



Principles of Compiler Construction

Prof. Wen-jun LI

School of Computer Science and Engineering

lnslwj@mail.sysu.edu.cn



Lecture 2. Formal Languages

1. **What** Is a Language, Formally?
2. **How** to Define a Specific Language?
3. Context-Free Grammar
4. Parse Tree
5. Ambiguity
6. Chomsky Hierarchy



Formal System

- Don't care the meaning of symbols.
 - Symbols can be substituted.
- Only depend on formally defined rules.
- Programs are formal systems.

1. What Is a Language, Formally?

- Alphabet $::=$ a **set** of **symbols**
 - Non-empty and finite
 - Symbol is a meta-definition
- String (sentence or word) $::=$ a **sequence** of symbols which are from the alphabet
 - Finite; Permit empty string ε [*epsilon*]
 - Manipulation: empty, length, prefix, suffix, substring, proper \sim , subsequence, concatenation, Kleene star (closure)
- Σ^* $::=$ the set of all possible strings on Σ
 - Mathematically, a monoid (semigroup with an identity) generated by Σ .

Formal Definition of Language

- Language ::= A **set** of strings
 - Maybe infinite, or an empty set $\{\}$ (or \emptyset)
- Language manipulation
 - Union: $L_1 \cup L_2$
 - Concatenation: $L_1 L_2$
 - Exponential: $L^0 = \{\epsilon\}$, $L^1 = L$, $L^2 = LL$, ...
 - Kleene closure (star): $L^* = \bigcup_{i=0}^{\infty} L^i$
 - Positive closure (plus): $L^+ = \bigcup_{i=1}^{\infty} L^i = L^* L$
- What is a meta-definition?

Mapping (Binding) to Practice

- Binding
 - Abstract: formal definition of a language
 - Concrete: programming languages
- **Discussion**
 - For C programs
 - Alphabet $\Sigma \rightarrow ?$
 - $\Sigma^* \rightarrow ?$
 - Language $\rightarrow ?$
 - A string in the language $\rightarrow ?$
 - For Java programs
 - Alphabet $\Sigma \rightarrow ?$
 - A string in the language $\rightarrow ?$

2. How to Define a Specific Language?

- Essentially, how to define a set:
 - Enumeration
 - { aa, ab, ba, bb }
 - Partial enumeration
 - { a, ab, abb, abbb, ... }
 - Predicate description
 - { $x \mid \mathbf{P}(x) \wedge x \in \Sigma^*$ }
- Problem: how to define the predicate **P**?
 - Described (by an expression)
 - Generated (by a grammar)
 - Recognized (by an automaton)

3. Context-Free Grammar (CFG)

- CFG is a 4-tuple (Σ, N, P, S)
 - Σ : alphabet, set of terminals
 - $N \cap \Sigma = \emptyset$: set of nonterminals
 - P : set of rewriting rules (productions)
 - Each production has the form $A \rightarrow \alpha$,
where $A \in N \wedge \alpha \in (\Sigma \cup N)^*$
 - $S \in N$: start (goal) symbol

An Example

- A Context-Free Grammar

- $G_1 = (\{a, b, c\}, \{A, B\}, \{A \rightarrow aB, A \rightarrow bB, A \rightarrow cB, B \rightarrow a, B \rightarrow b, B \rightarrow c\}, A)$

- Brief notation

$$A \rightarrow aB \mid bB \mid cB$$
$$B \rightarrow a \mid b \mid c$$

Derivation

- Derive (rewrite)

- $A \Rightarrow bB \Rightarrow bc$
- Direct derivation: $bB \Rightarrow bc$
- n-step derivation: $A \Rightarrow^2 bc$
- 0 or more steps derivation: $A \Rightarrow^* bc$
- 1 or more steps derivation: $A \Rightarrow^+ bc$
- Left-most derivation: $A \Rightarrow_{lm}^* bc$
 - Every step must be left-most derivation
- Right-most derivation: $A \Rightarrow_{rm}^* bc$
 - Every step must be right-most derivation
- Canonical derivation: right-most

Binding to Practice

- Arithmetic expressions

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

$$T \rightarrow T * F \mid F$$

$$F \rightarrow (E) \mid n$$

Each nonterminal
represents a subset
of the language

- Derivation of **$n+n*n$**

- **$E \Rightarrow E+T \Rightarrow E+T*F \Rightarrow T+T*F \Rightarrow F+T*F \Rightarrow F+F*F \Rightarrow F+F*n \Rightarrow F+n*n \Rightarrow n+n*n$**
- Not canonical!

- Canonical derivation of **$n+n*n$**

- **$E \Rightarrow E+T \Rightarrow E+T*F \Rightarrow E+T*n \Rightarrow E+F*n \Rightarrow E+n*n \Rightarrow T+n*n \Rightarrow F+n*n \Rightarrow n+n*n$**



Reduction

○ Reduce

- Derive in reverse order
- Direct reduction
- n-step reduction
- Left-most (canonical) reduction
- Right-most reduction

Sentences and Handles

- Sentential form
 - If $S \Rightarrow^* \alpha$, α is a sentential form of G .
 - Right (Left) sentential form
- Sentence
 - A sentential form with no nonterminals.
- Phrase
 - If $S \Rightarrow^* \alpha A \omega \wedge A \Rightarrow^+ \beta$, β is a phrase of $\alpha \beta \omega$ on A .
- Simple phrase
 - If $S \Rightarrow^* \alpha A \omega \wedge A \Rightarrow \beta$, β is a simple phrase.
- Handle
 - If $S \Rightarrow_{rm}^* \alpha A \omega \Rightarrow_{rm} \alpha \beta \omega$, $A \rightarrow \beta$ in the position following α is a handle of $\alpha \beta \omega$.

Define a Language by Grammar

- Language defined by a given CFG G
 - $L(G) = \{\alpha \mid S \Rightarrow^* \alpha \wedge \alpha \in \Sigma^*\}$
 - That is the set of all sentences of G
- Set of all sentential forms of G is
 - $SF(G) = \{\alpha \mid S \Rightarrow^* \alpha \wedge \alpha \in (\Sigma \cup N)^*\}$
- Thus we have $L(G) = SF(G) \cap \Sigma^*$.

BNF (Backus-Naur Form)

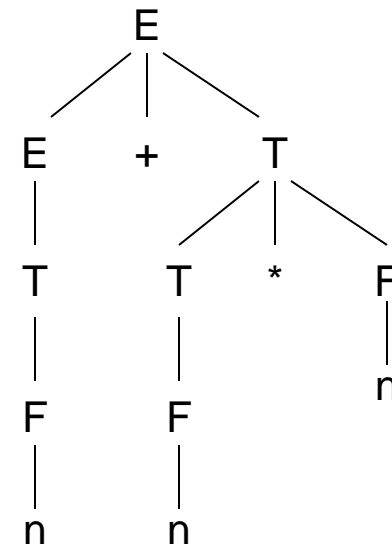
- Equivalent to CFG
 - Change " \rightarrow " to " $::=$ "
 - De facto standard to define a computer language
- EBNF (Extended BNF)
 - ISO/IEC 14977:1996(E). *The Standard Metalanguage Extended BNF*.
 - $=$, "terminals", [optional], {repetition}, (group), (*comment*), etc.

4. Parse Tree

- Graphic representation of derivations
 - Root: start symbol
 - Leafs: terminals **or** ϵ
 - Interior nodes: nonterminals
 - Offspring of a nonterminal: a production

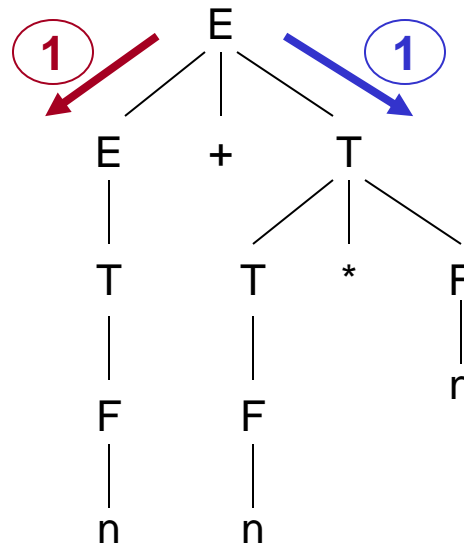
- An example

- $\mathbf{E} \Rightarrow \mathbf{E} + \mathbf{T} \Rightarrow \mathbf{E} + \mathbf{T} * \mathbf{F}$
 $\Rightarrow \mathbf{E} + \mathbf{T} * n \Rightarrow \mathbf{E} + \mathbf{F} * n$
 $\Rightarrow \mathbf{E} + n * n \Rightarrow \mathbf{T} + n * n$
 $\Rightarrow \mathbf{F} + n * n \Rightarrow n + n * n$



Insight into a Parse Tree

- An **abstraction** of derivations
 - Derivation order is discarded.



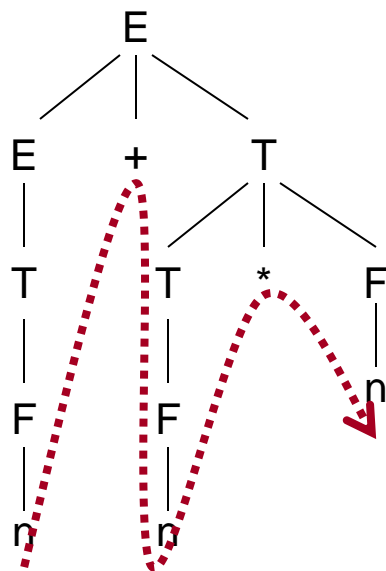
Sentences and Sentential Forms in a Parse Tree

- Sentence

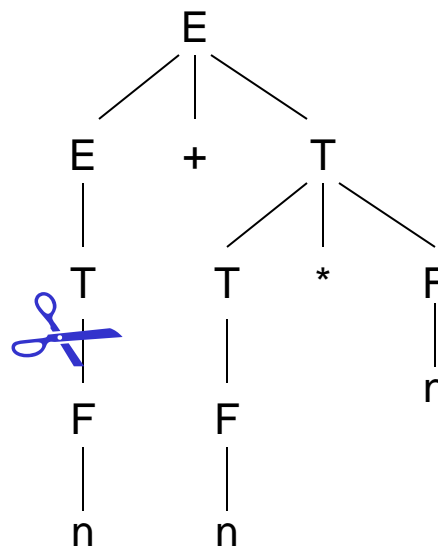
- Frontier of the parse tree, e.g. $n+n*n$

- Sentential Form

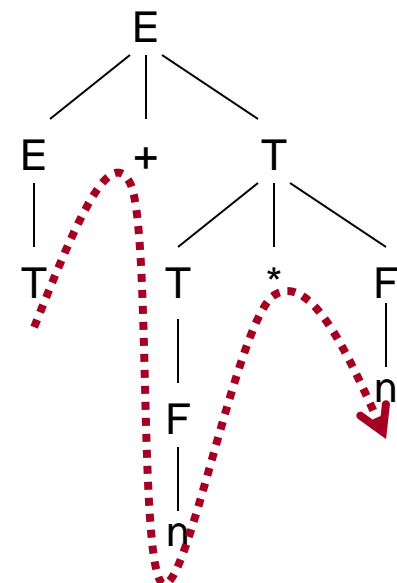
- Frontier of the pruned parse tree, e.g. $T+n*n$



Sentence



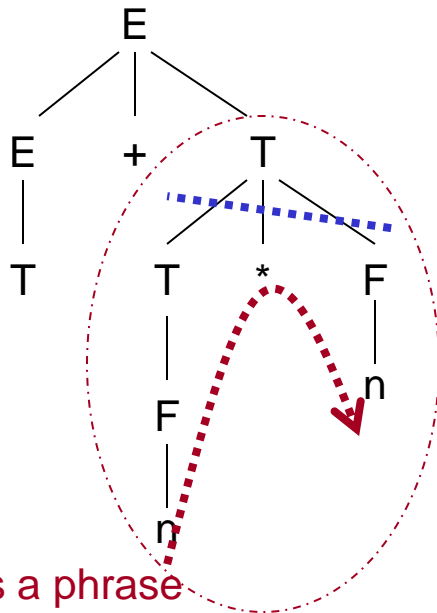
Pruning



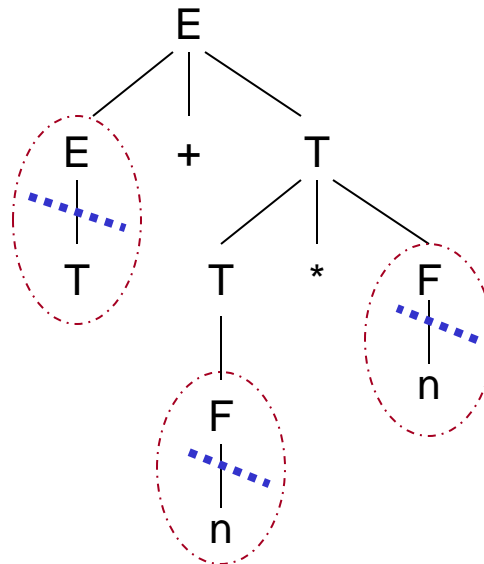
Sentential Form

Phrases and Handles in a Parse Tree

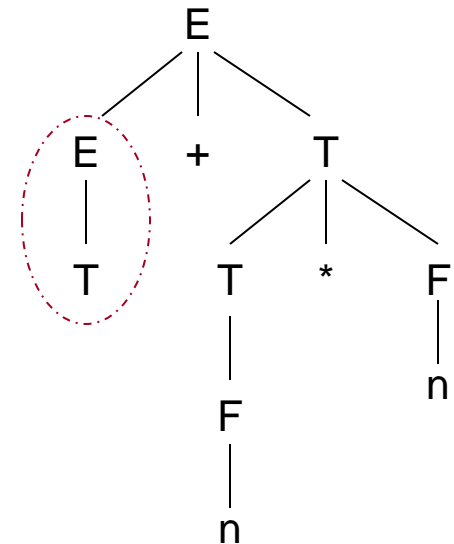
- For a given sentential form, e.g. $T+n*n$
 - Phrase: frontier of a nonterminal subtree
 - Simple phrase: frontier of a 2-level subtree
 - Handle: left-most simple phrase



$n*n$ is a phrase



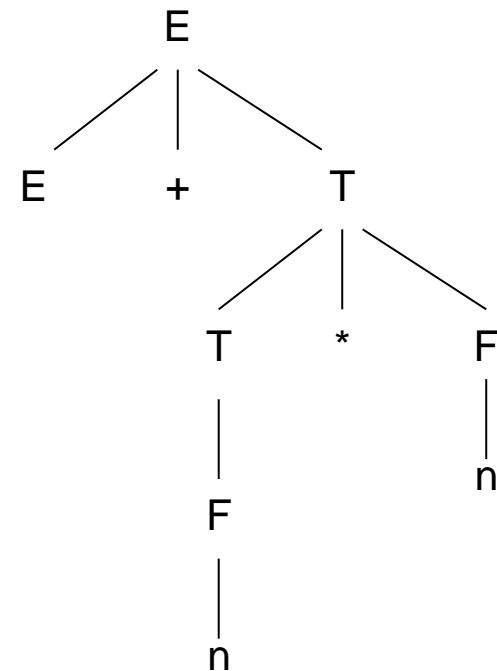
The only simple phrases are
 T , n , and n



Handle is unique: T

Discussions

- What is ...
 - The sentential form
 - All simple phrases
 - The only one handle



Parsing

- There are 2 questions that must be answered by parsing
 - Is the source program well-formed (legal in syntax)?
 - That is, is it a sentence of the language?
 - If it is, what's the syntax construction of the program?
 - Tree is an ideal data structure to present **part-of** relationship.
- 2 parsing strategies
 - Top-down vs. bottom-up

Top-Down Parsing

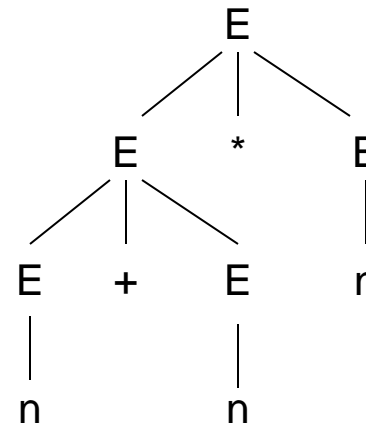
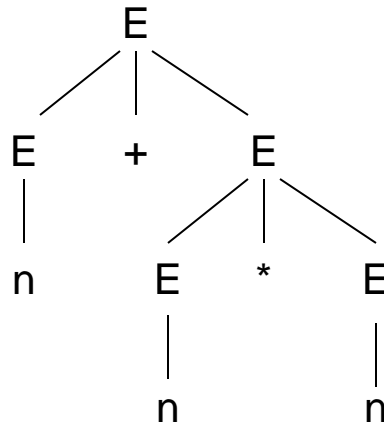
- Beginning at the start symbol
 - Usually expanding nonterminals in depth-first manner (predictive in nature)
 - Left-most derivation
 - Pre-order traversal of the parse tree
- Example: LL(k)
 - Read from **L**eft, and **L**eft-most derivation, with **k** lookaheads
 - Recursive descent (predictive) parsing

Bottom-Up Parsing

- Beginning from the terminal input string
 - Determining the production used to generate leaves
 - Right-most derivation in reverse order
 - Post-order traversal of the parse tree
- Example: LR(k)
 - Read from **L**eft, and **R**ight-most derivation, with **k** lookaheads
 - Parser generator yacc

5. Ambiguity

- A grammar with more than one parse tree for at least one sentence in the language
 - $E \rightarrow E + E \mid E * E \mid (E) \mid n$



2 parse trees for the sentence $n + n * n$

Ambiguous: Grammar vs. Language

- Ambiguous Language (inherent)
 - There exists no unambiguous grammar to define the language.
- Ambiguous Grammar (postnatal)
 - May be transformed to an equivalent unambiguous grammar.
 - The following two grammars define the same language
 - $E \rightarrow E + E \mid E * E \mid (E) \mid n$
 - $E \rightarrow E + T \mid T$
 $T \rightarrow T * F \mid F$
 $F \rightarrow (E) \mid n$

Binding to Practice

- Ambiguity in Programming Languages
 - Operators in an expression
 - Dangling `else` in embedded `if` statements
- Both of them are postnatal
 - Can be removed by additional disambiguation rules
 - Precedence and associativity
 - Inner "`if`" matched
 - Two approaches to resolving ambiguities
 - Rewrite the ambiguous grammar
 - Ad hoc solution: application of additional rules

6. Chomsky Hierarchy

Class	Grammar	Restriction	Recognizer
3	Regular	$A \rightarrow aB$ or $A \rightarrow a$, where $A, B \in N \wedge a \in \Sigma \cup \{\varepsilon\}$. $A \rightarrow \varepsilon$ permitted if A is the start symbol and does not appear on the right of any production.	Finite-State Automaton (FSA)
2	Context-Free	$A \rightarrow \alpha$, where $A \in N \wedge \alpha \in (\Sigma \cup N)^*$.	Push-Down Automaton (PDA)
1	Context-Sensitive	$\alpha \rightarrow \beta$, where $\alpha, \beta \in (\Sigma \cup N)^* \wedge \alpha \neq \varepsilon \wedge \alpha \leq \beta $. β can't be ε , unless α is the start symbol and does not appear on the right of any production.	Linear-Bounded Automaton (LBA)
0	Unrestricted	$\alpha \rightarrow \beta$, where $\alpha, \beta \in (\Sigma \cup N)^* \wedge \alpha \neq \varepsilon$.	Turing Machine (TM)

Useful in Practice

Useful in Theory

Regular Language

- Ability of regular languages
 - What is it able to define?
 - Identifiers
 - Decimal constants
 - What isn't it able to define?
 - Matched parentheses
- Satisfy the requirement of lexical descriptions



Context-Free Language

- Ability of context-free languages
 - What is it able to define?
 - Matched constructs
 - What isn't it able to define?
 - Only use declared variables
 - Matched parameter passing
- Satisfy the requirement of syntactic descriptions

Exercise 2.1

- Given the following grammar:

$$S \rightarrow (L) \mid a$$

$$L \rightarrow L , S \mid S$$

Construct a parse tree for the sentence
 $(a, ((a, a), (a, a)))$

Exercise 2.2

- Given the following grammar:

$\text{bexpr} \rightarrow \text{bexpr } \mathbf{or} \text{ bterm} \mid \text{bterm}$

$\text{bterm} \rightarrow \text{bterm } \mathbf{and} \text{ bfactor} \mid \text{bfactor}$

$\text{bfactor} \rightarrow \mathbf{not} \text{ bfactor} \mid (\text{bexpr}) \mid \mathbf{true} \mid \mathbf{false}$

Construct a parse tree for the sentence
not (true or false)

Exercise 2.3

- Is the grammar:
$$S \rightarrow a S b S \mid b S a S \mid \varepsilon$$
ambiguous? Why?

Further Reading

- Dragon Book, 2nd Edition (DBv2)
 - Comprehensive reading:
 - Section 2.1 – 2.3, 4.2.1 – 4.2.6 for the concepts of CFG, derivation, parse tree and ambiguity.
 - Section 4.5.1 – 4.5.2 for the concepts of reduction and handle.
 - Skip reading:
 - Section 2.4 for top-down parsing.
- Skip reading: definition of languages
 - BNF or EBNF of [Oberon-0](#), [Oberon](#), [OMG IDL](#), and [Lua](#) from our course website
- Skip reading: EBNF standard
 - ISO/IEC. The Standard Metalanguage EBNF.

Enjoy the Course!

