• FIRST/FOLLOW conflict

Back to problem 1

```
term \rightarrow ID | indexed_elem
indexed_elem \rightarrow ID [ expr ]
```

- FIRST(term) = { ID }
- FIRST(indexed_elem) = { ID }

• FIRST/FIRST conflict

Solution: left factoring

• Rewrite the grammar to be in LL(1)

term \rightarrow ID | indexed_elem indexed_elem \rightarrow ID [expr]

➡

term \rightarrow ID after_ID After_ID \rightarrow [expr] | ϵ

Intuition: just like factoring x*y + x*z into x*(y+z)

Left factoring – another example





- FIRST(S) = { a } FOLLOW(S) = { }
- FIRST(A) = { a , ε } FOLLOW(A) = { a }

• FIRST/FOLLOW conflict

Solution: substitution



Back to problem 3

 $E \rightarrow E$ - term | term

 Left recursion cannot be handled with a bounded lookahead

• What can we do?





• For our 3rd example:



LL(k) Parsers

- Recursive Descent
 - Manual construction
 - Uses recursion

- Wanted
 - A parser that can be generated automatically
 - Does not use recursion

Pushdown Automata (PDA)



Intuition: PDA

 An ε-NFA with the additional power to manipulate one stack





Intuition: PDA

- Think of an ε-NFA with the additional power that it can manipulate a stack
- PDA moves are determined by:
 - The current state (of its "ε-NFA")
 - The current input symbol (or ε)
 - The current symbol on top of its stack







Intuition: PDA

- Moves:
 - Change state
 - Replace the top symbol by 0...n symbols
 - 0 symbols = "pop" ("reduce")
 - 0 < symbols = sequence of "pushes" ("shift")
- Nondeterministic choice of next move



PDA Formalism

Non terminals

- PDA = (Q, Σ , Γ , δ , q_0 , \$, F): **Tokens**
 - Q: finite set of states
 - Σ: Input symbols alphabet
 - Γ: stack symbols alphabet
 - $-\delta$: transition function
 - $-q_0$: start state
 - \$: start symbol
 - F: set of final states



The Transition Function

- $\delta(q, a, X) = \{ (p_1, \sigma_1), \dots, (p_n, \sigma_n) \}$
 - Input: triplet
 - A state $q \in Q$
 - An input symbol $a \in \Sigma$ or ε
 - A stack symbol $X \in \Gamma$
 - Output: set of 0 ... k actions of the form (p, σ)
 - A state $p \in Q$
 - σ a sequence $X_1 \cdots X_n \in \Gamma^*$ of stack symbols



Actions of the PDA

- Say $(p, \sigma) \in \delta(q, a, X)$
 - If the PDA is in state q and X is the top symbol
 and a is at the front of the input
 - Then it can
 - Change the state to *p*.
 - Remove *a* from the front of the input
 - (but a may be ε).
 - Replace *X* on the top of the stack by σ.



Example: Deterministic PDA

- Design a PDA to accept $\{0^n1^n \mid n > 1\}$.
- The states:
 - q = We have not seen 1 so far
 - start state
 - p = we have seen at least one 1 and no 0s since
 - f = final state; accept.



Example: Stack Symbols

- \$ = start symbol.
 - Also marks the bottom of the stack,
 - Indicates when we have counted the same number of 1's as 0's.
- X = "counter"
 - used to count the number of 0s we saw



Example: Transitions

- $\delta(q, 0, \$) = \{(q, X\$)\}.$
- $\delta(q, 0, X) = \{(q, XX)\}.$

 These two rules cause one X to be pushed onto the stack for each 0 read from the input.

•
$$\delta(q, 1, X) = \{(p, \epsilon)\}.$$

- When we see a 1, go to state p and pop one X.

•
$$\delta(p, 1, X) = \{(p, \epsilon)\}.$$

Pop one X per 1.

• $\delta(p, \epsilon, \$) = \{(f, \$)\}.$



Accept at bottom.

















0 0 0 1 1 1 p X X X \$















Example: Non Deterministic PDA

A PDA that accepts palindromes
 - L {pp' ∈ Σ* | p'=reverse(p)}



LL(k) parsing via pushdown automata

- Pushdown automaton uses
 - Prediction stack
 - Input stream
 - Transition table
 - nonterminals x tokens -> production alternative
 - Entry indexed by nonterminal N and token t contains the alternative of N that must be predicated when current input starts with t

LL(k) parsing via pushdown automata

- Two possible moves
 - Prediction
 - When top of stack is nonterminal N, pop N, lookup table[N,t]. If table[N,t] is not empty, push table[N,t] on prediction stack, otherwise – syntax error
 - Match
 - When top of prediction stack is a terminal T, must be equal to next input token t. If (t == T), pop T and consume t. If (t ≠ T) syntax error
- Parsing terminates when prediction stack is empty
 - If input is empty at that point, success. Otherwise, syntax error

Example transition table

- (1) $E \rightarrow LIT$
- (2) $E \rightarrow (E OP E)$
- (3) $E \rightarrow \text{not } E$
- (4) LIT \rightarrow true
- (5) LIT \rightarrow false
- (6) $OP \rightarrow and$
- (7) $OP \rightarrow or$

Nonterminals

(8) $OP \rightarrow xor$

Which rule should be used

Input tokens \$ false not true and or xor 2 Ε 3 1 1 5 LIT 4 OP 7 8 6

Model of non-recursive predictive parser



Running parser example

aacbb\$

Input suffix	Stack content	Move	
aacbb\$	A\$	predict(A,a) = $A \rightarrow aAb$	
aacbb\$	aAb\$	match(a,a)	
acbb\$	Ab\$	predict(A,a) = $A \rightarrow aAb$	
acbb\$	aAbb\$	match(a,a)	
cbb\$	Abb\$	predict(A,c) = $A \rightarrow c$	
cbb\$	cbb\$	match(c,c)	
bb\$	bb\$	match(b,b)	
b\$	b\$	match(b,b)	
\$	\$	match(\$,\$) – success	

	а	b	С
А	$A \rightarrow aAb$		$A \rightarrow c$

Erorrs

Handling Syntax Errors

- Report and locate the error
- Diagnose the error
- Correct the error
- Recover from the error in order to discover more errors

without reporting too many "strange" errors
Error Diagnosis

- Line number
 - may be far from the actual error
- The current token
- The expected tokens
- Parser configuration

Error Recovery

- Becomes less important in interactive environments
- Example heuristics:
 - Search for a semi-column and ignore the statement
 - Try to "replace" tokens for common errors
 - Refrain from reporting 3 subsequent errors
- Globally optimal solutions
 - For every input w, find a valid program w' with a "minimal-distance" from w

Illegal input example

abcbb\$

$$A \rightarrow aAb \mid c$$

Input suffix	Stack content	Move
abcbb\$	A\$	predict(A,a) = $A \rightarrow aAb$
abcbb\$	aAb\$	match(a,a)
bcbb\$	Ab\$	predict(A,b) = ERROR

	а	b	С
А	$A \rightarrow aAb$		$A \rightarrow c$

Error handling in LL parsers

 $S \rightarrow a c \mid b S$

Input suffix	Stack content	Move
c\$	S\$	predict(S,c) = ERROR

• Now what?

c\$

Predict b S anyway "missing token b inserted in line XXX"

	а	b	С
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error handling in LL parsers

c\$

$S \rightarrow a$ (c b S
---------------------	---------

Input suffix	Stack content	Move
bc\$	S\$	predict(b,c) = $S \rightarrow bS$
bc\$	bS\$	match(b,b)
c\$	S\$	Looks familiar?

• Result: infinite loop

	а	b	С
S	$S \rightarrow a c$	$S \rightarrow b S$	

Error handling and recovery

- x = a * (p+q * (-b * (r-s);
 - Where should we report the error?
 - The valid prefix property

The Valid Prefix Property

- For every prefix tokens
 - $t_1, t_2, ..., t_i$ that the parser identifies as legal:
 - there exists tokens t_{i+1}, t_{i+2}, ..., t_n such that t₁, t₂, ..., t_n is a syntactically valid program
- If every token is considered as single character:
 - For every prefix word u that the parser identifies as legal there exists w such that u.w is a valid program

Recovery is tricky

• Heuristics for dropping tokens, skipping to semicolon, etc.

Building the Parse Tree

Adding semantic actions

- Can add an action to perform on each production rule
- Can build the parse tree
 - Every function returns an object of type Node
 - Every Node maintains a list of children
 - Function calls can add new children

Building the parse tree

```
Node E() {
  result = new Node();
  result.name = "E";
  if (current \in {TRUE, FALSE}) // E \rightarrow LIT
    result.addChild(LIT());
  else if (current == LPAREN) // E \rightarrow ( E OP E )
    result.addChild(match(LPAREN));
    result.addChild(E());
    result.addChild(OP());
    result.addChild(E());
    result.addChild(match(RPAREN));
  else if (current == NOT) // E \rightarrow not E
    result.addChild(match(NOT));
    result.addChild(E());
  else error;
    return result;
```

Parser for Fully Parenthesized Expers

```
static int Parse Expression(Expression **expr p) {
 Expression *expr = *expr_p = new_expression() ;
  /* try to parse a digit */
 if (Token.class == DIGIT) {
        expr->type='D'; expr->value=Token.repr -'0';
        get next token();
         return 1; }
/* try parse parenthesized expression */
if (Token.class == '(') {
      expr->type='P'; get next token();
      if (!Parse Expression(&expr->left)) Error("missing expression");
      if (!Parse_Operator(&expr->oper)) Error("missing operator");
      if (Token.class != ')') Error("missing )");
      get next token();
      return 1; }
return 0;
```

}

Earley Parsing



Jay Earley, PhD

Earley Parsing

• Invented by Jay Earley [PhD. 1968]

Handles arbitrary context free grammars
 – Can handle ambiguous grammars

- Complexity O(N³) when N = |input|
- Uses dynamic programming
 - Compactly encodes ambiguity

Dynamic programming

- Break a problem P into subproblems P₁...P_k
 - Solve P by combining solutions for $P_1...P_k$
 - Memo-ize (store) solutions to subproblems instead of re-computation
- Bellman-Ford shortest path algorithm
 - Sol(x,y,i) = minimum of
 - Sol(x,y,i-1)
 - Sol(t,y,i-1) + weight(x,t) for edges (x,t)

Earley Parsing

- Dynamic programming implementation of a recursive descent parser
 - S[N+1] Sequence of sets of "Earley states"
 - N = |INPUT|
 - Earley state (item) s is a sentential form + aux info
 - S[i] All parse tree that can be produced (by a RDP) after reading the first i tokens
 - S[i+1] built using S[0] ... S[i]

Earley Parsing

- Parse arbitrary grammars in O(|input|³)
 - O(|input|²) for unambigous grammer
 - Linear for most LR(k) langaues (next lesson)
- Dynamic programming implementation of a recursive descent parser
 - S[N+1] Sequence of sets of "Earley states"
 - N = |INPUT|
 - Earley states is a sentential form + aux info
 - S[i] All parse tree that can be produced (by an RDP) after reading the first i tokens
 - S[i+1] built using S[0] ... S[i]

Earley States

- s = < constituent, back >
 - constituent (dotted rule) for $A \rightarrow \alpha \beta$
 - $A \rightarrow \circ \alpha \beta$ predicated constituents
 - $A \rightarrow \alpha \bullet \beta$ in-progress constituents
 - $A \rightarrow \alpha \beta \bullet$ completed constituents
 - back previous Early state in derivation

Earley States

- s = < constituent, back >
 - constituent (dotted rule) for $A \rightarrow \alpha \beta$
 - $A \rightarrow \circ \alpha \beta$ predicated constituents
 - $A \rightarrow \alpha \bullet \beta$ in-progress constituents
 - $A \rightarrow \alpha \beta \bullet$ completed constituents
 - back previous Early state in derivation

Earley Parser

```
Input = x[1...N]
S[0] = \langle E' \rightarrow \bullet E, 0 \rangle; S[1] = ... S[N] = \{\}
for i = 0 ... N do
    until S[i] does not change do
       foreach s \in S[i]
           if s = \langle A \rightarrow \dots \bullet a \dots, b \rangle and a = x[i+1] then
                                                                                                                // scan
               S[i+1] = S[i+1] \cup \{ < A \rightarrow ..., b > \}
           if s = \langle A \rightarrow ..., b \rangle and X \rightarrow \alpha then
                                                                                                                // predict
               S[i] = S[i] \cup \{\langle X \rightarrow \bullet \alpha, i \rangle \}
           if s = \langle A \rightarrow ... \bullet, b \rangle and \langle X \rightarrow ... \bullet A..., k \rangle \in S[b] then // complete
                S[i] = S[i] \cup \{\langle X \rightarrow ..., k \rangle\}
```

Example



FIGURE 1. Earley sets for the grammar $E \rightarrow E + E \mid n$ and the input n + n. Items in bold are ones which correspond to the input's derivation.