

Defining Program Syntax

Syntax And Semantics

- **Programming language syntax:** how programs look, their form and structure
 - Syntax is defined using a formal grammar
- **Programming language semantics:** what programs do, their behavior and meaning
 - Semantics is harder to define—more on this in Chapter 23

Outline

- Grammar and parse tree examples
- BNF and parse tree definitions
- Constructing grammars
- Phrase structure and lexical structure
- Other grammar forms

An English Grammar

A sentence $\langle S \rangle$ is a noun phrase $\langle NP \rangle$, a verb $\langle V \rangle$, and a noun phrase $\langle NP \rangle$.

$$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$$

A noun phrase $\langle NP \rangle$ is an article $\langle A \rangle$ and a noun $\langle N \rangle$.

$$\langle NP \rangle ::= \langle A \rangle \langle N \rangle$$

A verb $\langle V \rangle$ is...

$$\langle V \rangle ::= \mathbf{loves} \mid \mathbf{hates} \mid \mathbf{eats}$$

An article $\langle A \rangle$ is...

$$\langle A \rangle ::= \mathbf{a} \mid \mathbf{the}$$

A noun $\langle N \rangle$ is...

$$\langle N \rangle ::= \mathbf{dog} \mid \mathbf{cat} \mid \mathbf{rat}$$

How The Grammar Works

- The grammar is a set of rules that say how to build a tree—a *parse tree*
- $\langle S \rangle$ at the root of the tree
- The grammar's rules define how children can be added at any point in the tree
- For instance, $\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$

defines nodes $\langle NP \rangle$, $\langle V \rangle$, and $\langle NP \rangle$, in that order, as children of $\langle S \rangle$

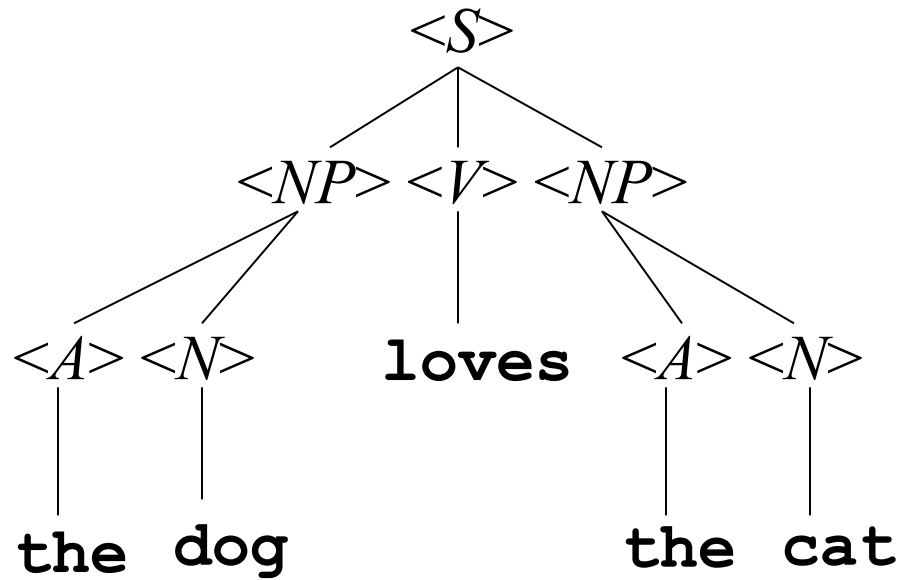
Parse Derivation

$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$
 $\langle NP \rangle ::= \langle A \rangle \langle N \rangle$
 $\langle V \rangle ::= \mathbf{loves} \mid \mathbf{hates} \mid \mathbf{eats}$
 $\langle A \rangle ::= \mathbf{a} \mid \mathbf{the}$
 $\langle N \rangle ::= \mathbf{dog} \mid \mathbf{cat} \mid \mathbf{rat}$

One **derivation** that $\langle S \rangle =$ **the dog loves the cat** is produced by the grammar rules:

$\langle S \rangle = \langle NP \rangle \quad \langle V \rangle \quad \langle NP \rangle$
 $= \langle A \rangle \langle N \rangle \langle V \rangle \quad \langle NP \rangle$
 $= \langle A \rangle \langle N \rangle \langle V \rangle \quad \langle A \rangle \langle N \rangle$
 $= \langle A \rangle \langle N \rangle \mathbf{loves} \langle A \rangle \langle N \rangle$
 $= \mathbf{the} \langle N \rangle \mathbf{loves} \langle A \rangle \langle N \rangle$
 $= \mathbf{the} \mathbf{dog} \mathbf{loves} \langle A \rangle \langle N \rangle$
 $= \mathbf{the} \mathbf{dog} \mathbf{loves} \mathbf{the} \langle N \rangle$
 $= \mathbf{the} \mathbf{dog} \mathbf{loves} \mathbf{the} \mathbf{cat}$

Parse Tree: the dog loves the cat

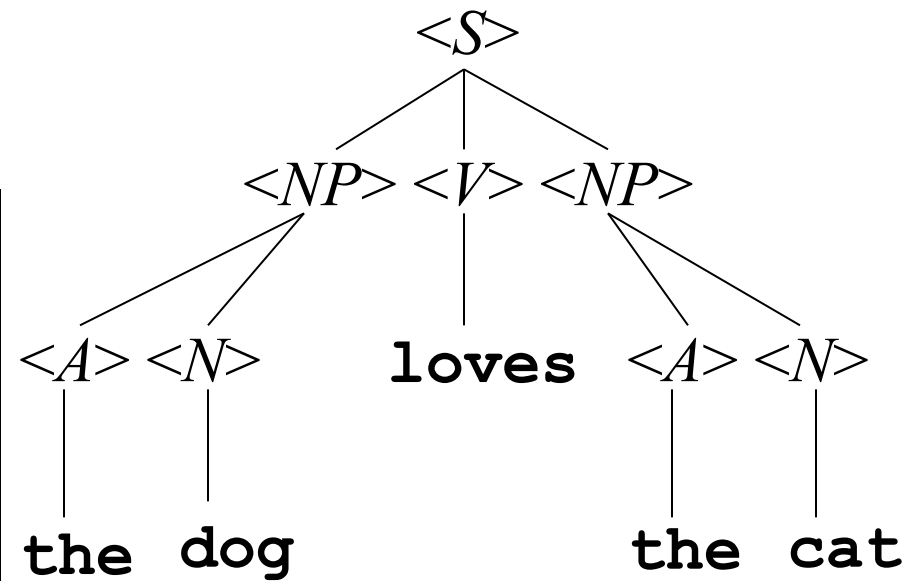


$\langle S \rangle = \langle NP \rangle \langle V \rangle \langle NP \rangle$
 $= \langle A \rangle \langle N \rangle \text{ loves } \langle A \rangle \langle N \rangle$
 $= \text{the dog loves the cat}$

$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$
 $\langle NP \rangle ::= \langle A \rangle \langle N \rangle$
 $\langle V \rangle ::= \text{loves} \mid \text{hates} \mid \text{eats}$
 $\langle A \rangle ::= \text{a} \mid \text{the}$
 $\langle N \rangle ::= \text{dog} \mid \text{cat} \mid \text{rat}$

Exercise 1

$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$
 $\langle NP \rangle ::= \langle A \rangle \langle N \rangle$
 $\langle V \rangle ::= \mathbf{loves} \mid \mathbf{hates} \mid \mathbf{eats}$
 $\langle A \rangle ::= \mathbf{a} \mid \mathbf{the}$
 $\langle N \rangle ::= \mathbf{dog} \mid \mathbf{cat} \mid \mathbf{rat}$



1. Which of the following are valid $\langle S \rangle$?
 - the dog hates the dog
 - dog loves the cat
 - loves the dog the cat
2. Parse:
 - a cat eats the rat
 - the dog loves cat

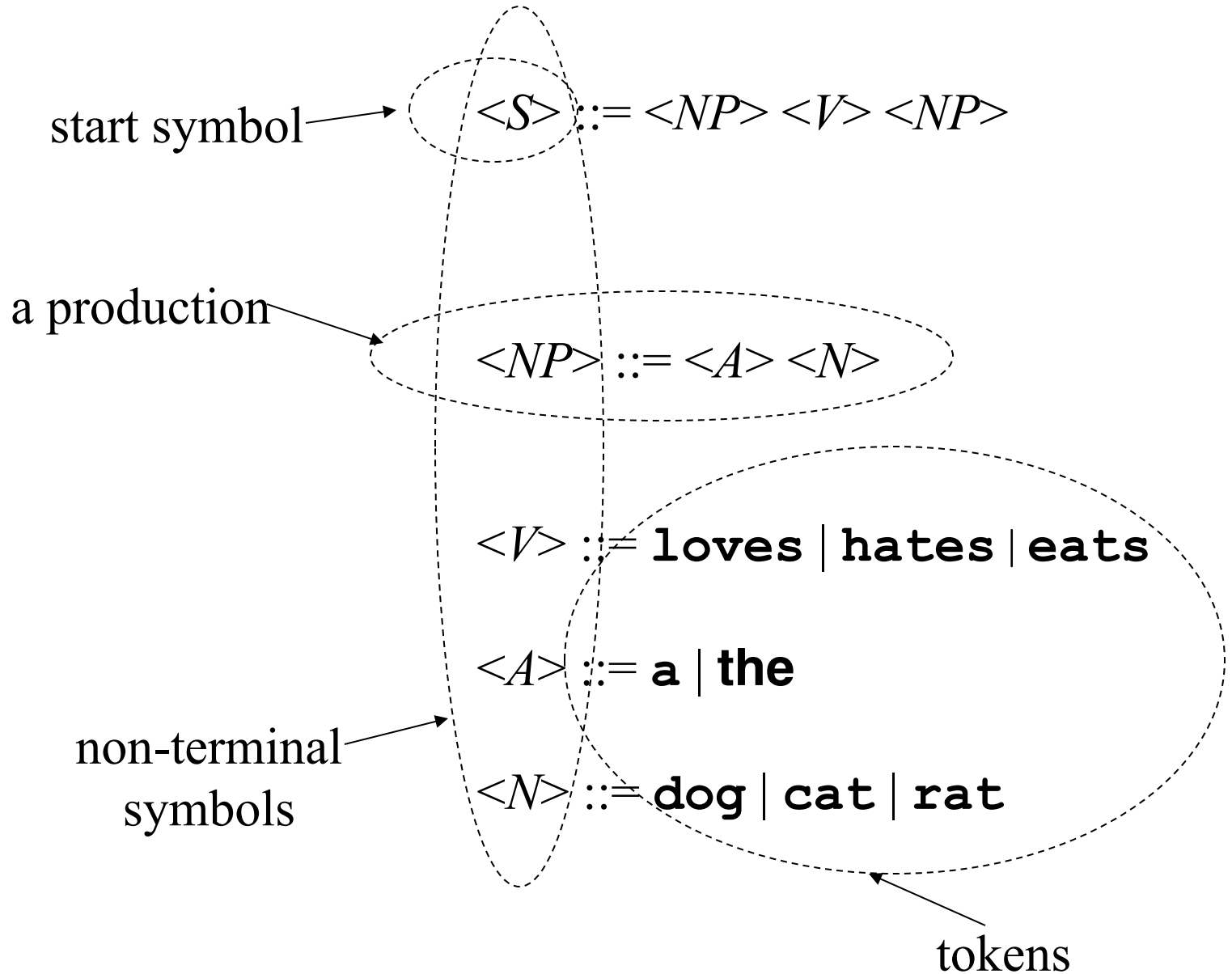
Outline

- Grammar and parse tree examples
- **BNF and parse tree definitions**
- Constructing grammars
- Phrase structure and lexical structure
- Other grammar forms

BNF Grammar Definition

- *Backus Naur Form* grammar consists of four parts:
 - The set of *tokens*
 - The set of *non-terminal symbols*
 - The *start symbol*
 - The set of *productions*

BNF Grammar Definitions Explained



Definition, Continued

- The *tokens* are the smallest units of syntax
 - Strings of one or more characters of program text
 - They are atomic: not treated as being composed from smaller parts
- The *non-terminal symbols* stand for larger pieces of syntax
 - They are strings enclosed in angle brackets, as in $\langle NP \rangle$
 - They are not strings that occur literally in program text
 - The grammar says how they can be expanded into strings of tokens
- The *start symbol* is the particular non-terminal that forms the root of any parse tree for the grammar

Definition, Continued

- The *productions* are the tree-building rules
- Each one has a left-hand side, the separator $::=$, and a right-hand side
 - The left-hand side is a single non-terminal
 - The right-hand side is a sequence of one or more things, each of which can be either a token or a non-terminal
- A production gives one possible way of building a parse tree: it permits the non-terminal symbol on the left-hand side to have the symbols on the right-hand side, in order, as its children in a parse tree

Alternatives (OR)

- When there is more than one production with the same left-hand side, an abbreviated form can be used
- In BNF grammar:
 - Gives the left-hand side (symbol),
 - the separator $::=$,
 - and then a list of possible right-hand sides separated by the special symbol $|$

Example

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle) \mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c}$$

Note that there are six productions in this grammar.
It is equivalent to this one:

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle$$
$$\langle exp \rangle ::= \langle exp \rangle * \langle exp \rangle$$
$$\langle exp \rangle ::= (\langle exp \rangle)$$
$$\langle exp \rangle ::= \mathbf{a}$$
$$\langle exp \rangle ::= \mathbf{b}$$
$$\langle exp \rangle ::= \mathbf{c}$$

Empty

- The special non-terminal $\langle empty \rangle$ is for places where you want the grammar to generate nothing
- For example, this grammar defines a typical if-then construct with an optional else part:

$$\begin{aligned} \langle if-stmt \rangle & ::= \mathbf{if} \langle expr \rangle \mathbf{then} \langle stmt \rangle \langle else-part \rangle \\ \langle else-part \rangle & ::= \mathbf{else} \langle stmt \rangle \mid \langle empty \rangle \end{aligned}$$

Grammar Parse Derivation

- Begin with a start symbol
- Choose a production with start symbol on left-hand side
- Replace start symbol with the right-hand side of that production

$\langle S \rangle ::= \langle NP \rangle \langle V \rangle \langle NP \rangle$
$\langle NP \rangle ::= \langle A \rangle \langle N \rangle$
$\langle V \rangle ::= \mathbf{loves} \mid \mathbf{hates} \mid \mathbf{eats}$
$\langle A \rangle ::= \mathbf{a} \mid \mathbf{the}$
$\langle N \rangle ::= \mathbf{dog} \mid \mathbf{cat} \mid \mathbf{rat}$

a cat eats the rat

1. Choose a non-terminal S in resulting string
2. Choose a production P with non-terminal S on its left-hand side
3. Replace S with the right-hand side of P
4. Repeat process until no non-terminals remain.

$\langle S \rangle$	=	$\langle NP \rangle$		$\langle V \rangle$		$\langle NP \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	$\langle V \rangle$		$\langle NP \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	$\langle V \rangle$	$\langle A \rangle$	$\langle N \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	$\langle N \rangle$	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	cat	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	cat	eats	the	$\langle N \rangle$
	=	a	cat	eats	the	rat

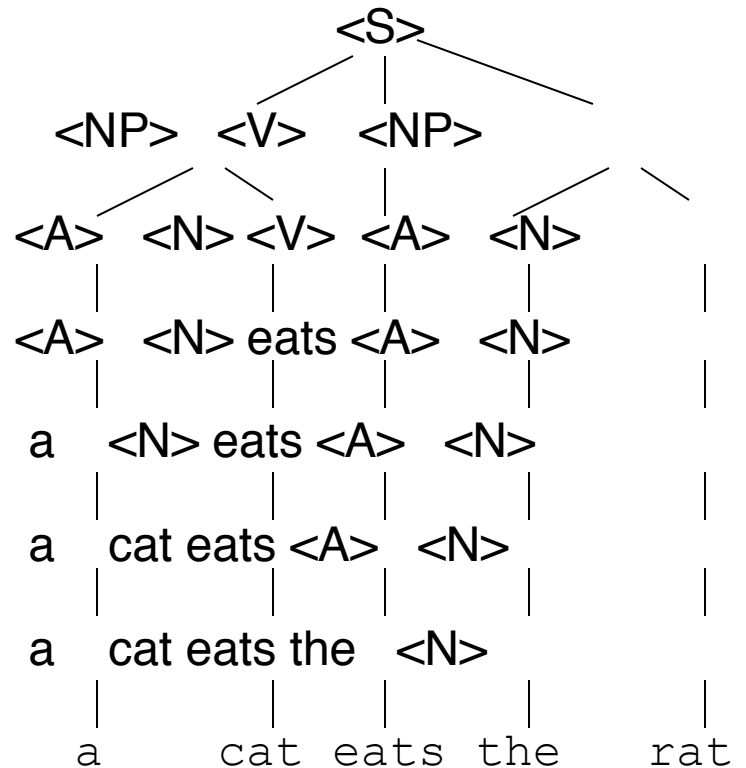
Parse Trees

- To build a parse tree, put the start symbol at the root
- Add children to every non-terminal, *following any one of the productions for that non-terminal in the grammar*
- Done when all the leaves are tokens
- Read off leaves from left to right—that is the string derived by the tree

$\langle S \rangle$::=	$\langle NP \rangle$	$\langle V \rangle$	$\langle NP \rangle$		
$\langle NP \rangle$::=	$\langle A \rangle$	$\langle N \rangle$			
$\langle V \rangle$::=	loves		hates		eats
$\langle A \rangle$::=	a		the		
$\langle N \rangle$::=	dog		cat		rat

$\langle S \rangle =$ a cat eats the rat

$\langle S \rangle$	=	$\langle NP \rangle$		$\langle V \rangle$		$\langle NP \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	$\langle V \rangle$		$\langle NP \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	$\langle V \rangle$	$\langle A \rangle$	$\langle N \rangle$
	=	$\langle A \rangle$	$\langle N \rangle$	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	$\langle N \rangle$	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	cat	eats	$\langle A \rangle$	$\langle N \rangle$
	=	a	cat	eats	the	$\langle N \rangle$
	=	a	cat	eats	the	rat



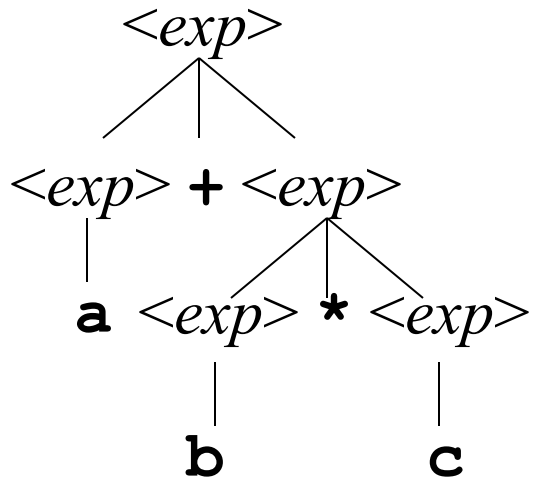
A Programming Language Grammar

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle) \mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c}$$

- An expression can be:
 - the sum of two expressions,
 - or the product of two expressions,
 - or a parenthesized subexpression,
 - or **a**,
 - or **b**,
 - or **c**

Parse and Parse Tree: $a+b*c$

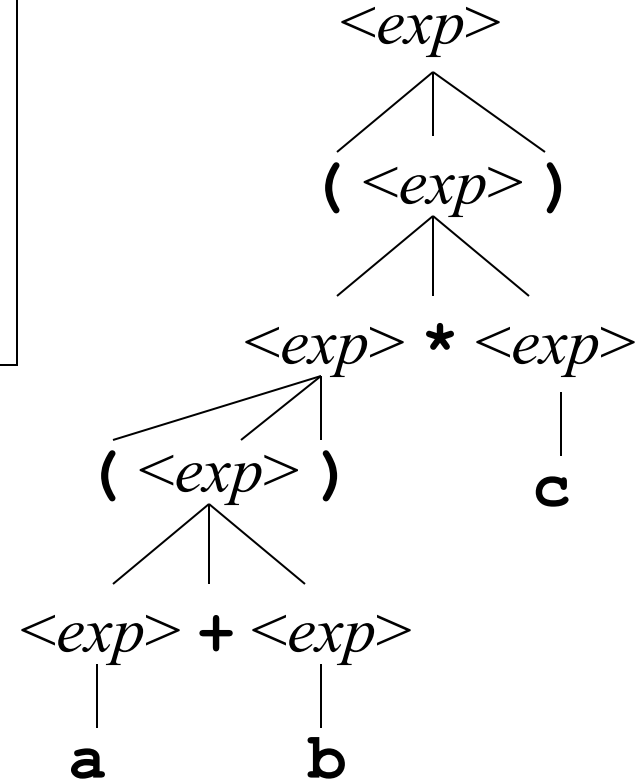
$\langle \text{exp} \rangle = \langle \text{exp} \rangle + \langle \text{exp} \rangle$
 $= a + \langle \text{exp} \rangle$
 $= a + \langle \text{exp} \rangle * \langle \text{exp} \rangle$
 $= a + b * c$



$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \mid (\langle \text{exp} \rangle) \mid a \mid b \mid c$

Parse and Parse Tree: $((a+b)*c)$

$\langle \text{exp} \rangle = (\langle \text{exp} \rangle)$
 $= (\langle \text{exp} \rangle * \langle \text{exp} \rangle)$
 $= ((\langle \text{exp} \rangle) * \langle \text{exp} \rangle)$
 $= ((\langle \text{exp} \rangle) * c)$
 $= ((\langle \text{exp} \rangle + \langle \text{exp} \rangle) * c)$
 $= ((a + b) * c)$



$\langle \text{exp} \rangle ::= \langle \text{exp} \rangle + \langle \text{exp} \rangle \mid \langle \text{exp} \rangle * \langle \text{exp} \rangle \mid (\langle \text{exp} \rangle) \mid a \mid b \mid c$

Exercise 2

$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle) \mid a \mid b \mid c$
--

Give the parse tree for each of these strings:

a. a+b

b. a*b+c

c. (a+b)*c

Compiler Note

- What we just did is *parsing*: trying to find a parse tree for a given string
- That's what compilers do for every program you try to compile: try to build a parse tree for your program, using the grammar for whatever language you used
- Take a course in compiler construction to learn about algorithms for doing this efficiently

Language Definition

- We use *grammars* to define the syntax of programming languages
- The language defined by a grammar is the set of all strings that can be derived by some parse tree for the grammar
- As in the previous example, that set is often infinite (though grammars are finite)
- Constructing grammars is a little like programming...

Outline

- Grammar and parse tree examples
- BNF and parse tree definitions
- **Constructing grammars**
- Phrase structure and lexical structure
- Other grammar forms

Constructing Grammars

- Most important trick: divide and conquer
- Example: the language of Java declarations:
 - a type name,
 - a list of variables separated by commas,
 - and a semicolon
- Each variable can optionally be followed by an initializer:

```
float a;  
boolean a,b,c;  
int a=1, b, c=1+2;
```

Example, Continued

```
int a=1, b, c=1+2;
```

- Easy if we postpone defining the comma-separated list of variables with initializers:

<var-dec> ::= *<type-name>* *<declarator-list>* ;

- Primitive type names are easy enough too:

<type-name> ::= **boolean** | **byte** | **short** | **int**
 | **long** | **char** | **float** | **double**

- (Note: skipping constructed types: class names, interface names, and array types)

Example, Continued

- That leaves the comma-separated list of variables with initializers
- Again, postpone defining variables with initializers, and just do the comma-separated list part:

```
int a=1, b, c=1+2;
```

<var-dec> ::= <type-name> <declarator-list> ;

<declarator-list> ::= <declarator>

| <declarator> , <declarator-list>

Example, Continued

```
int a=1, b, c=1+2;
```

- That leaves the variables with initializers:

$\langle var-dec \rangle ::= \langle type-name \rangle \langle declarator-list \rangle ;$

$\langle declarator-list \rangle ::= \langle declarator \rangle$
| $\langle declarator \rangle , \langle declarator-list \rangle$

$\langle declarator \rangle ::= \langle variable-name \rangle$
| $\langle variable-name \rangle = \langle expr \rangle$

- For full Java, we would need to allow pairs of square brackets after the variable name
- There is also a syntax for array initializers
- And definitions for $\langle variable-name \rangle$ and $\langle expr \rangle$

Grammar Construction Example

Construct a grammar in BNF for each language:

1. $\langle \text{digit} \rangle$ as a character 0-9.

$$\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

2. $\langle \text{unsigned} \rangle$ as the set of all strings with one or more $\langle \text{digit} \rangle$. Note the left-recursion.

$$\langle \textit{unsigned} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{unsigned} \rangle \langle \textit{digit} \rangle$$

3. $\langle \text{signed} \rangle$ as the set of all strings starting with $-$ or $+$ and followed by an $\langle \text{unsigned} \rangle$.

$$\langle \textit{signed} \rangle ::= +\langle \textit{unsigned} \rangle \mid -\langle \textit{unsigned} \rangle$$

Grammar Construction Example

1. $\langle \text{digit} \rangle$ as a character 0-9.

$$\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$$

2. $\langle \text{unsigned} \rangle$ as the set of all strings with one or more $\langle \text{digit} \rangle$. Note the left-recursion.

$$\langle \textit{unsigned} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{unsigned} \rangle \langle \textit{digit} \rangle$$

Parse 321 $\langle \textit{unsigned} \rangle$

/ \

$\langle \textit{unsigned} \rangle$ $\langle \textit{digit} \rangle$

/ \ \

$\langle \textit{unsigned} \rangle$ $\langle \textit{digit} \rangle$ 1

/ \

$\langle \textit{digit} \rangle$ 2

/

3

Exercise 3

$\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$\langle \textit{unsigned} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{unsigned} \rangle \langle \textit{digit} \rangle$
--

$\langle \textit{signed} \rangle ::= +\langle \textit{unsigned} \rangle \mid -\langle \textit{unsigned} \rangle$
--

Construct a grammar in BNF for each language:

1. $\langle \textit{integer} \rangle$ as the set of all strings of $\langle \textit{signed} \rangle$ or $\langle \textit{unsigned} \rangle$.
2. $\langle \textit{decimal} \rangle$ as the set of all strings of $\langle \textit{integer} \rangle$ followed by a '.' and optionally followed by an $\langle \textit{unsigned} \rangle$.
3. $\langle \textit{2or3digits} \rangle$ as the set of all strings of two or three $\langle \textit{digit} \rangle$.
4. $\langle \textit{2's} \rangle$ as the set of all strings of one or more 2's.
5. $\langle \textit{1+2's} \rangle$ as the set of all strings beginning with '1' and followed by any number of 2's.
6. $\langle \textit{2's+1} \rangle$ as the set of all strings beginning with any number of 2's and followed by a '1'.
7. $\langle \textit{AdigitBs} \rangle$ as the set of all strings beginning with 'A' and optionally followed by any number of $\langle \textit{digit} \rangle$ or 'B'.

Outline

- Grammar and parse tree examples
- BNF and parse tree definitions
- Constructing grammars
- **Phrase structure and lexical structure**
- Other grammar forms

Where Do Tokens Come From?

$\langle \textit{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$\langle \textit{unsigned} \rangle ::= \langle \textit{digit} \rangle \mid \langle \textit{unsigned} \rangle \langle \textit{digit} \rangle$

- Tokens are pieces of program text that we choose not to think of as being built from smaller pieces
- Identifiers (**count**), keywords (**if**), operators (**==**), constants (**123.4**), etc.
- Programs stored in files are just sequences of characters
- How is such a file divided into a sequence of tokens?

Lexical Structure And Phrase Structure

- *Phrase structure*: how a program is built from a sequence of tokens

```
<if-stmt> ::= if <expr> then <stmt> <else-part>  
<else-part> ::= else <stmt> | <empty>
```

- *Lexical structure*: how tokens are built from a sequence of characters

```
<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9
```

```
<unsigned> ::= <digit> | <unsigned> <digit>
```

One Grammar For Both

- You could do it all with one grammar by using characters as the only tokens
- Not done in practice: things like white space and comments would make the grammar too messy to be readable

$\langle \text{if-stmt} \rangle ::= \mathbf{if} \langle \text{white-space} \rangle \langle \text{expr} \rangle \langle \text{white-space} \rangle$
 $\mathbf{then} \langle \text{white-space} \rangle$
 $\langle \text{stmt} \rangle \langle \text{white-space} \rangle \langle \text{else-part} \rangle$

$\langle \text{else-part} \rangle ::= \mathbf{else} \langle \text{white-space} \rangle \langle \text{stmt} \rangle \mid \langle \text{empty} \rangle$

Separate Grammars

- Usually there are two separate grammars
 - One says how to construct a sequence of tokens from a file of characters
 - One says how to construct a parse tree from a sequence of tokens

$\langle \textit{program-file} \rangle ::= \langle \textit{end-of-file} \rangle \mid \langle \textit{element} \rangle \langle \textit{program-file} \rangle$

$\langle \textit{element} \rangle ::= \langle \textit{token} \rangle \mid \langle \textit{one-white-space} \rangle \mid \langle \textit{comment} \rangle$

$\langle \textit{one-white-space} \rangle ::= \langle \textit{space} \rangle \mid \langle \textit{tab} \rangle \mid \langle \textit{end-of-line} \rangle$

$\langle \textit{token} \rangle ::= \langle \textit{identifier} \rangle \mid \langle \textit{operator} \rangle \mid \langle \textit{constant} \rangle \mid \dots$

Separate Compiler Passes

- The *scanner* reads the input file and divides it into tokens according to the first grammar
- The scanner discards white space and comments
- The *parser* constructs a parse tree (or at least goes through the motions—more about this later) from the token stream according to the second grammar

Exercise 4

```
<space> ::=  
<digit> ::= 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9  
<unsigned> ::= <digit> | <unsigned> <digit>  
<signed> ::= +<unsigned> | -<unsigned>  
<integer> ::= <signed> | <unsigned>  
<decimal> ::= <integer>.<unsigned> | <integer> .  
<operator> ::= + | == | =  
<identifier> ::= x | y  
<constant> ::= <integer> | <decimal>  
<keyword> ::= if | then | endif
```

List the *scanner* output from the following:

```
if x == 3.14 then y = x + y endif
```

Historical Note #1

- Early languages sometimes did not separate lexical structure from phrase structure
 - Early Fortran and Algol dialects allowed spaces anywhere, even in the middle of a keyword
 - Other languages like PL/I allow keywords to be used as identifiers
- This makes them harder to scan and parse
- It also reduces readability

Historical Note #2

- Some languages have a *fixed-format* lexical structure—column positions are significant
 - One statement per line (i.e. per card)
 - First few columns for statement label
 - Etc.
- Early dialects of Fortran, Cobol, and Basic
- Almost all modern languages are *free-format*: column positions are ignored

Outline

- Grammar and parse tree examples
- BNF and parse tree definitions
- Constructing grammars
- Phrase structure and lexical structure
- **Other grammar forms**

Other Grammar Forms

- BNF variations
- EBNF variations
- Syntax diagrams

BNF Variations

- Some use \rightarrow or $=$ instead of $::=$
- Some leave out the angle brackets and use a distinct typeface for tokens
- Some allow single quotes around tokens, for example to distinguish ' | ' as a token from | as a meta-symbol

EBNF Variations

- Additional syntax to simplify some grammar chores:
 - $\{x\}$ or x^* to mean zero or more repetitions of x
 - x^+ to mean one or more repetitions of x
 - $[x]$ to mean x is optional (i.e. $x \mid \langle empty \rangle$)
 - $()$ for grouping
 - $|$ anywhere to mean a choice among alternatives
 - Quotes around tokens, if necessary, to distinguish from all these meta-symbols

EBNF Examples

$\langle \text{if-stmt} \rangle ::= \mathbf{if} \langle \text{expr} \rangle \mathbf{then} \langle \text{stmt} \rangle [\mathbf{else} \langle \text{stmt} \rangle]$

$\langle \text{stmt-list} \rangle ::= \{ \langle \text{stmt} \rangle ; \}$

$\langle \text{thing-list} \rangle ::= \{ (\langle \text{stmt} \rangle \mid \langle \text{declaration} \rangle) ; \}$

$\langle \text{digit} \rangle ::= 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$

$\langle \text{unsigned} \rangle ::= \langle \text{digit} \rangle^+$

$\langle \text{signed} \rangle ::= (+ \mid -) \langle \text{unsigned} \rangle$

- Anything that extends BNF this way is called an Extended BNF: EBNF
- There are many variations

Exercise 5

Construct a grammar in EBNF for each language:

1. `<unsigned>` as the set of all strings with one or more `<digit>`.
2. `<signed>` as the set of all strings starting with `-` or `+` and followed by an `<unsigned>`.
3. `<integer>` as the set of all strings of `<signed>` or `<unsigned>`.
4. `<decimal>` as the set of all strings of `<integer>` followed by a `'.'` and optionally followed by an `<unsigned>`.

EBNF

Extensions

`{x}` or `x*` to mean zero or more repetitions of `x`

`x+` to mean one or more repetitions of `x`

`[x]` to mean `x` is optional (i.e. `x | <empty>`)

`()` for grouping

`|` anywhere to mean a choice among alternatives

Exercise 5 continued

$\{x\}$ or x^* to mean zero or more repetitions of x

x^+ to mean one or more repetitions of x

$[x]$ to mean x is optional (i.e. $x \mid \langle \text{empty} \rangle$)

$()$ for grouping

$|$ anywhere to mean a choice among alternatives

6. $\langle \text{identifier} \rangle$ as the set of all strings starting with $\langle \text{alpha} \rangle$ and followed by zero or more $\langle \text{alpha} \rangle$ or $\langle \text{digit} \rangle$.
7. $\langle 1+2' \text{ s} \rangle$ as the set of all strings beginning with '1' and followed by any number of 2' s.
8. $\langle \text{AdigitBs} \rangle$ as the set of all strings beginning with 'A' and optionally followed by any number of $\langle \text{digit} \rangle$ or 'B'.
9. Scientific notation (e.g. 1.2E-13)

Syntax Diagrams

- Syntax diagrams (“railroad diagrams”)
- Start with an EBNF grammar
- A simple production is just a chain of boxes (for nonterminals) and ovals (for terminals):

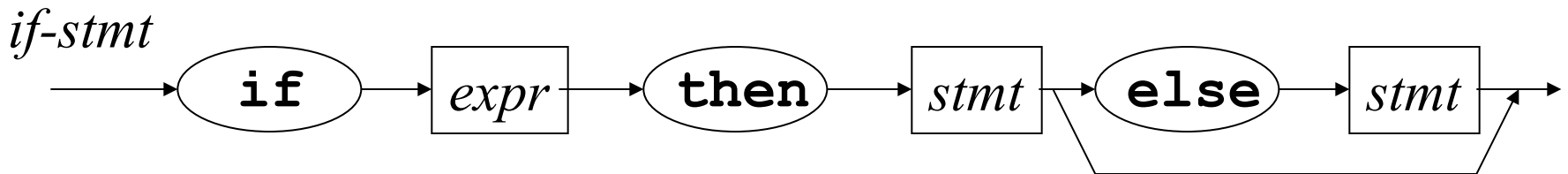
$\langle \text{if-stmt} \rangle ::= \mathbf{if} \langle \text{expr} \rangle \mathbf{then} \langle \text{stmt} \rangle \mathbf{else} \langle \text{stmt} \rangle$



Bypasses

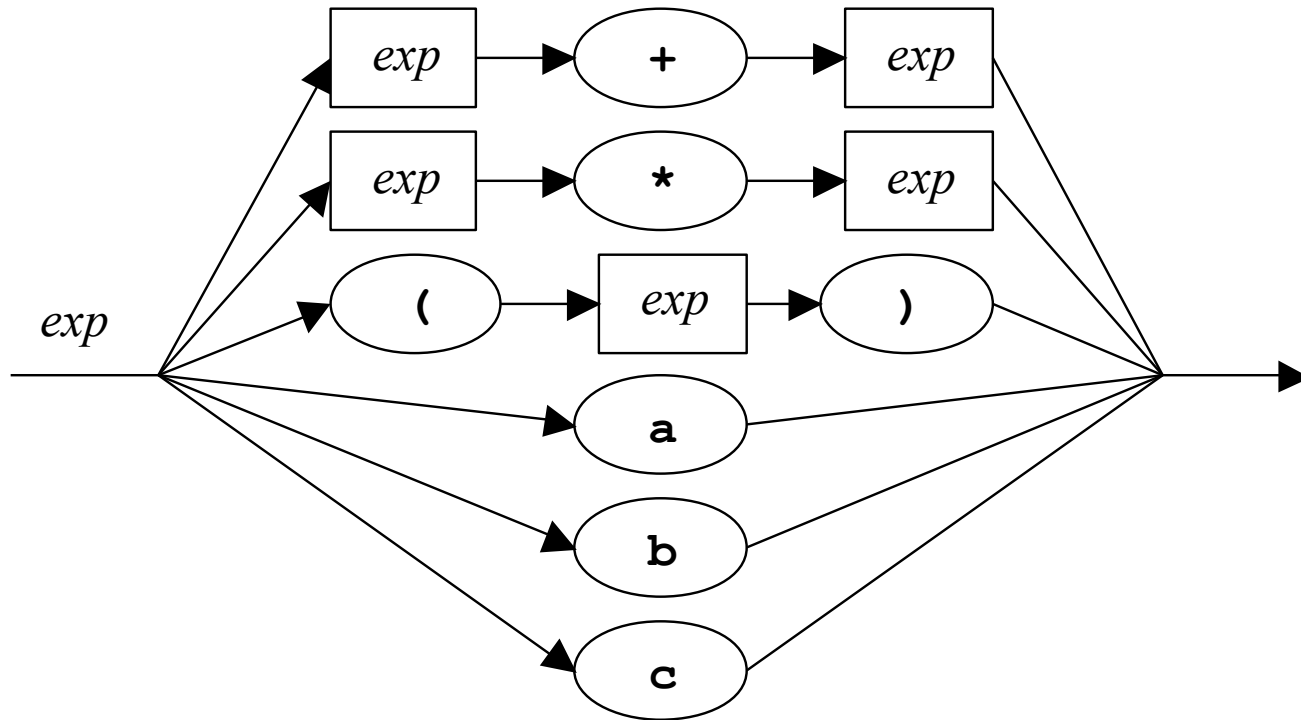
- Square-bracket pieces from the EBNF get paths that bypass them

$\langle if-stmt \rangle ::= \mathbf{if} \langle expr \rangle \mathbf{then} \langle stmt \rangle [\mathbf{else} \langle stmt \rangle]$



Branching

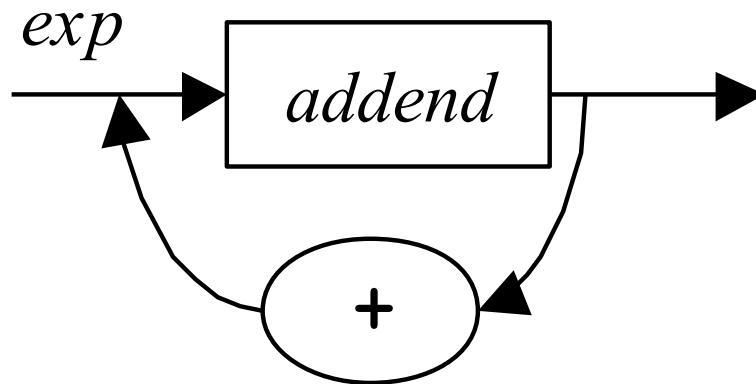
- Use branching for multiple productions

$$\langle exp \rangle ::= \langle exp \rangle + \langle exp \rangle \mid \langle exp \rangle * \langle exp \rangle \mid (\langle exp \rangle)$$
$$\mid \mathbf{a} \mid \mathbf{b} \mid \mathbf{c}$$


Loops

- Use loops for EBNF curly brackets

$\langle exp \rangle ::= \langle addend \rangle \{ + \langle addend \rangle \}$



Syntax Diagrams, Pro and Con

- Easier for people to read casually
- Harder to read precisely: what will the parse tree look like?
- Harder to make machine readable (for automatic parser-generators)

Formal Context-Free Grammars

- In the study of formal languages and automata, grammars are expressed in yet another notation:

$S \rightarrow aSb \mid X$ S is a string of symbols **a S b** or **X**.
 $X \rightarrow cX \mid \epsilon$ X is a string of symbols **c X** or empty

- These are called *context-free grammars* because children of a node only depend on that node's non-terminal symbol; not on the context of neighboring nodes in the tree. Simpler to define and compile.
- Context sensitive language elements include scope but is not generally part of a grammar.
- Other kinds of grammars are also studied: *regular grammars* (weaker), *context-sensitive grammars* (stronger), etc.

Many Other Variations

- BNF and EBNF ideas are widely used
- Exact notation differs, in spite of occasional efforts to get uniformity
- But as long as you understand the ideas, differences in notation are easy to pick up

Example

WhileStatement:

while (*Expression*) *Statement*

DoStatement:

do *Statement* while (*Expression*);

ForStatement:

for (*ForInit*_{opt} ; *Expression*_{opt} ; *ForUpdate*_{opt})
Statement

[from *The Java™ Language Specification*,
James Gosling et. al.]

Conclusion

- We use grammars to define programming language syntax, both lexical structure and phrase structure
- Connection between theory and practice
 - Two grammars, two compiler passes
 - Parser-generators can write code for those two passes automatically from grammars

Conclusion, Continued

- Multiple audiences for a grammar
 - Novices want to find out what legal programs look like
 - Experts—advanced users and language system implementers—want an exact, detailed definition
 - Tools—parser and scanner generators—want an exact, detailed definition in a particular, machine-readable form