

Attribute Grammars and Abstract Syntax Trees

Attribute Grammars

Attribute grammar: Context-free grammar extended with *attributes* associated to the symbols of the grammar. The attributes provide context-sensitive information that can be transported across the syntax tree.

Attribute grammars were proposed by Donald Knuth to formalize the semantics of a context-free language. They are useful for:

- semantic analysis (e.g., type checking)
- intermediate-code generation
- language interpretation

Attribute grammars can produce *decorated* syntax trees.

Attribute Grammars

Example: write a grammar for $L = \{a^n b^n c^n \mid n \geq 1\}$.
(it is known that no context-free grammar exists)

First attempt:

L : A B C

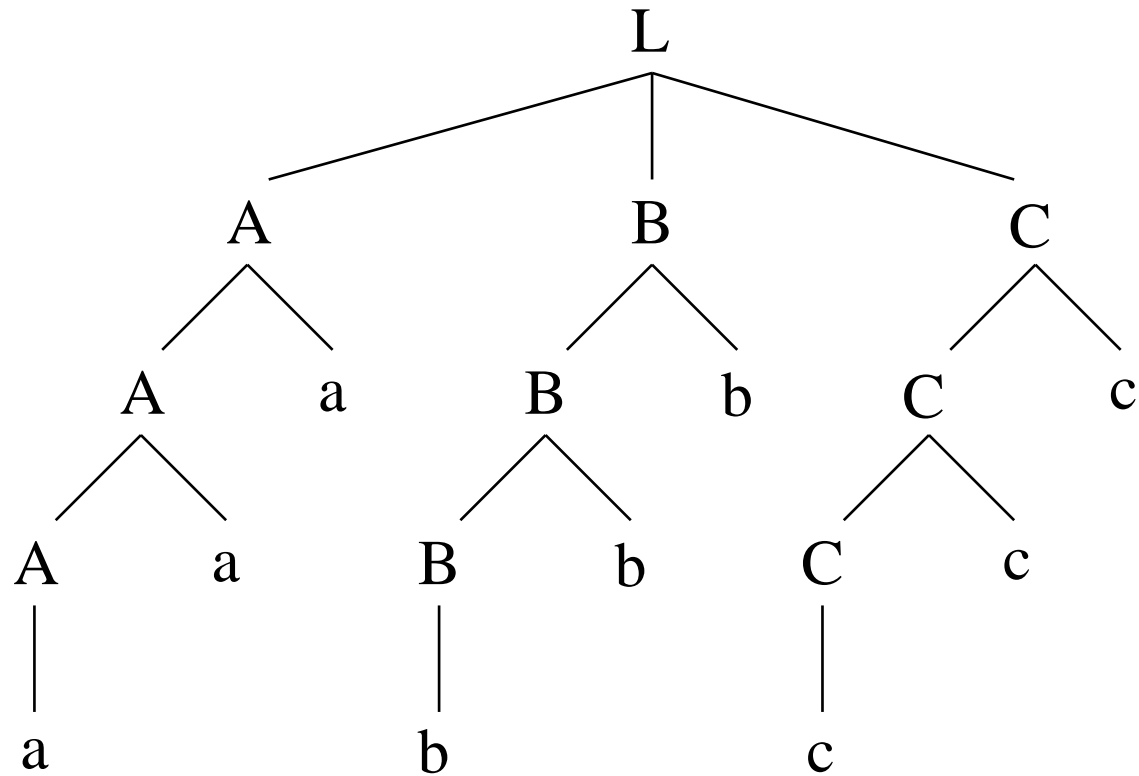
A : a | A a

B : b | B b

C : c | C c

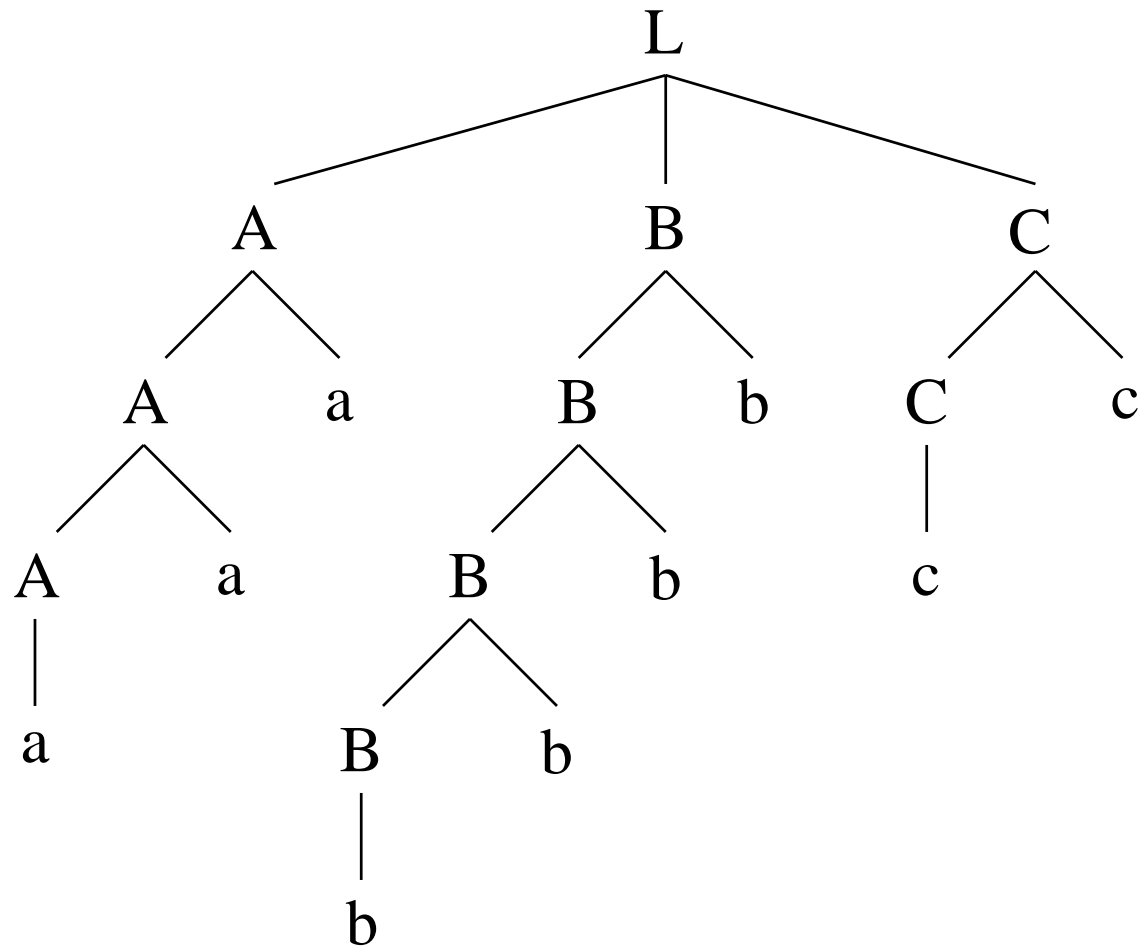
It recognizes $\{a^n b^m c^p \mid n \geq 1, m \geq 1, p \geq 1\}$.

Attribute Grammars



OK!

Attribute Grammars



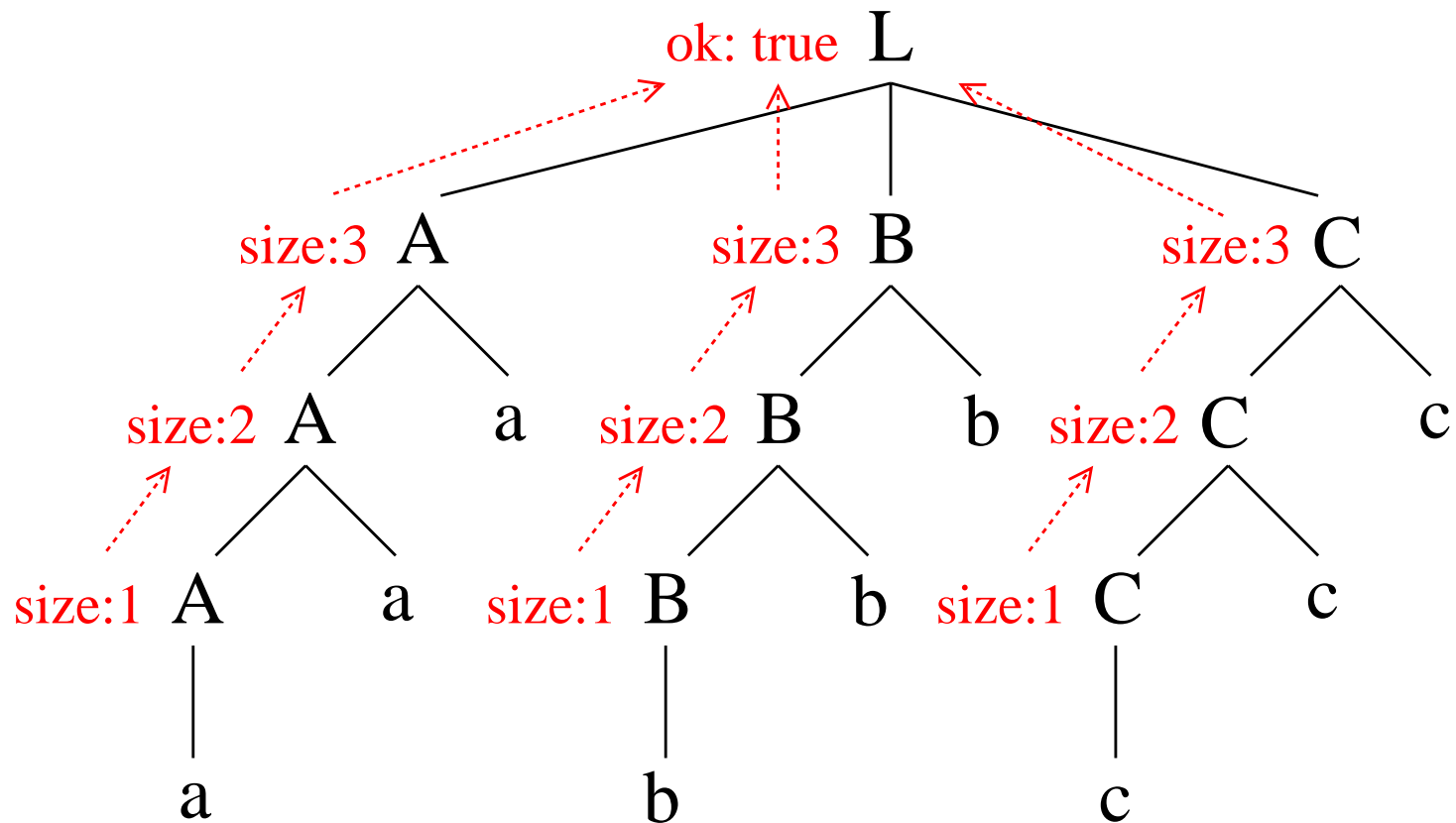
KO

Synthesized Attributes

Augment the grammar with a *synthesized* attribute associated to **A**, **B** and **C** that indicates the length of the sequence. Another synthesized attribute (**ok**) is associated to **L** to indicate that the sequence is correct.

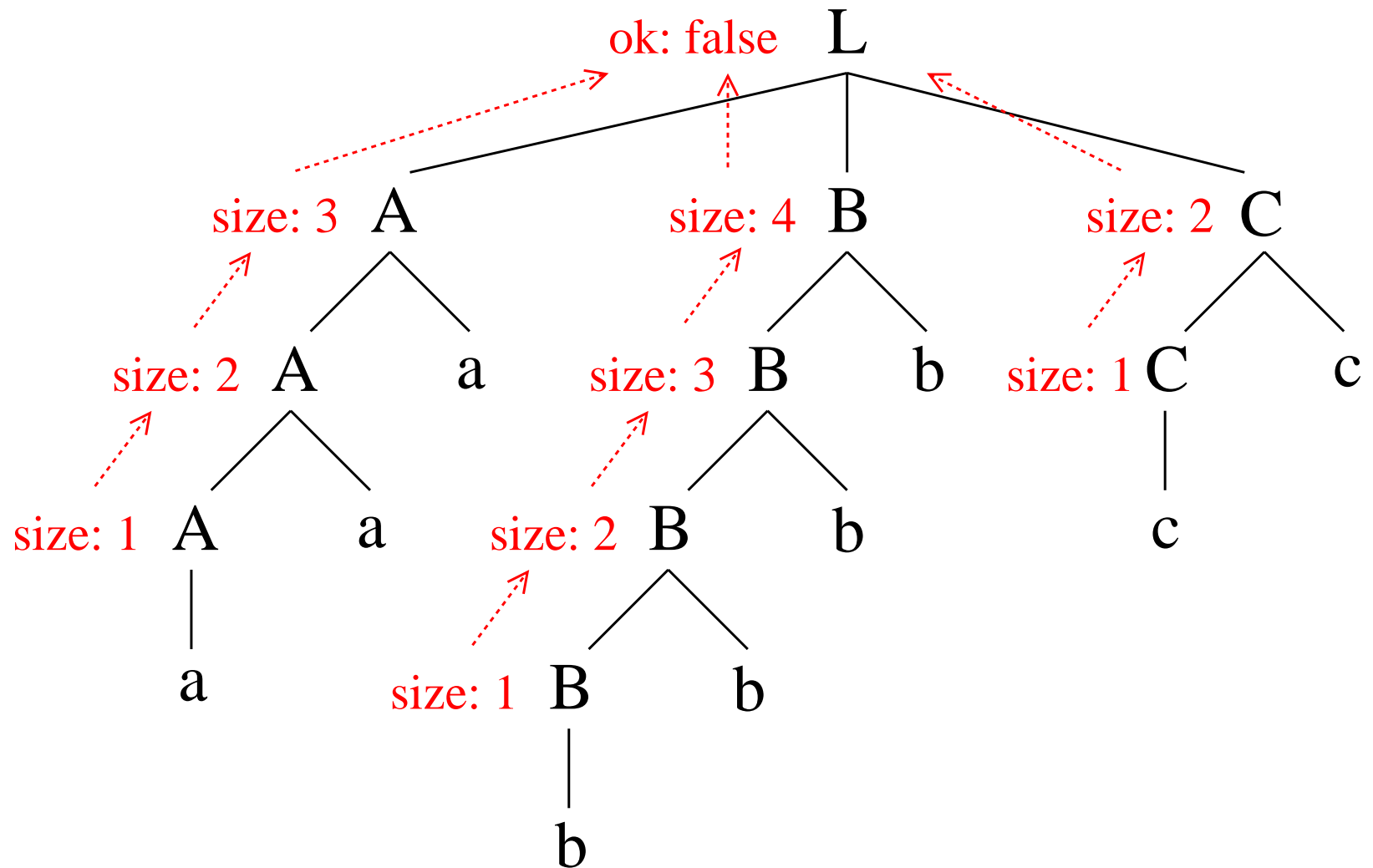
```
L : A B C { L.ok = (A.size == B.size)  
               and (A.size == C.size) }  
  
A : a      { A.size = 1 }  
    | A' a { A.size = A'.size + 1 }  
  
B : b      { B.size = 1 }  
    | B' b { B.size = B'.size + 1 }  
  
C : c      { C.size = 1 }  
    | C' c { C.size = C'.size + 1 }
```

Synthesized Attributes



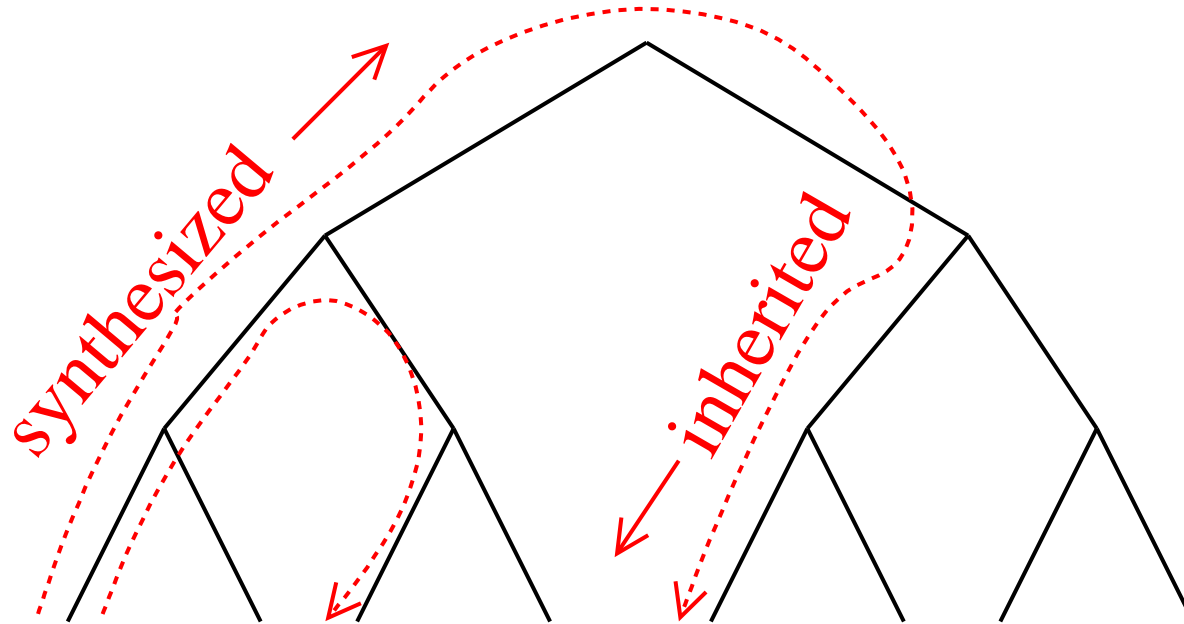
Synthesized attributes are calculated from RHS to LHS of the production rules (bottom-up in the parse tree).

Synthesized Attributes



Inherited Attributes

Sometimes it is necessary to move some information up to some node in the tree and transfer it down to some other part of the tree.



Inherited Attributes

Inherited attributes are propagated to the symbols of the RHS of the production rules (top-down or left-to-right in the parse tree).

The grammar can be augmented with an inherited attribute (**InhSize**) for the symbols **B** and **C**.

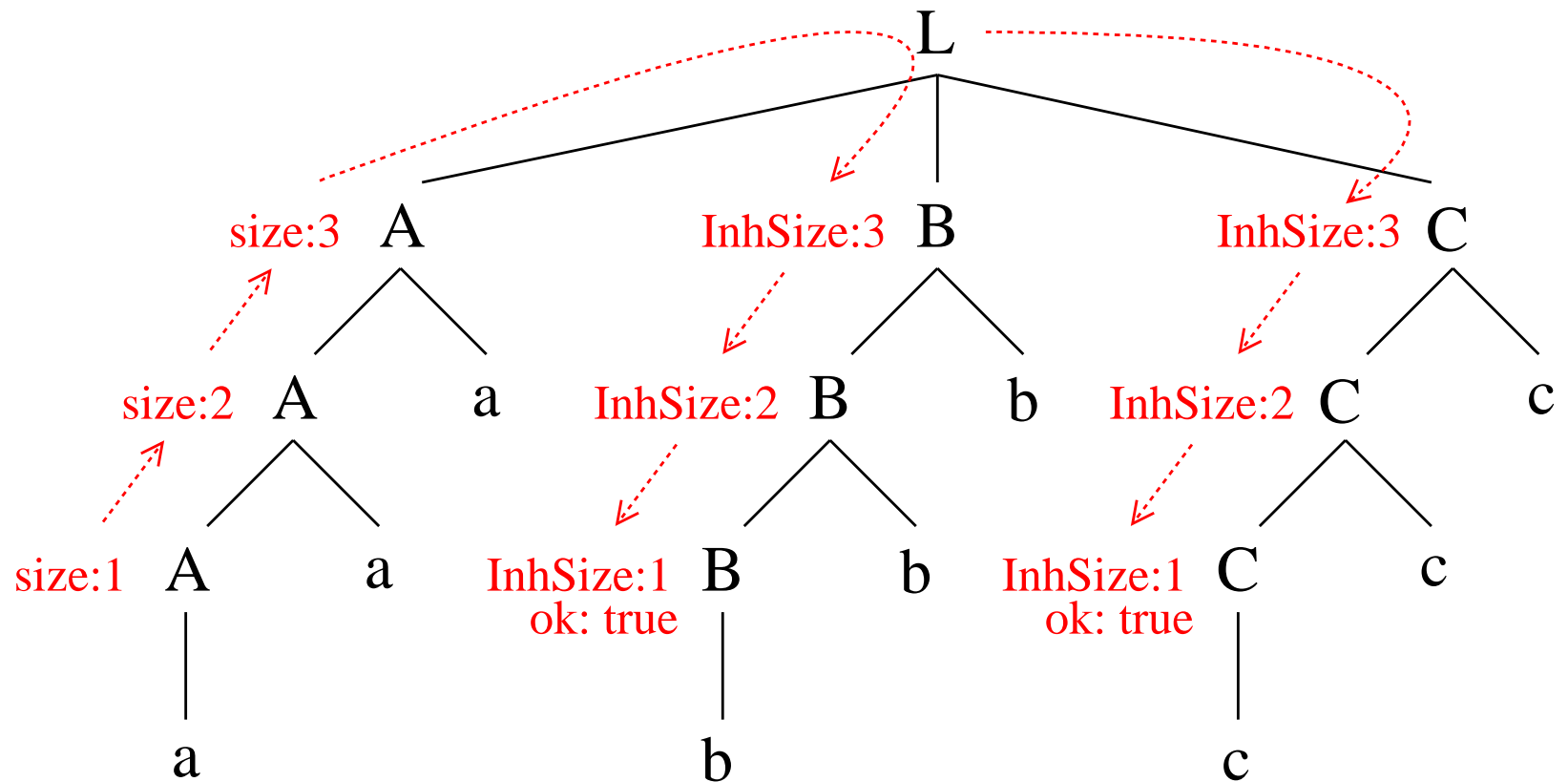
```
L : A B C { B.InhSize = A.size
           C.InhSize = A.size }

A : a      { A.size = 1 }
  | A' a   { A.size = A'.size + 1 }

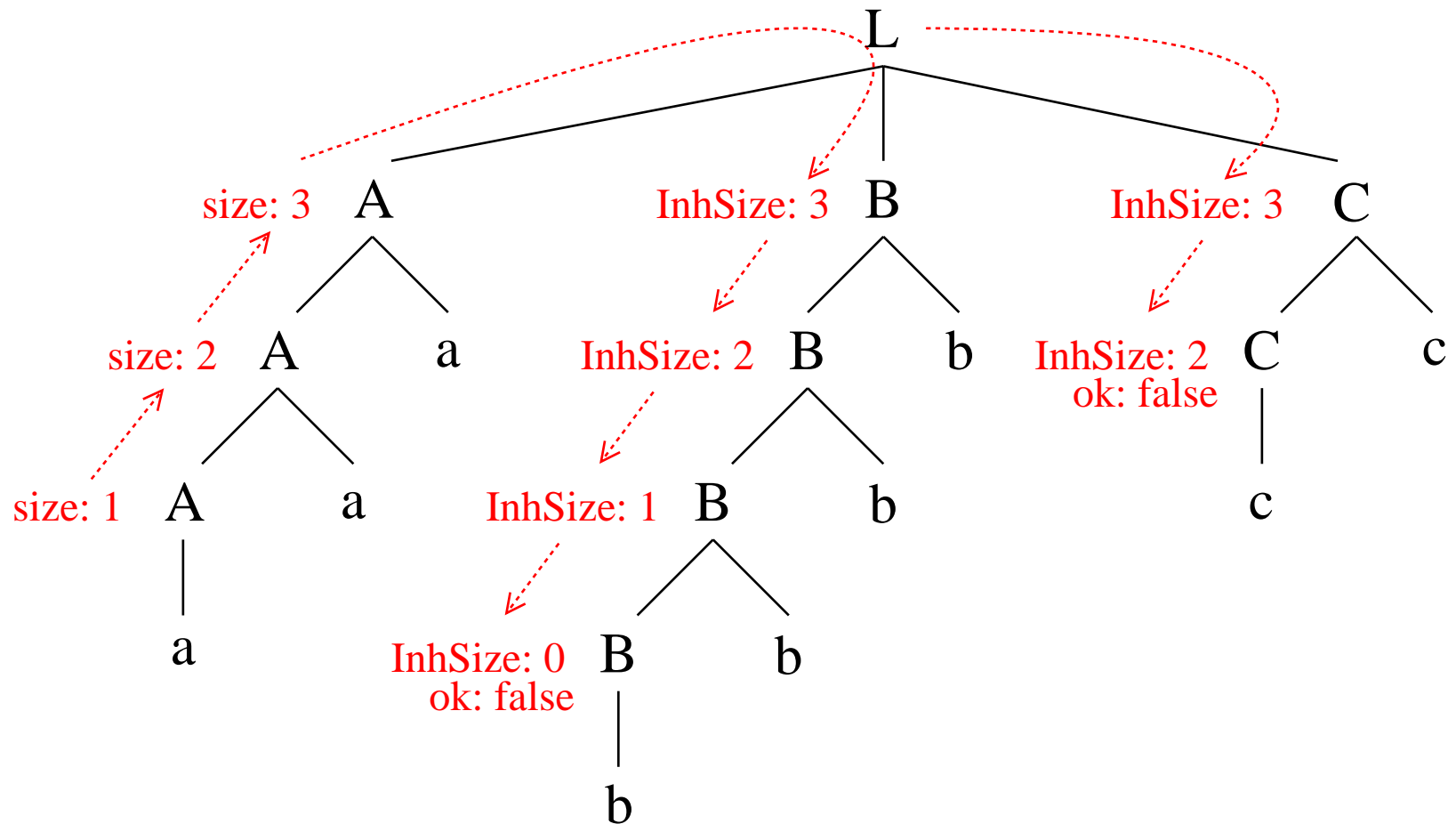
B : b      { B.ok = (B.InhSize == 1) }
  | B' b   { B'.InhSize = B.InhSize - 1 }

C : c      { C.ok = (C.InhSize == 1) }
  | C' c   { C'.InhSize = C.InhSize - 1 }
```

Inherited Attributes



Inherited Attributes



A simple calculator

PRODUCTION	SEMANTIC RULES
$L \rightarrow E$	$L.val = E.val$
$E \rightarrow E_1 + T$	$E.val = E_1.val + T.val$
$E \rightarrow T$	$E.val = T.val$
$T \rightarrow T_1 * F$	$T.val = T_1.val * F.val$
$T \rightarrow F$	$T.val = F.val$
$F \rightarrow (E)$	$F.val = E.val$
$F \rightarrow \mathbf{digit}$	$F.val = \mathbf{digit.lexval}$

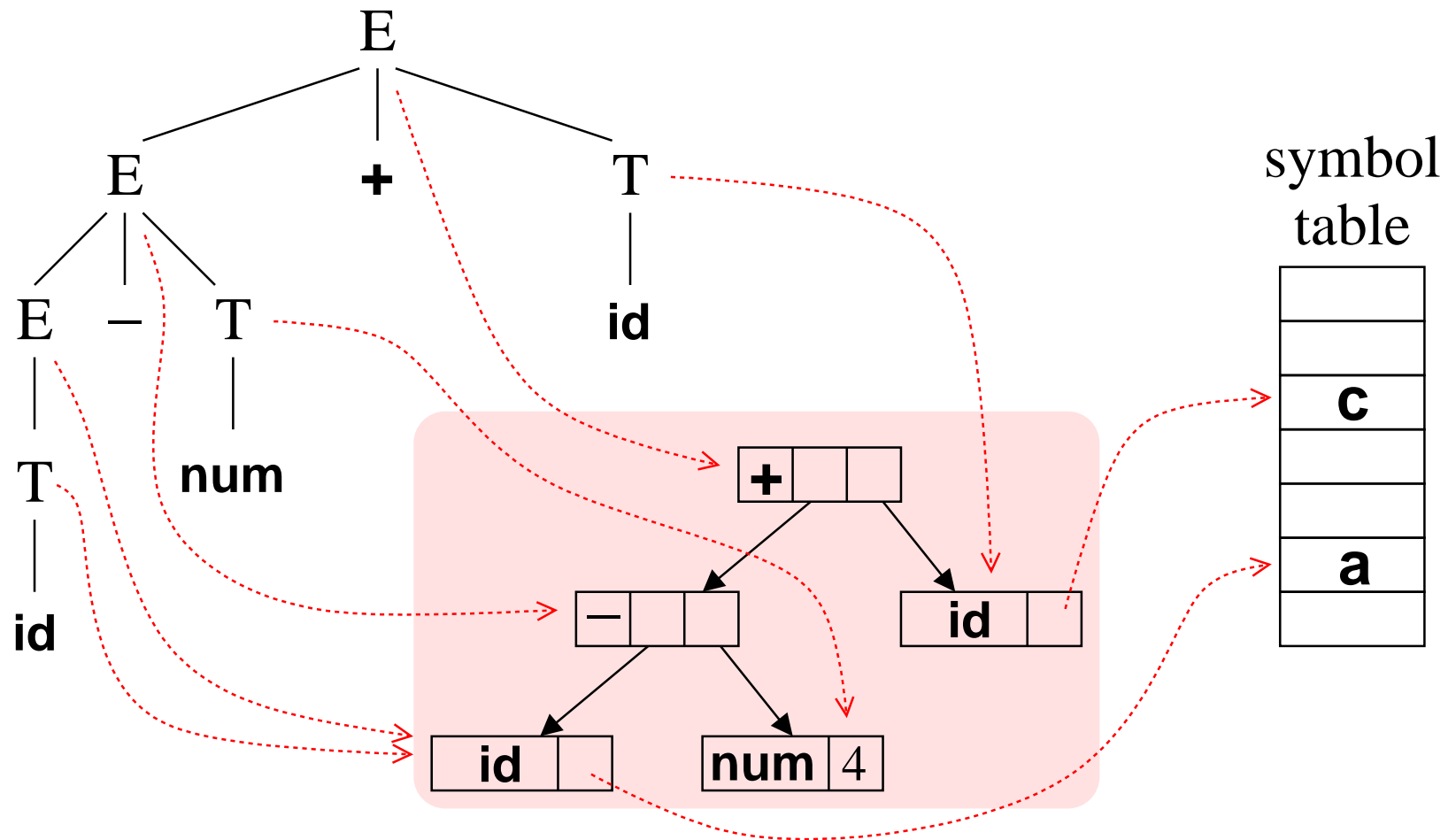
Lexical analyzers provide synthesized attributes (*lexval*) with the values of the terminal symbols.

Constructing syntax trees

PRODUCTION	SEMANTIC RULES
$E \rightarrow E_1 + T$	$E.n = \mathbf{new\ Node}('+', E_1.n + T.n)$
$E \rightarrow E_1 - T$	$E.n = \mathbf{new\ Node}('-', E_1.n + T.n)$
$E \rightarrow T$	$E.n = T.n$
$T \rightarrow (E)$	$T.n = E.n$
$T \rightarrow \mathbf{id}$	$T.n = \mathbf{new\ Leaf}(\mathbf{id}, \mathbf{id.entry})$
$T \rightarrow \mathbf{num}$	$T.n = \mathbf{new\ Leaf}(\mathbf{num}, \mathbf{num.val})$

Constructing syntax trees

Parse and syntax tree for: **a - 4 + c**



Abstract Syntax Trees

Parsing and Syntax Trees

Parsing decides if the program is part of the language.

Not that useful: we want more than a yes/no answer.

Like most, ANTLR parsers can include *actions*: pieces of code that run when a rule is matched.

Top-down parsers: actions executed during parsing rules.

Bottom-up parsers: actions executed when rule is “reduced.”

Actions

Simple languages can be interpreted with parser actions.

```
class CalcParser extends Parser;
```

```
expr returns [int r] { int a; r=0; }  
  : r=mexpr ("+" a=mexpr { r += a; } ) * EOF ;
```

```
mexpr returns [int r] { int a; r=0; }  
  : r=atom ("*" a=atom { r *= a; } ) * ;
```

```
atom returns [int r] { r=0; }  
  : i:INT  
  { r = Integer.parseInt(i.getText()); } ;
```

Actions

In a top-down parser, actions are executed during the matching routines.

Actions can appear anywhere within a rule: before, during, or after a match.

```
rule { /* before */  
    : A { /* during */ } B  
    | C D { /* after */ } ;
```

Bottom-up parsers restricted to running actions only after a rule has matched.

Implementing Actions

Nice thing about top-down parsing: grammar is essentially imperative.

Action code simply interleaved with rule-matching.

Easy to understand what happens when.

Implementing Actions

```
expr returns [int r] { int a; r=0; }  
  : r=mexpr ("+" a=mexpr { r += a; } ) * EOF ;  
  
public final int expr() {           // What ANTLR builds  
    int r; int a; r=0;  
    r=mexpr();  
    while ((LA(1)==PLUS)) {         // ( ) *  
        match(PLUS);               // "+"  
        a=mexpr();                 // a=mexpr  
        r += a;                     // { r += a; }  
    }  
    match(Token.EOF_TYPE);  
    return r;  
}
```

Actions

Usually, actions build a data structure that represents the program.

Separates parsing from translation.

Makes modification easier by minimizing interactions.

Allows parts of the program to be analyzed in different orders.

Actions

Bottom-up parsers can only build bottom-up data structures.

Children known first, parents later.

→ Constructor for any object can require knowledge of children, but not of parent.

Context of an object only established later.

Top-down parsers can build both kinds of data structures.

What To Build?

Typically, an Abstract Syntax Tree that represents the program.

Represents the syntax of the program almost exactly, but easier for later passes to deal with.

Punctuation, whitespace, other irrelevant details omitted.

Abstract vs. Concrete Trees

Like scanning and parsing, objective is to discard irrelevant details.

E.g., comma-separated lists are nice syntactically, but later stages probably just want lists.

AST structure almost a direct translation of the grammar.

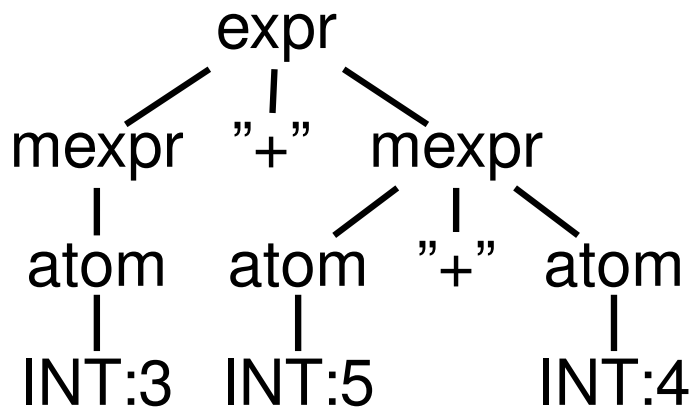
Abstract vs. Concrete Trees

`expr : mexpr ("+" mexpr) * ;`

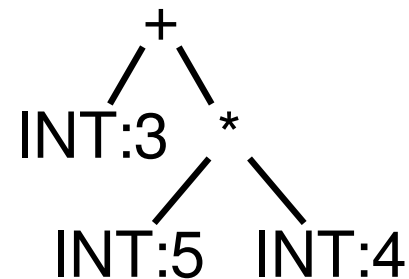
`mexpr : atom ("*" atom) * ;`

`atom : INT ;`

`3 + 5 * 4`



Concrete Parse Tree

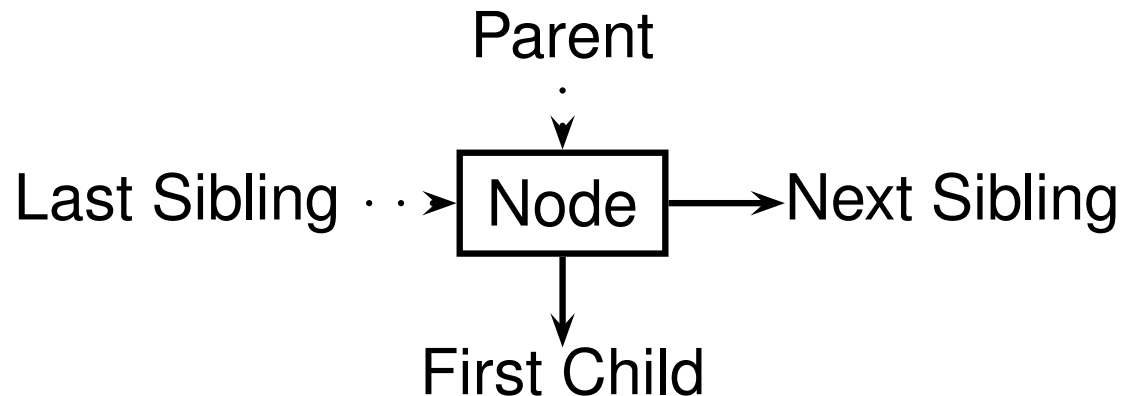


Abstract Syntax Tree

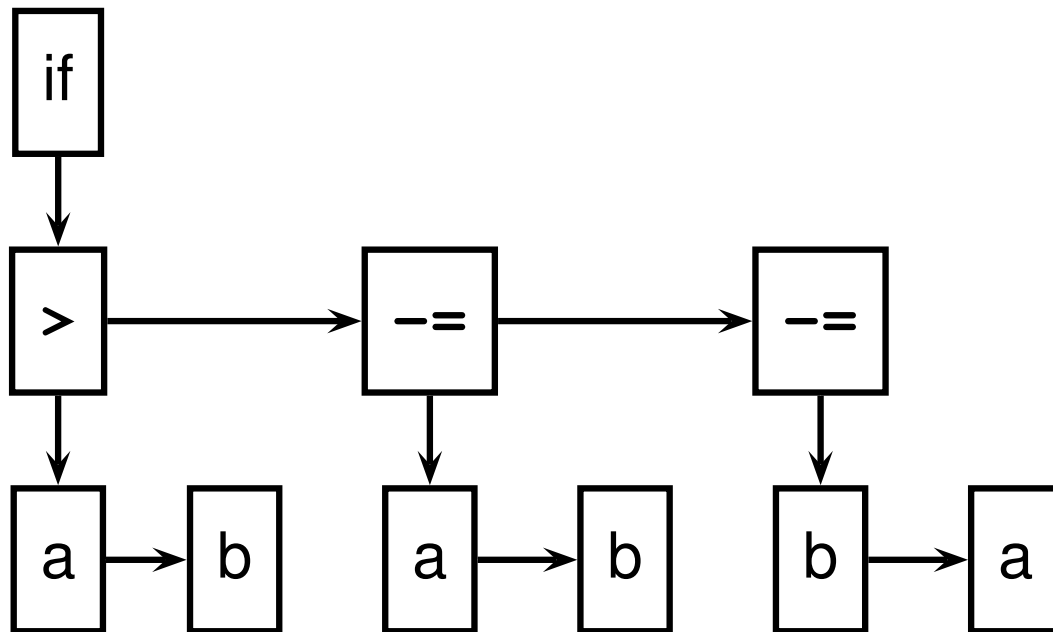
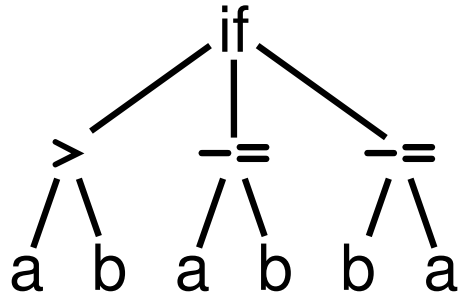
Implementing ASTs

Most general implementation: ASTs are n -ary trees.

Each node holds a token and pointers to its first child and next sibling:

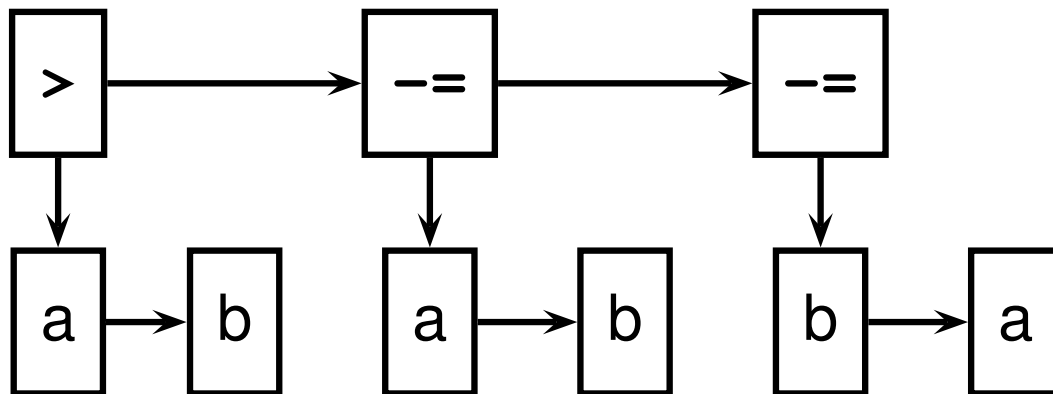
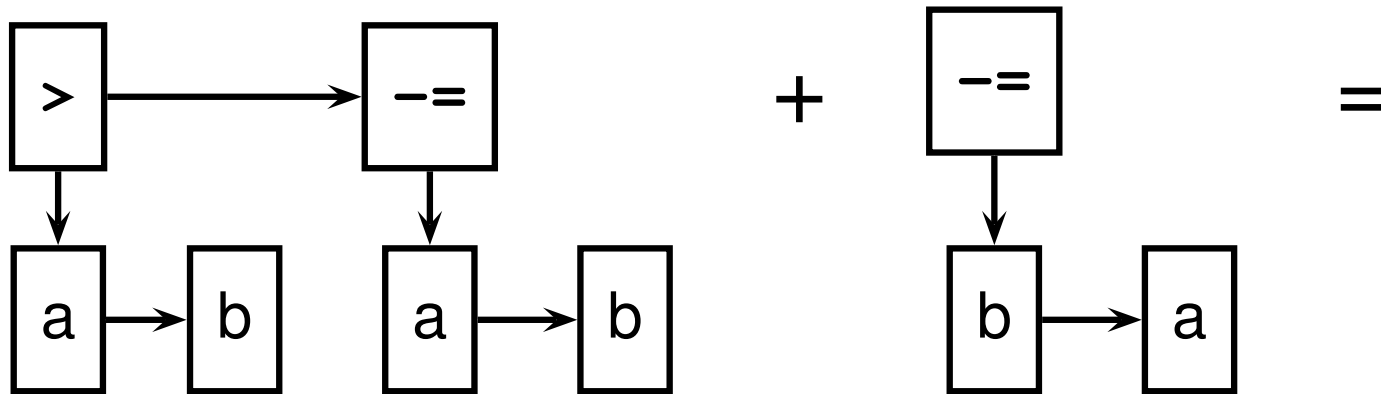


Example of AST structure



Typical AST Operations

Create a new node; Append a subtree as a child.



Comment on Generic ASTs

Is this general-purpose structure too general?

Not very object-oriented: whole program represented with one type.

Alternative: Heterogeneous ASTs: one class per object.

```
class BinOp {  
    int operator; Expr left, right;  
};  
class IfThen {  
    Expr predicate; Stmt thenPart, elsePart;  
};
```

Heterogeneous ASTs

Advantage: avoid switch statements when walking tree.

Disadvantage: each analysis requires another method.

```
class BinOp {  
    int operator; Expr left, right;  
    void typeCheck() { ... };  
    void constantProp() { ... };  
    void buildThreeAddr() { ... };  
};
```

Analyses spread out across class files.

Classes become littered with analysis code, additional annotations.

Comment on Generic ASTs

ANTLR offers a compromise:

It can automatically generate tree-walking code.

→ It generates the big switch statement.

Each analysis can have its own file.

Still have to modify each analysis if the AST changes.

→ Choose the AST structure carefully.

Building ASTs

The Obvious Way to Build ASTs

```
class ASTNode {  
    ASTNode( Token t ) { ... }  
    void appendChild( ASTNode c ) { ... }  
    void appendSibling( ASTNode C) { ... }  
}
```

```
stmt returns [ASTNode n]  
: 'if' p=expr 'then' t=stmt 'else' e=stmt  
  { n = new ASTNode(new Token("IF"));  
    n.appendChild(p);  
    n.appendChild(t);  
    n.appendChild(e); } ;
```

The Obvious Way

Putting code in actions that builds ASTs is traditional and works just fine.

But it's tedious.

Fortunately, ANTLR can automate this process.

Building an AST Automatically with ANTLR

```
class TigerParser extends Parser;  
options {  
    buildAST=true;  
}
```

By default, each matched token becomes an AST node.

Each matched token or rule is made a sibling of the AST for the rule.

After a token, ^ makes the node a root of a subtree.

After a token, ! prevents an AST node from being built.

Automatic AST Construction

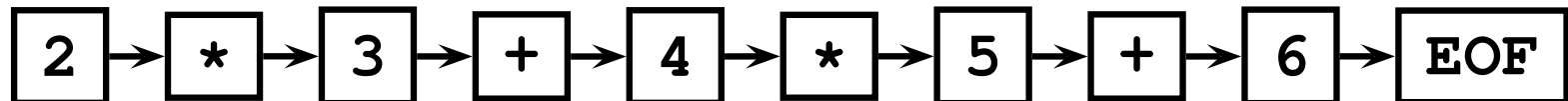
Running

```
class CalcParser extends Parser;  
  options { buildAST=true; }  
expr : mexpr ('+' mexpr)* EOF ;  
mexpr : atom ('*' atom)* ;  
atom : INT ;
```

on

2*3+4*5+6

gives



AST Construction with Annotations

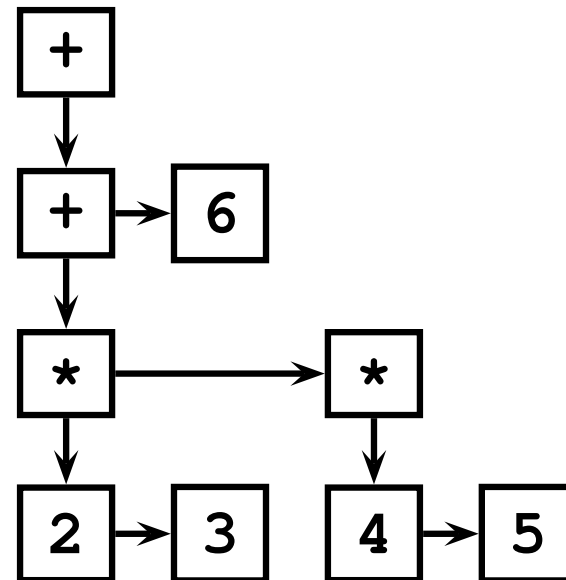
Running

```
class CalcParser extends Parser;  
  options { buildAST=true; }  
expr : mexpr ('+' ^ mexpr)* EOF! ;  
mexpr : atom ('*' ^ atom)* ;  
atom : INT ;
```

on

$2 * 3 + 4 * 5 + 6$

gives



Choosing AST Structure

Designing an AST Structure

Sequences of things

Removing unnecessary punctuation

Additional grouping

How many token types?

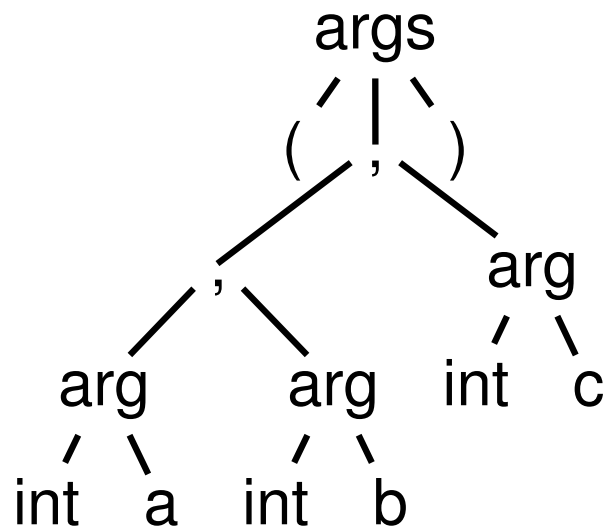
Sequences of Things

Comma-separated lists are common

```
int gcd(int a, int b, int c)
```

```
args : "(" ( arg ("," arg)* )? ")" ;
```

A concrete parse tree:



Drawbacks:

Many unnecessary nodes

Branching suggests recursion

Harder for later routines to get the data they want

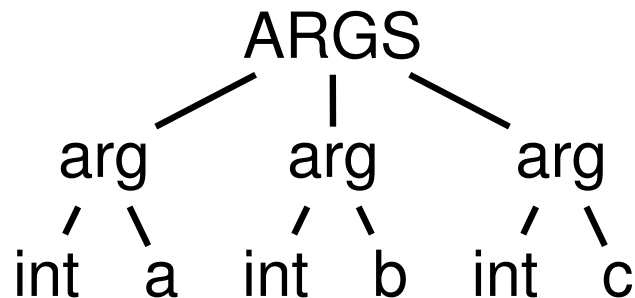
Sequences of Things

Better to choose a simpler structure for the tree.

Punctuation irrelevant; build a simple list.

```
int gcd(int a, int b, int c)
```

```
args : "(" ! ( arg ("," ! arg)* ) ? ")" !  
      { #args = #([ARGS], args); } ;
```



What's going on here?

```
args : "(" ! ( arg ( "," ! arg ) * ) ? ")" !  
      { #args = #([ARGS], args); } ;
```

Rule generates a sequence of arg nodes.

Node generation suppressed for punctuation (parens, commas).

Action uses ANTLR's terse syntax for building trees.

```
{ #args = # ( [ARGS] , args ) ; } ;
```

```
graph TD
    A["#args"] --- B["set the args tree to a new tree whose root is a node of type ARGS and whose child is the old args tree"]
    C["# ("] --- B
    D["[ARGS]"] --- B
    E["args"] --- B
    F[") ; } ;"] --- B
```

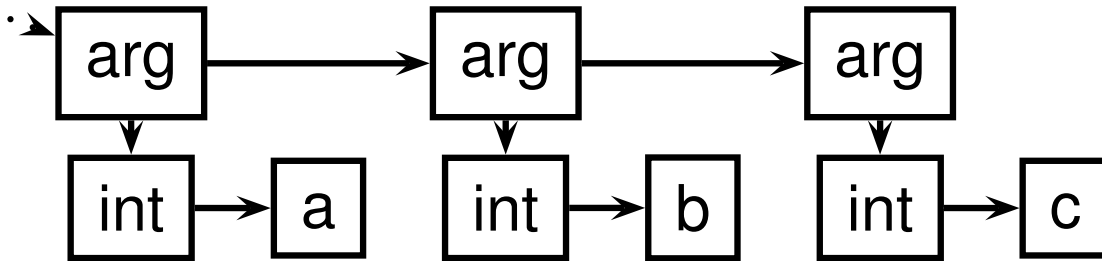
“set the args tree to a new tree whose root is a node of type ARGS and whose child is the old args tree”

What's going on here?

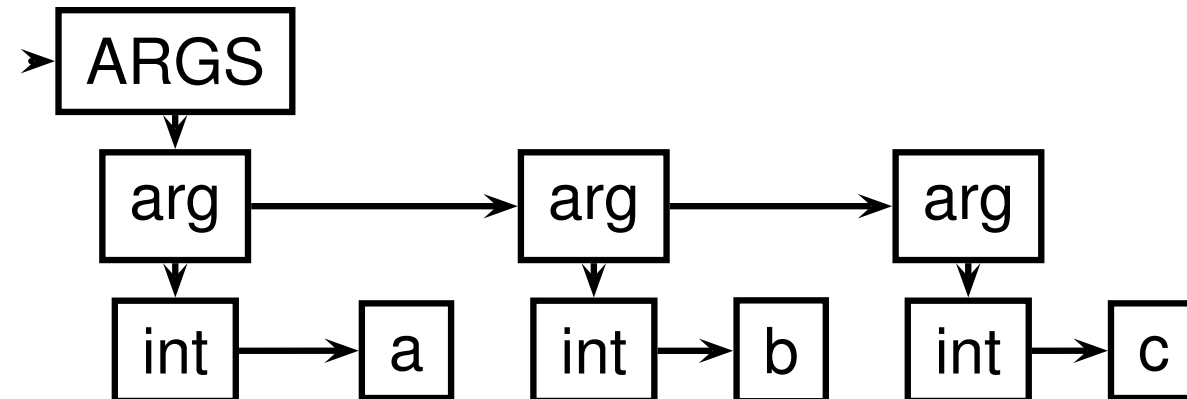
```
(int a, int b, int c)
```

```
args : "(" ! ( arg ("," ! arg)* ) ? ")" !  
      { #args = #([ARGS], args); } ;
```

#args



#args



Removing Unnecessary Punctuation

Punctuation makes the syntax readable, unambiguous.

Information represented by structure of the AST

Things typically omitted from an AST

- Parentheses
Grouping and precedence/associativity overrides
- Separators (commas, semicolons)
Mark divisions between phrases
- Extra keywords
while-do, if-then-else (one is enough)

Additional Grouping

The Tiger language from Appel's book allows mutually recursive definitions only in uninterrupted sequences:

```
let
```

```
    function f1() = ( f2() ) /* OK */
```

```
    function f2() = ( ... )
```

```
in ... end
```

```
let
```

```
    function f1() = ( f2() ) /* Error */
```

```
    var foo := 42 /* splits group */
```

```
    function f2() = ( ... )
```

```
in ... end
```

Grouping

Convenient to group sequences of definitions in the AST.
Simplifies later static semantic checks.

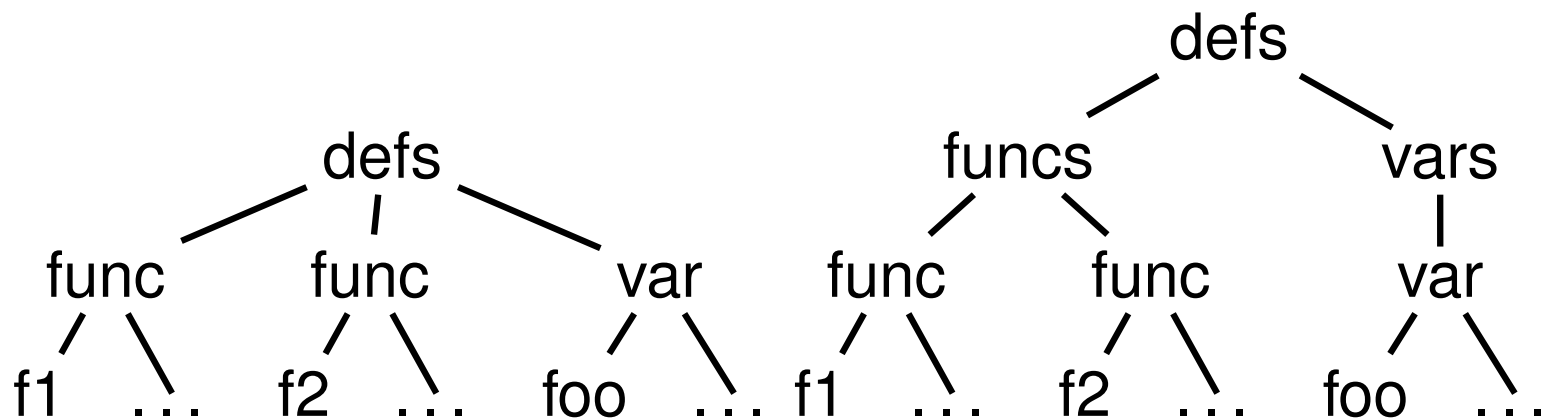
```
let
```

```
    function f1() = ( ... )
```

```
    function f2() = ( ... )
```

```
    var foo := 42
```

```
in ... end
```



Grouping

Identifying and building sequences of definitions a little tricky in ANTLR.

Obvious rules

```
defs    : ( funcs | vars | types ) * ;
```

```
funcs   : ( func ) + ;
```

```
vars    : ( var  ) + ;
```

```
types   : ( type ) + ;
```

are ambiguous: Maximum-length sequences or minimum-length sequences?

Grouping

Hint: Use ANTLR's **greedy** option to disambiguate this.

The greedy flag decides whether repeating a rule takes precedence when an outer rule could also work.

```
string : (dots)* ;  
dots  : ("." )+ ;
```

When faced with a period, the second rule can repeat itself or exit.

The Greedy Option

Setting greedy true makes “dots” as long as possible

```
string : (dots)* ;  
dots : ( options greedy=true; : ".")+ ;
```

Setting greedy false makes each “dots” a single period

```
string : (dots)* ;  
dots : ( options greedy=false; : ".")+ ;
```

How Many Types of Tokens?

Since each token is a type plus some text, there is some choice.

Generally, want each “different” construct to have a different token type.

Different types make sense when each needs different analysis.

Arithmetic operators usually not that different.

For the assignment, you need to build a node of type “BINOP” for every binary operator. The text indicates the actual operator.

Exercise

Write an attribute grammar to specify binary numerals. A binary numeral is a non-empty sequence of binary digits followed by a period and another non-empty sequence of binary digits. The grammar must calculate an attribute that returns the value of the binary numeral (a real number).

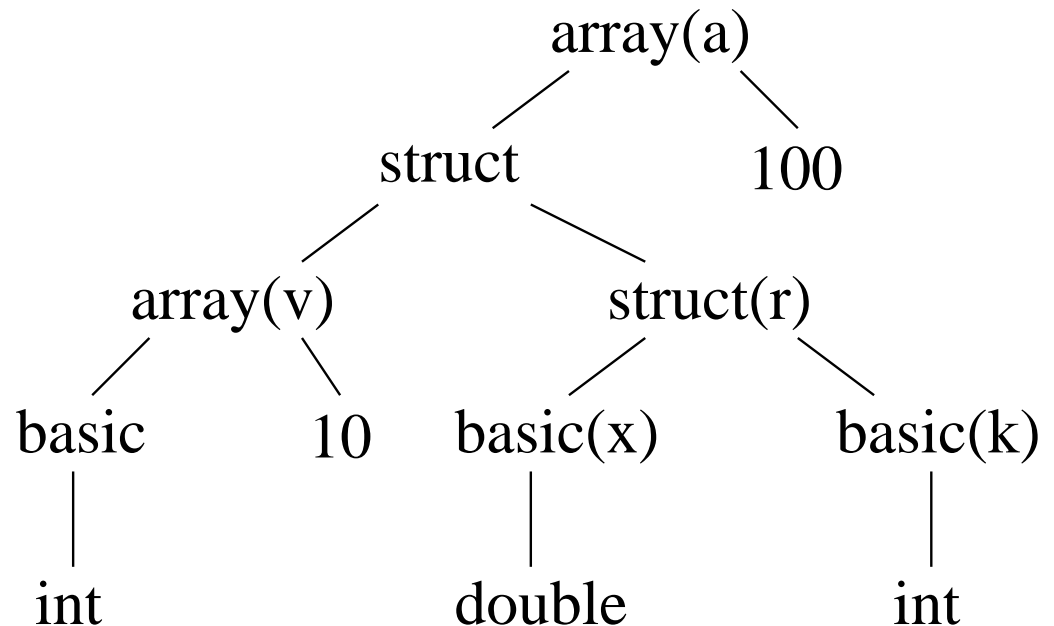
Use the following context-free grammar:

$$Num \rightarrow Digits '.' Digits$$
$$Digits \rightarrow Bit \mid Digits Bit$$
$$Bit \rightarrow 0 \mid 1$$

Exercise

Write an attribute grammar to construct a syntax tree for data types with arrays and structs (see example). Assume that only **int** and **double** basic types are used.

```
struct {  
    int [10] v;  
    struct {  
        double x;  
        int k;  
    } r;  
} [100] a;
```



Synthesize an attribute that indicates the size of the data structure, assuming that an **int** takes 4 bytes and a **double** takes 8 bytes.