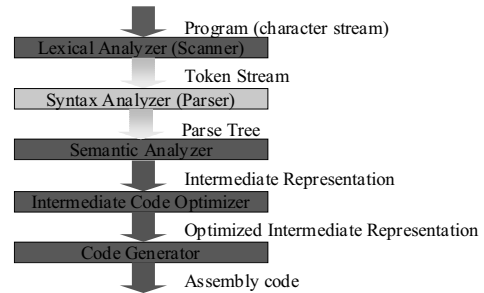


6.035

Fall 2000

Lecture 4: Top Down Parsing

Parsing



Martin Rinard

2

6.035 ©MIT Fall 2000

Last Lecture

- Focus: Specifying languages
- Generative approach
- Context-Free Grammars
 - Ambiguity
 - Abstract versus Concrete Syntax

Martin Rinard

3

6.035 ©MIT Fall 2000

This Lecture

- Determining if a program is syntactically correct
- Building a parse tree from a sequence of tokens
- Top-down, recursive descent approach
- Structure of parsing program matches structure of grammar

Martin Rinard

4

6.035 ©MIT Fall 2000

Basic Approach

- Start with goal symbol
- Build a leftmost derivation
 - If leftmost symbol is nonterminal, choose a production and apply it
 - If leftmost symbol is terminal, match against input
 - If all terminals match, have found a parse!
 - Key: find correct productions for nonterminals

Martin Rinard

5

6.035 ©MIT Fall 2000

Grammar for Parsing Example

$Goal \rightarrow Expr$

$Expr \rightarrow Expr + Term$

$Expr \rightarrow Expr - Term$

$Expr \rightarrow Term$

$Term \rightarrow Term * Num$

$Term \rightarrow Term / Num$

$Term \rightarrow Num$

$Num \rightarrow 0$

$Num \rightarrow 1$

$Num \rightarrow 2$

Martin Rinard

6

6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

graph BT
    GP[Current Position in Parse Tree] --> G[Goal]
        
```

Remaining Input
2-2*2

Sentential Form
Goal

Martin Rinard 7 6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

graph BT
    GP[Current Position in Parse Tree] --> E[Expr]
    G[Goal] --> E
        
```

Remaining Input
2-2*2

Sentential Form
Expr

Applied Production
Goal → *Expr*

Martin Rinard 8 6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

graph TD
    G[Goal] --> E1[Expr]
    G --> M[-]
    G --> T1[Term]
    E1 --> E2[Expr]
        
```

Remaining Input
2-2*2

Sentential Form
Expr - Term

Applied Production
Expr → *Expr - Term*

Martin Rinard 9 6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

graph TD
    G[Goal] --> E1[Expr]
    G --> M[-]
    G --> T1[Term]
    E1 --> E2[Expr]
    E2 --> T2[Term]
        
```

Remaining Input
2-2*2

Sentential Form
Term - Term

Applied Production
Expr → *Term*

Martin Rinard 10 6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

graph TD
    G[Goal] --> E1[Expr]
    G --> M[-]
    G --> T1[Term]
    E1 --> E2[Expr]
    E2 --> T2[Term]
    T2 --> N[Num]
        
```

Remaining Input
2-2*2

Sentential Form
Num - Term

Applied Production
Term → *Num*

Martin Rinard 11 6.035 ©MIT Fall 2000

Parsing Example

Parse Tree

```

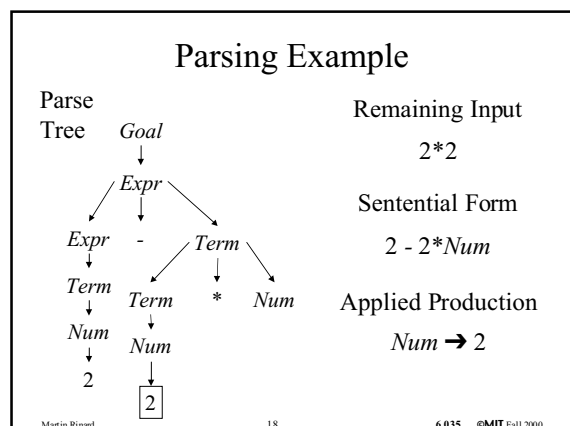
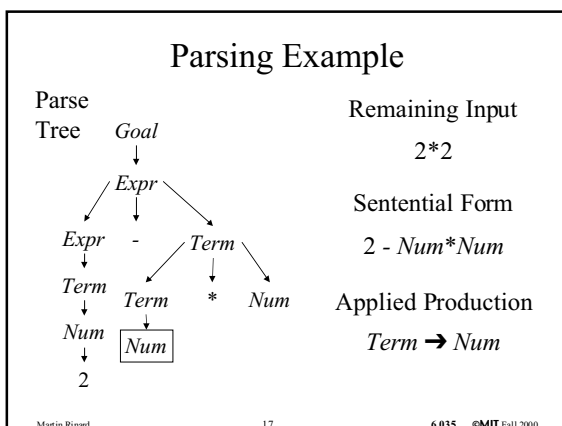
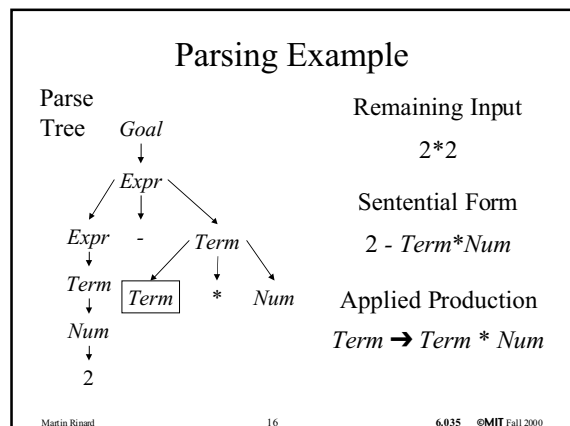
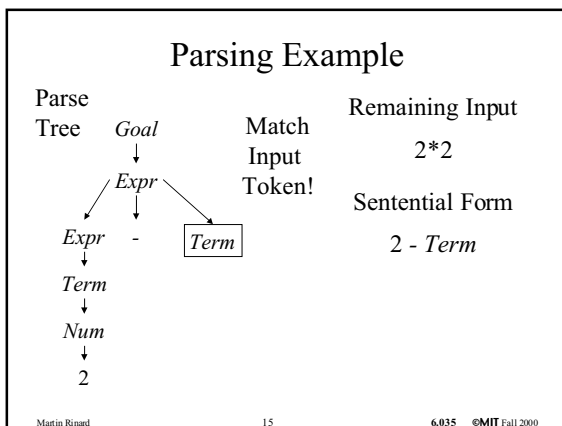
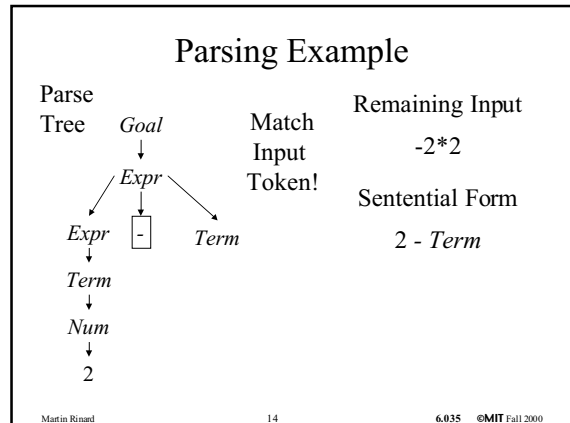
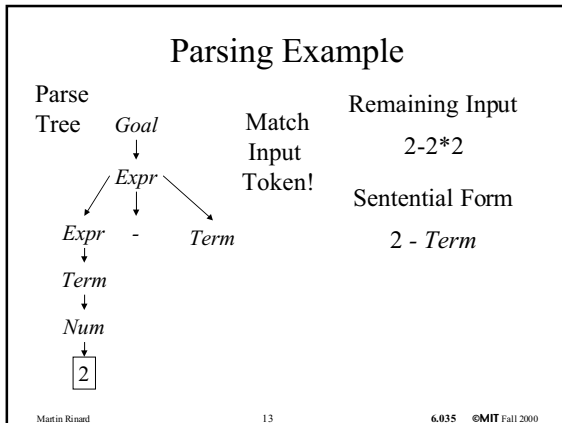
graph TD
    G[Goal] --> E1[Expr]
    G --> M[-]
    G --> T1[Term]
    E1 --> E2[Expr]
    E2 --> T2[Term]
    T2 --> N[Num]
    N --> 2[2]
        
```

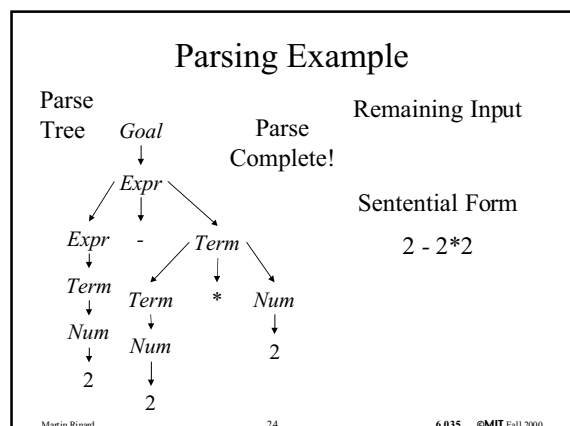
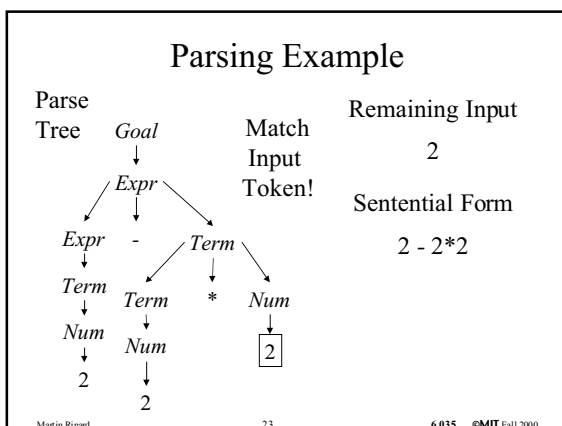
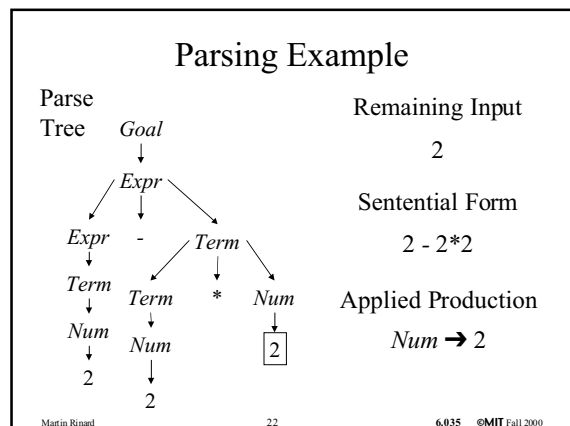
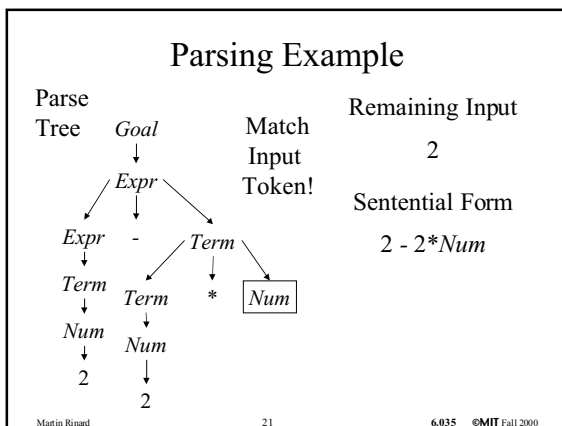
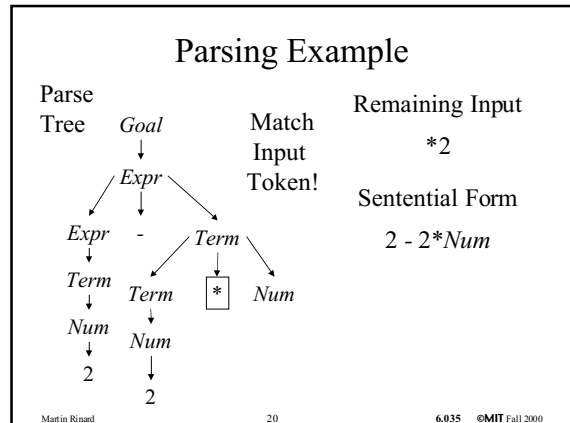
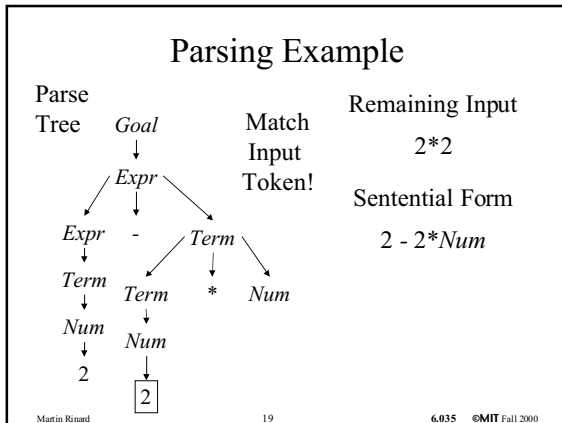
Remaining Input
2-2*2

Sentential Form
2 - Term

Applied Production
Num → *2*

Martin Rinard 12 6.035 ©MIT Fall 2000





Summary

- Three Actions (Mechanisms)
 - Apply production to expand current nonterminal in parse tree
 - Match current terminal (consuming input)
 - Accept the parse as correct
- Parser generates preorder traversal of parse tree
 - visit parents before children
 - visit siblings from left to right

Martin Roud

25

6.035 ©MIT Fall 2000

Policy Problem

- Which production to use for each nonterminal?
- Classical Separation of Policy and Mechanism
- One Approach: Backtracking
 - Treat it as a search problem
 - At each choice point, try next alternative
 - If it is clear that current try fails, go back to previous choice and try something different
- General technique for searching
- Used a lot in classical AI and natural language parsing

Martin Roud

26

6.035 ©MIT Fall 2000

Backtracking Example

Parse Tree *Goal*

Remaining Input
2-2*2

Sentential Form
Goal

Martin Roud

27

6.035 ©MIT Fall 2000

Backtracking Example

Parse Tree
Goal
↓
Expr

Remaining Input
2-2*2

Sentential Form
Expr

Applied Production
Goal → *Expr*

Martin Roud

28

6.035 ©MIT Fall 2000

Backtracking Example

Parse Tree
Goal
↓
Expr
↙ ↓ ↘
Expr + *Term*

Remaining Input
2-2*2

Sentential Form
Expr - *Term*

Applied Production
Expr → *Expr* + *Term*

Martin Roud

29

6.035 ©MIT Fall 2000

Backtracking Example

Parse Tree
Goal
↓
Expr
↙ ↓ ↘
Expr + *Term*
↓
Term

Remaining Input
2-2*2

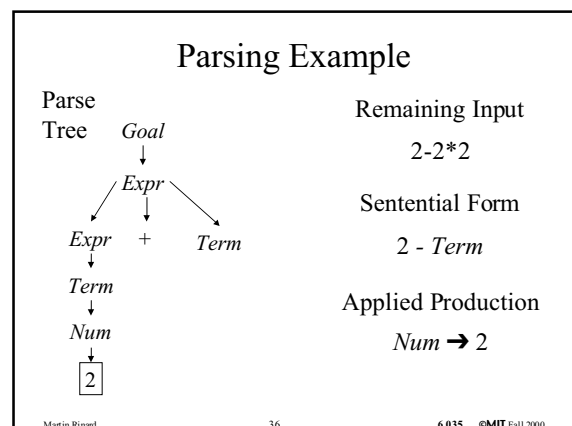
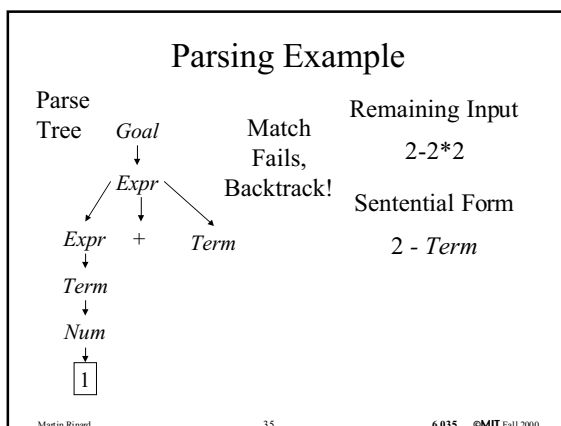
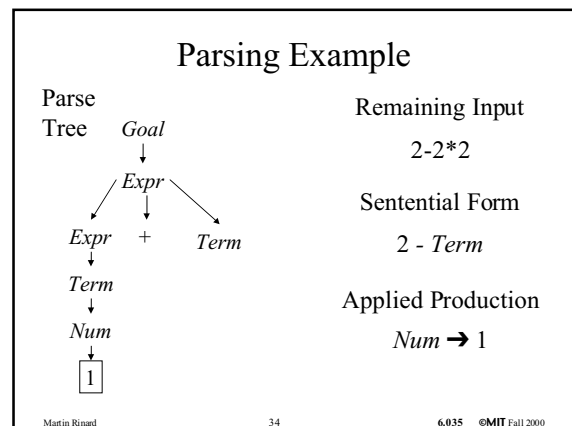
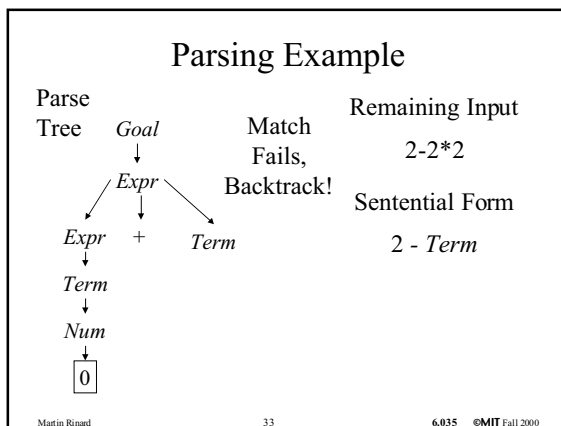
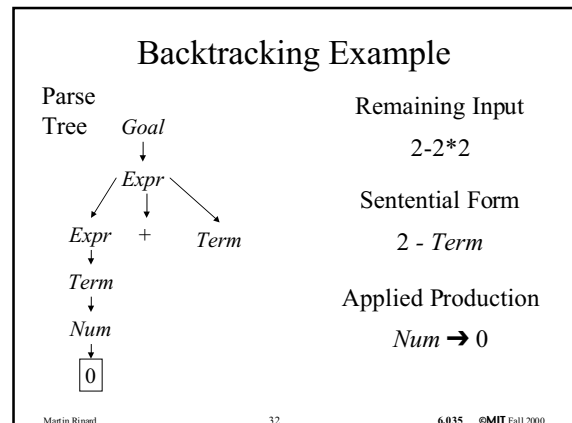
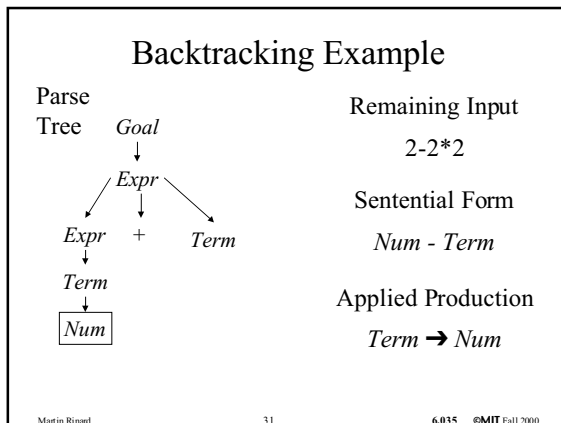
Sentential Form
Term - *Term*

Applied Production
Expr → *Term*

Martin Roud

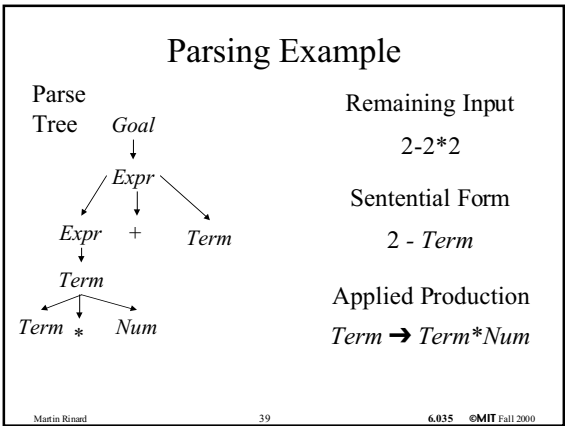
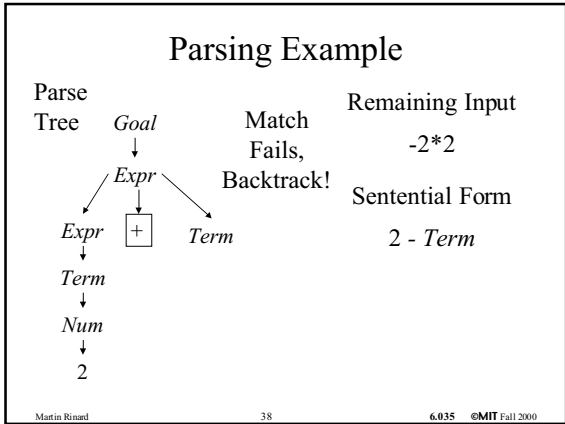
30

6.035 ©MIT Fall 2000



Parsing Example

<p>Parse Tree</p> <pre> Goal ↓ Expr / \ Expr + Term ↓ Term ↓ Num ↓ [2] </pre>	<p>Match Input Token!</p>	<p>Remaining Input</p> <p style="text-align: center;">2-2*2</p> <p>Sentential Form</p> <p style="text-align: center;">2 - Term</p>
--------------------------------------------------------------------------------------------------------------------	---------------------------	------------------------------------------------------------------------------------------------------------------------------------



Left Recursion + Top-Down Parsing = Infinite Loop

- Example Production: $Term \rightarrow Term * Num$
- Potential parsing steps:

The diagram illustrates the potential for infinite loops in top-down parsing due to left recursion. It shows two parse trees for the production $Term \rightarrow Term * Num$.

The left tree shows a single step where $Term$ is expanded to $Term * Num$. The root node is $Term$, which has three children: $Term$, $*$, and Num . The leftmost $Term$ child is boxed.

The right tree shows a recursive expansion where the new $Term$ is further expanded to $Term * Num$, leading to an infinite loop. The root node is $Term$, which has three children: $Term$, $*$, and Num . The leftmost $Term$ child is further expanded into another $Term$, $*$, and Num node, with the new $Term$ child also boxed.

General Problem with Search

- Must ensure that every step makes some sort of meaningful progress
- Handled in various ways in various contexts
- For parsing, hack grammar to remove left recursion

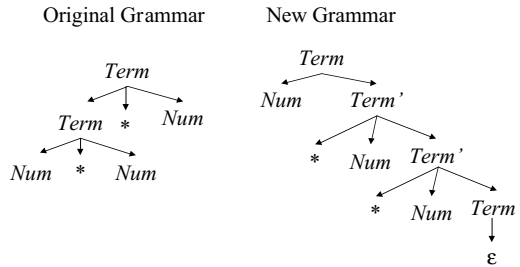
Martin Rinard 41 6.035 ©MIT Fall 2000

Hacked Grammar

Original Grammar Fragment	New Grammar Fragment
$Term \rightarrow Term * Num$	$Term \rightarrow Num Term'$
$Term \rightarrow Term / Num$	$Term' \rightarrow * Num Term'$
$Term \rightarrow Num$	$Term' \rightarrow / Num Term'$
	$Term' \rightarrow \epsilon$

Martin Rinard 42 6.035 ©MIT Fall 2009

Parse Tree Comparisons



Martin Rinaud

43

6.035 ©MIT Fall 2000

Predictive Parsing

- Alternative to backtracking
- Useful for programming languages, which can be designed to make parsing easier
- Basic idea
 - Look ahead in input stream
 - Decide which production to apply based on next tokens in input stream
 - We will use one token of lookahead

Martin Rinaud

44

6.035 ©MIT Fall 2000

Predictive Parsing Example Grammar

$Goal \rightarrow Expr$	$Term \rightarrow Num Term'$
$Expr \rightarrow Term Expr'$	$Term' \rightarrow * Num Term'$
$Expr' \rightarrow + Expr'$	$Term' \rightarrow / Num Term'$
$Expr' \rightarrow - Expr'$	$Term' \rightarrow \epsilon$
$Expr' \rightarrow \epsilon$	$Num \rightarrow 0$
	$Num \rightarrow 1$
	$Num \rightarrow 2$

Martin Rinaud

45

6.035 ©MIT Fall 2000

Parsing Choice Points

- Assume Num is current position in parse tree
 - Must choose one of
 - $Num \rightarrow 0$
 - $Num \rightarrow 1$
 - $Num \rightarrow 2$
- Look at next token and see which one it is
 - If token is 0, choose $Num \rightarrow 0$
 - If token is 1, choose $Num \rightarrow 1$
 - If token is 2, choose $Num \rightarrow 2$

Martin Rinaud

46

6.035 ©MIT Fall 2000

More Choice Points

- Assume $Term'$ is current position in parse tree
 - $Term' \rightarrow * Num Term'$
 - $Term' \rightarrow / Num Term'$
 - $Term' \rightarrow \epsilon$
- If next token is $*$ or $/$, apply appropriate production
- Otherwise, apply $Term' \rightarrow \epsilon$

Martin Rinaud

47

6.035 ©MIT Fall 2000

Nonterminals

- What about productions with nonterminals?
 - $NT \rightarrow NT_1 \alpha_1$
 - $NT \rightarrow NT_2 \alpha_2$
- Must choose based on possible first terminals that NT_1 and NT_2 can generate
- What if NT_1 or NT_2 can generate ϵ ?
 - Must choose based on α_1 and α_2

Martin Rinaud

48

6.035 ©MIT Fall 2000

First(β)

- $T \in \text{First}(\beta)$ if T can appear as the first symbol in a derivation starting from β
- Start with concept of NT deriving ϵ
- Two rules
 - $NT \rightarrow \epsilon$ implies NT derives ϵ
 - $NT \rightarrow NT_1 \dots NT_n$ and for all $1 \leq i \leq n$ NT_i derives ϵ implies NT derives ϵ

Martin Rind

49

6.035 ©MIT Fall 2000

Fixed Point Algorithm for Derives

for all nonterminals NT
 set NT derives ϵ to be false
 for all productions of the form $NT \rightarrow \epsilon$
 set NT derives ϵ to be true
 while (some NT derives ϵ changed in last iteration)
 for all productions of the form $NT \rightarrow NT_1 \dots NT_n$
 if (for all $1 \leq i \leq n$ NT_i derives ϵ)
 set NT derives ϵ to be true

Martin Rind

50

6.035 ©MIT Fall 2000

Rules for First

- 1) $T \in \text{First}(T)$
 - 2) $\text{First}(S) \subseteq \text{First}(S\beta)$
 - 3) NT derives ϵ implies $\text{First}(\beta) \subseteq \text{First}(NT\beta)$
 - 4) $NT \rightarrow S\beta$ implies $\text{First}(S\beta) \subseteq \text{First}(NT)$
- Notation
 - T is a terminal, NT is a nonterminal, S is a terminal or nonterminal, and β is a sequence of terminals or nonterminals

Martin Rind

51

6.035 ©MIT Fall 2000

Rules + Request Generate System of Subset Inclusion Constraints

Grammar
 $Term' \rightarrow * Num Term'$
 $Term' \rightarrow / Num Term'$
 $Term' \rightarrow \epsilon$ $Num \rightarrow 0$
 $Num \rightarrow 1$ $Num \rightarrow 2$

Request: What is $\text{First}(Term')$?

Rules
 1) $T \in \text{First}(T)$
 2) $\text{First}(S) \subseteq \text{First}(S\beta)$
 3) NT derives ϵ implies $\text{First}(\beta) \subseteq \text{First}(NT\beta)$
 4) $NT \rightarrow S\beta$ implies $\text{First}(S\beta) \subseteq \text{First}(NT)$

Constraints
 $\text{First}(* Num Term') \subseteq \text{First}(Term')$
 $\text{First}(/ Num Term') \subseteq \text{First}(Term')$
 $\text{First}(\epsilon) \subseteq \text{First}(Term')$
 $\text{First}(0) \subseteq \text{First}(Num)$
 $\text{First}(1) \subseteq \text{First}(Num)$
 $\text{First}(2) \subseteq \text{First}(Num)$

Martin Rind

52

6.035 ©MIT Fall 2000

Constraint Propagation Algorithm

Constraints	Solution
$\text{First}(* Num Term')$	$\text{First}(Term') = \{\}$
$\subseteq \text{First}(Term')$	$\text{First}(* Num Term') = \{\}$
$\text{First}(/ Num Term')$	$\text{First}(/ Num Term') = \{\}$
$\subseteq \text{First}(Term')$	$\text{First}(\epsilon) = \{\}$
$\text{First}(\epsilon) \subseteq \text{First}(* Num Term')$	$\text{First}(/) = \{\}$
$\text{First}(/) \subseteq \text{First}(/ Num Term')$	
$* \in \text{First}(\epsilon)$	Initialize Sets to $\{\}$
$/ \in \text{First}(/)$	Propagate Constraints Until Fixed Point

Martin Rind

53

6.035 ©MIT Fall 2000

Constraint Propagation Algorithm

Constraints	Solution
$\text{First}(* Num Term')$	$\text{First}(Term') = \{\}$
$\subseteq \text{First}(Term')$	$\text{First}(* Num Term') = \{\}$
$\text{First}(/ Num Term')$	$\text{First}(/ Num Term') = \{\}$
$\subseteq \text{First}(Term')$	$\text{First}(\epsilon) = \{\epsilon\}$
$\text{First}(\epsilon) \subseteq \text{First}(* Num Term')$	$\text{First}(/) = \{/ \}$
$\text{First}(/) \subseteq \text{First}(/ Num Term')$	
$* \in \text{First}(\epsilon)$	
$/ \in \text{First}(/)$	

Martin Rind

54

6.035 ©MIT Fall 2000

Constraint Propagation Algorithm

Constraints	Solution
$\text{First}(* \text{ Num Term }')$	$\text{First}(\text{Term }') = \{\}$
$\subseteq \text{First}(\text{Term }')$	$\text{First}(* \text{ Num Term }') = \{*\}$
$\text{First}(/ \text{ Num Term }')$	$\text{First}(/ \text{ Num Term }') = \{/ \}$
$\subseteq \text{First}(\text{Term }')$	$\text{First}(*') = \{*\}$
$\text{First}(*') \subseteq \text{First}(* \text{ Num Term }')$	$\text{First}(/) = \{/ \}$
$\text{First}(/) \subseteq \text{First}(/ \text{ Num Term }')$	
$* \in \text{First}(*')$	
$/ \in \text{First}(/)$	

Martin Rinard

55

6.035 ©MIT Fall 2000

Constraint Propagation Algorithm

Constraints	Solution
$\text{First}(* \text{ Num Term }')$	$\text{First}(\text{Term }') = \{*, /\}$
$\subseteq \text{First}(\text{Term }')$	$\text{First}(* \text{ Num Term }') = \{*\}$
$\text{First}(/ \text{ Num Term }')$	$\text{First}(/ \text{ Num Term }') = \{/ \}$
$\subseteq \text{First}(\text{Term }')$	$\text{First}(*') = \{*\}$
$\text{First}(*') \subseteq \text{First}(* \text{ Num Term }')$	$\text{First}(/) = \{/ \}$
$\text{First}(/) \subseteq \text{First}(/ \text{ Num Term }')$	
$* \in \text{First}(*')$	
$/ \in \text{First}(/)$	

Martin Rinard

56

6.035 ©MIT Fall 2000

General Idea in Computer Science

- Systems of subset inclusion constraints
- Typically, a problem generates a system
- Outline of solution algorithm
 - Initialize all sets to $\{\}$
 - Use a constraint propagation algorithm that iteratively update sets by satisfying one constraint at a time
 - Terminate when reach a fixed point.

Martin Rinard

57

6.035 ©MIT Fall 2000

Predictive Parsing + Hand Coding = Recursive Descent Parser

- One procedure per nonterminal
- That procedure examines the current input symbol
- Recursively calls procedures for RHS of chosen production
- Procedures return true if parse succeeded, false otherwise

Martin Rinard

58

6.035 ©MIT Fall 2000

Example

```

Term()
  return(Num(), TermPrime())
TermPrime()
  if (token = *) || (token = /)
    token = NextToken();
    return (Num() && TermPrime())
  else return(true)
Num()
  if (token = 1) || (token == 2) || (token == 3)
    token = NextToken(); return(true)
  else return(false)

```

Martin Rinard

59

6.035 ©MIT Fall 2000

Building A Parse Tree

- Have each procedure return the section of the parse tree for the part of the string it parsed
- Use exceptions to make code structure clean

Martin Rinard

60

6.035 ©MIT Fall 2000

Building A Parse Tree in Example

```
Term()
    return(new TermNode(Num(), TermPrime()))
TermPrime()
    if (token = *) || (token = /)
        t = token; token = NextToken();
        return (new TermPrimeNode(t, Num(), TermPrime()));
    else return(new TermPrimeNode( $\epsilon$ ));
Num()
    if (token = 1) || (token == 2) || (token == 3)
        t = token; token = NextToken(); return(new NumNode(t));
    else throw(SyntaxError);
```

Martin Rinard

61

6.035 ©MIT Fall 2000

Summary

- Top-Down Parsing
- Use Lookahead to Avoid Backtracking
- Parser is Hand-Coded Set of Mutually Recursive Procedures
- General Ideas in Computer Science
 - Backtracking as a Search Mechanism
 - Systems of Subset Inclusion Constraints
 - Fixed-Point Solution Algorithms

Martin Rinard

62

6.035 ©MIT Fall 2000