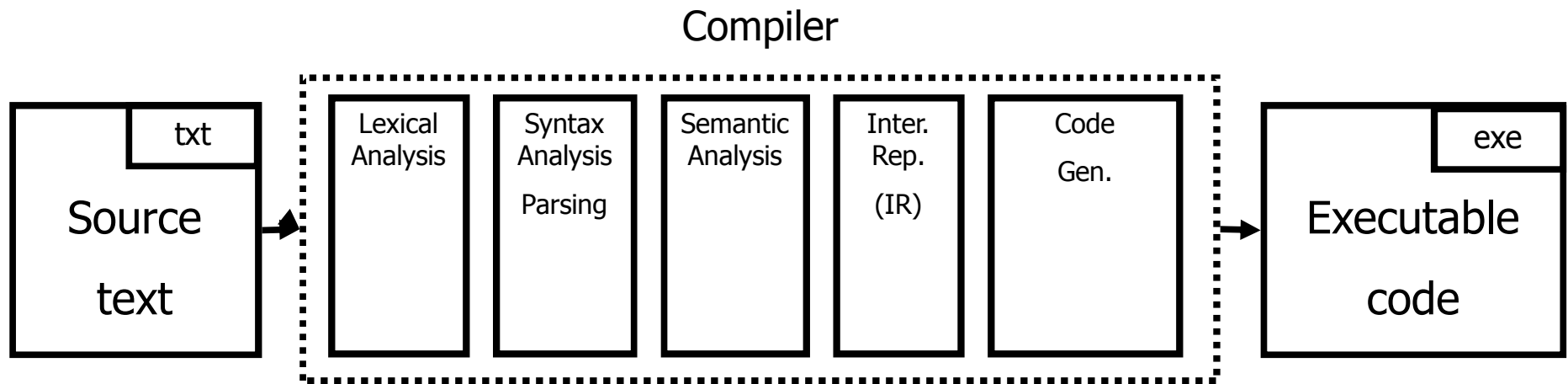


Lecture 07 – attribute grammars + intro to IR

THEORY OF COMPILATION

Eran Yahav

You are here



Last Week: Types

- What is a type?
 - Simplest answer: a set of values
 - Integers, real numbers, booleans, ...
- Why do we care?
 - Safety
 - Guarantee that certain errors cannot occur at runtime
 - Abstraction
 - Hide implementation details
 - Documentation
 - Optimization

Last Week: Type System

- A type system of a programming language is a way to define how “good” program behave
 - Good programs = well-typed programs
 - Bad programs = not well typed
- Type checking
 - Static typing – most checking at compile time
 - Dynamic typing – most checking at runtime
- Type inference
 - Automatically infer types for a program (or show that there is no valid typing)

Strongly vs. weakly typed

- Coercion
- Strongly typed
 - C, C++, Java
- Weakly typed
 - Perl, PHP
- (YMMV, not everybody agrees on this classification)

Output: 73

warning: initialization makes integer from pointer without a cast

perl

```
$a=31;  
$b="42x";  
$c=$a+$b;  
print $c;
```

C

```
main() {  
    int a=31;  
    char b[3]="42x";  
    int c=a+b;  
}
```

error: Incompatible type for declaration. Can't convert java.lang.String to int

Java

```
public class... {  
    public static void main() {  
        int a=31;  
        String b ="42x";  
        int c=a+b;  
    }  
}
```

Last week: how does this magic happen?

- We probably need to go over the AST?
- how does this relate to the clean formalism of the parser?

Syntax Directed Translation

- The parse tree (syntax) is used to drive the translation
- Semantic attributes
 - Attributes attached to grammar symbols
- Semantic actions
 - How to update the attributes when a production is used in a derivation
- Attribute grammars

Attribute grammars

- Attributes
 - Every grammar symbol has attached attributes
 - Example: Expr.type
- Semantic actions
 - Every production rule can define how to assign values to attributes
 - Example:
Expr \rightarrow Expr + Term
Expr.type = Expr1.type when (Expr1.type == Term.type)
Error otherwise

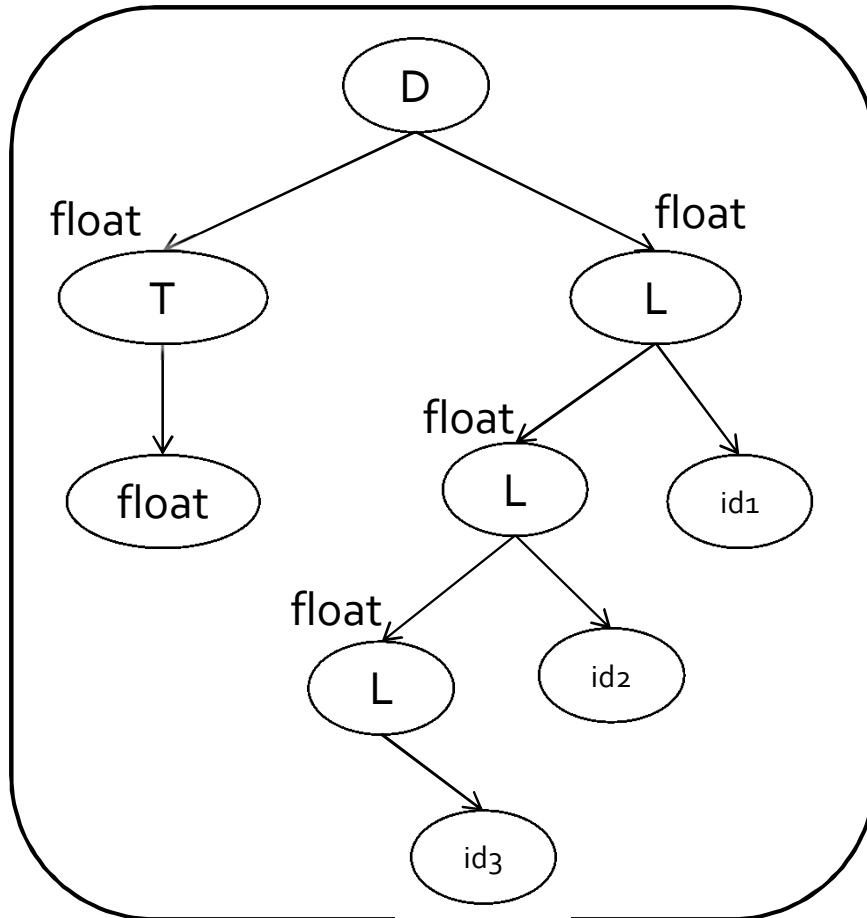
Indexed symbols

- Add indexes to distinguish repeated grammar symbols
- Does not affect grammar
- Used in semantic actions

- $\text{Expr} \rightarrow \text{Expr} + \text{Term}$
Becomes
 $\text{Expr} \rightarrow \text{Expr}_1 + \text{Term}$

Example

float x,y,z



Production	Semantic Rule
$D \rightarrow T L$	$L.in = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$L_1.in = L.in$ $addType(id.entry, L.in)$
$L \rightarrow id$	$addType(id.entry, L.in)$

Attribute Evaluation

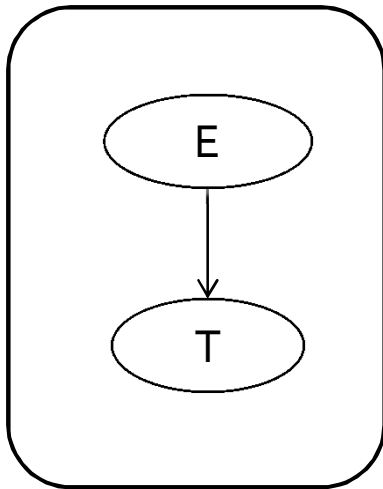
- Build the AST
- Fill attributes of terminals with values derived from their representation
- Execute evaluation rules of the nodes to assign values until no new values can be assigned
 - In the right order such that
 - No attribute value is used before its available
 - Each attribute will get a value only once

Dependencies

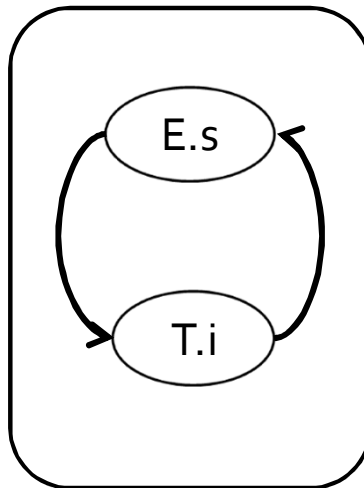
- A semantic equation $a = b_1, \dots, b_m$ requires computation of b_1, \dots, b_m to determine the value of a
- The value of a depends on b_1, \dots, b_m
 - We write $a \leftarrow b_i$

Cycles

- Cycle in the dependence graph
- May not be able to compute attribute values



AST



Dependence
graph

$$\begin{aligned} E.S &= T.i \\ T.i &= E.S + 1 \end{aligned}$$

Attribute Evaluation

- Build the AST
- Build dependency graph
- Compute evaluation order using topological ordering
- Execute evaluation rules based on topological ordering

- Works as long as there are no cycles

Building Dependency Graph

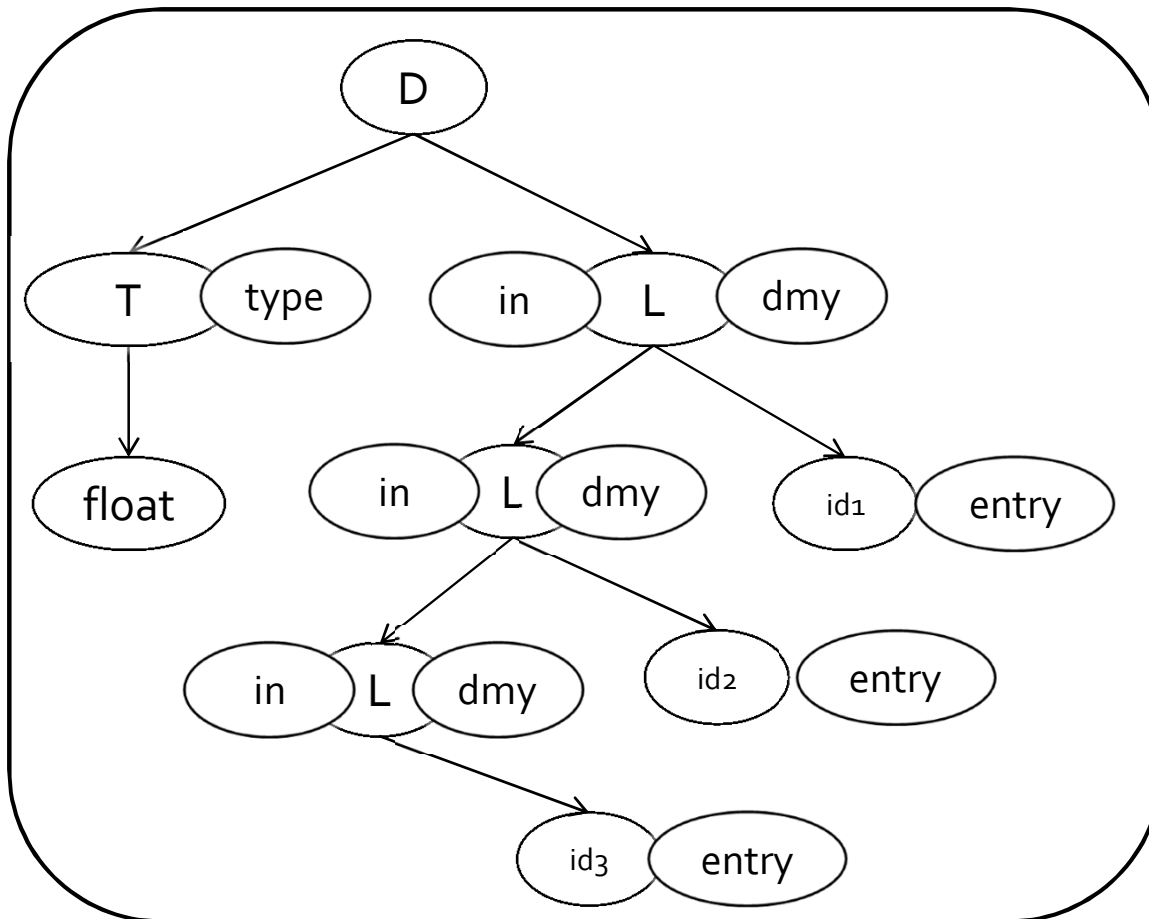
- All semantic equations take the form

$$\text{attr1} = \text{func1}(\text{attr1.1}, \text{attr1.2}, \dots)$$
$$\text{attr2} = \text{func2}(\text{attr2.1}, \text{attr2.2}, \dots)$$

- Actions with side effects use a dummy attribute
- Build a directed dependency graph G
 - For every attribute a of a node n in the AST create a node $n.a$
 - For every node n in the AST and a semantic action of the form $b = f(c_1, c_2, \dots, c_k)$ add edges of the form (c_i, b)

Example

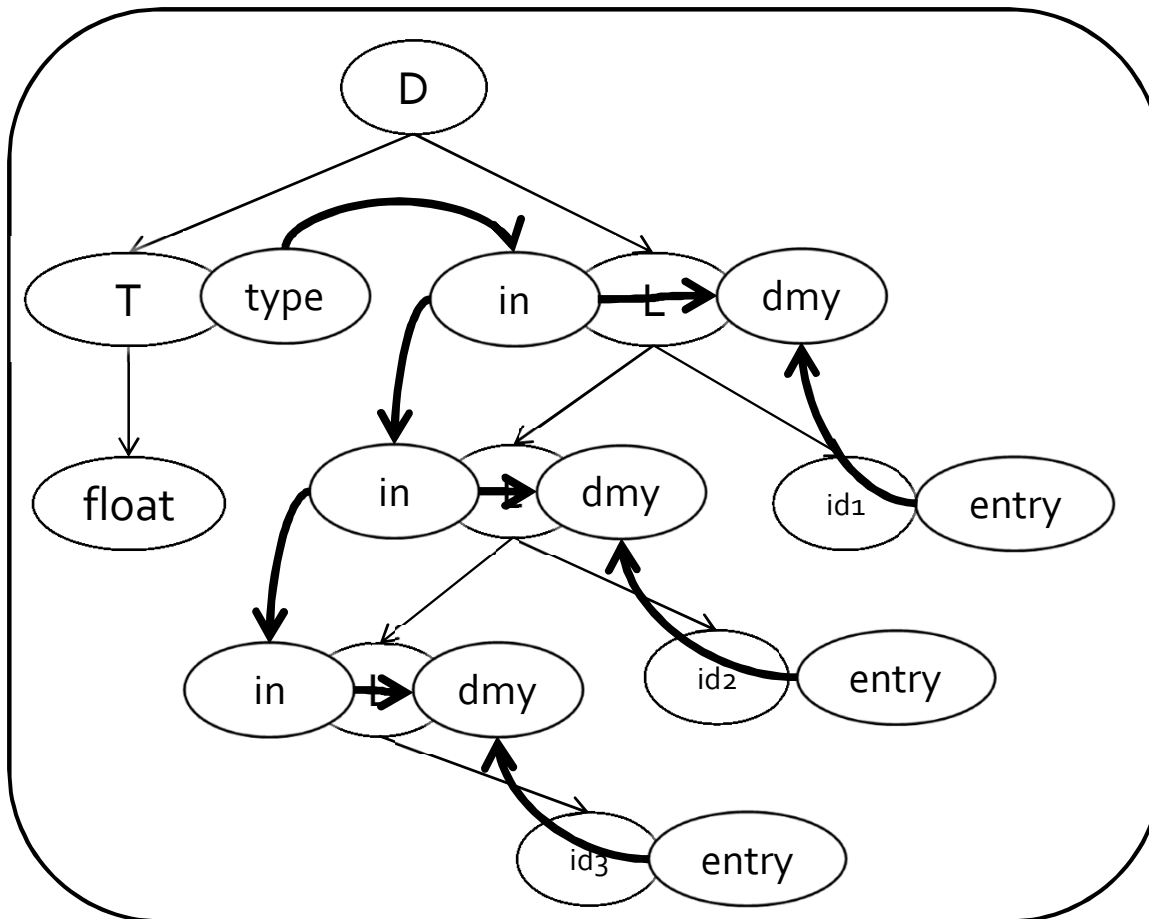
float x,y,z



Prod.	Semantic Rule
$D \rightarrow T L$	$L.in = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$L_1.in = L.in$ $addType(id.entry, L.in)$
$L \rightarrow id$	$addType(id.entry, L.in)$

Example

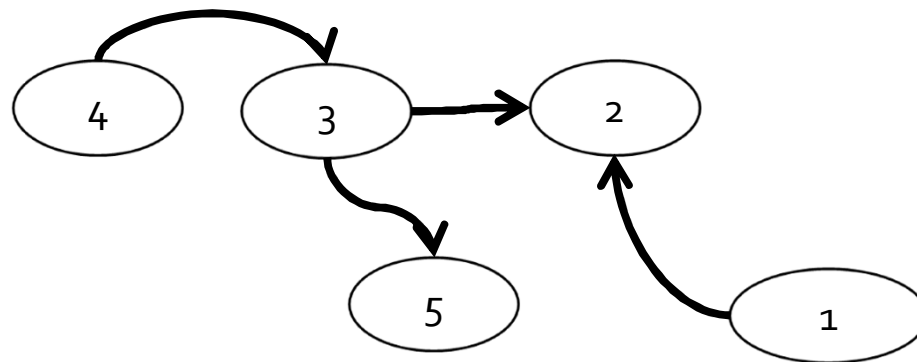
float x,y,z



Prod.	Semantic Rule
$D \rightarrow T L$	$L.in = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$L_1.in = L.in$ $addType(id.entry, L.in)$
$L \rightarrow id$	$addType(id.entry, L.in)$

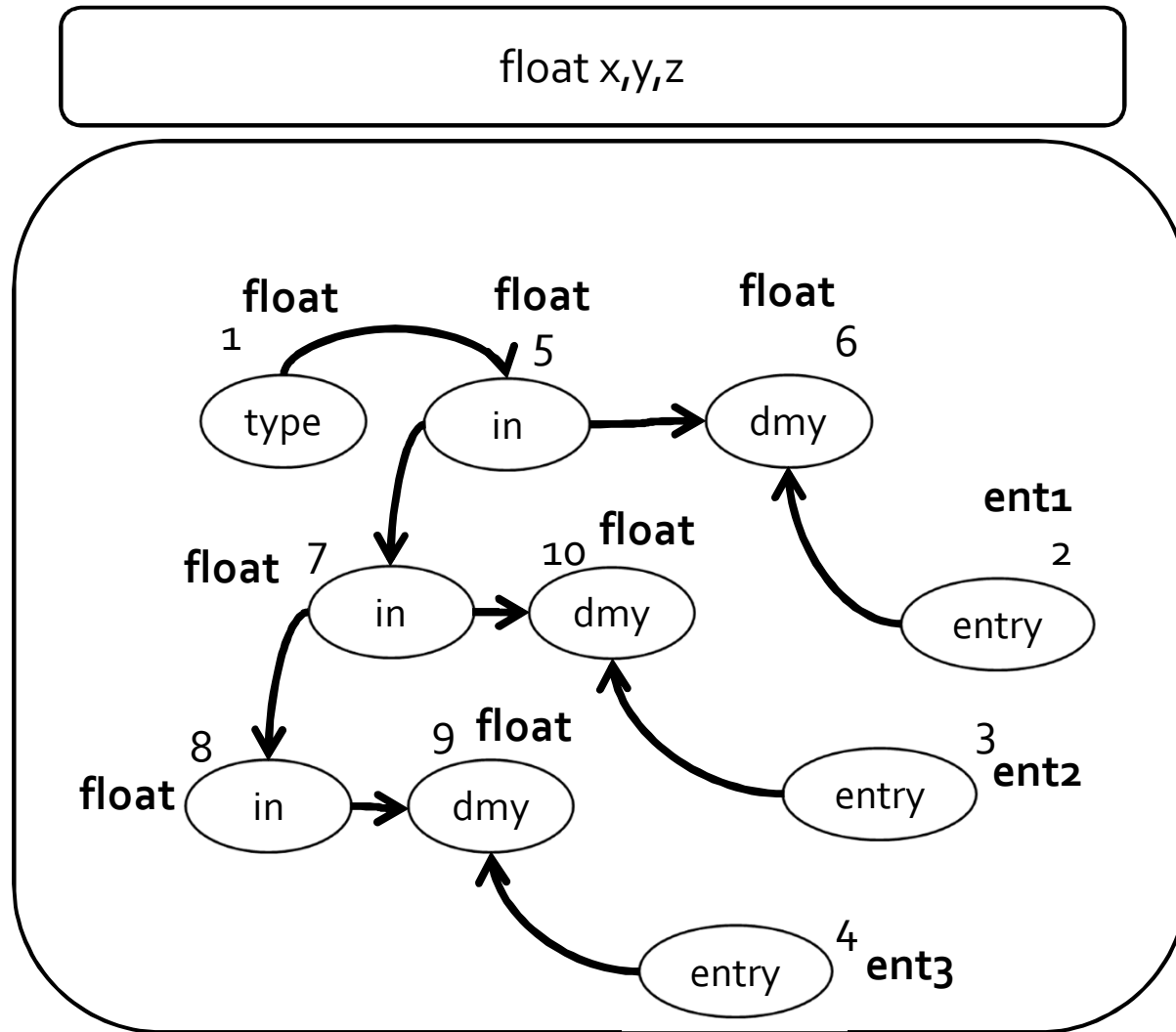
Topological Order

- For a graph $G=(V,E)$, $|V|=k$
- Ordering of the nodes v_1, v_2, \dots, v_k such that for every edge $(v_i, v_j) \in E$, $i < j$



Example topological orderings: 1 4 3 2 5, 4 1 3 5 2

Example



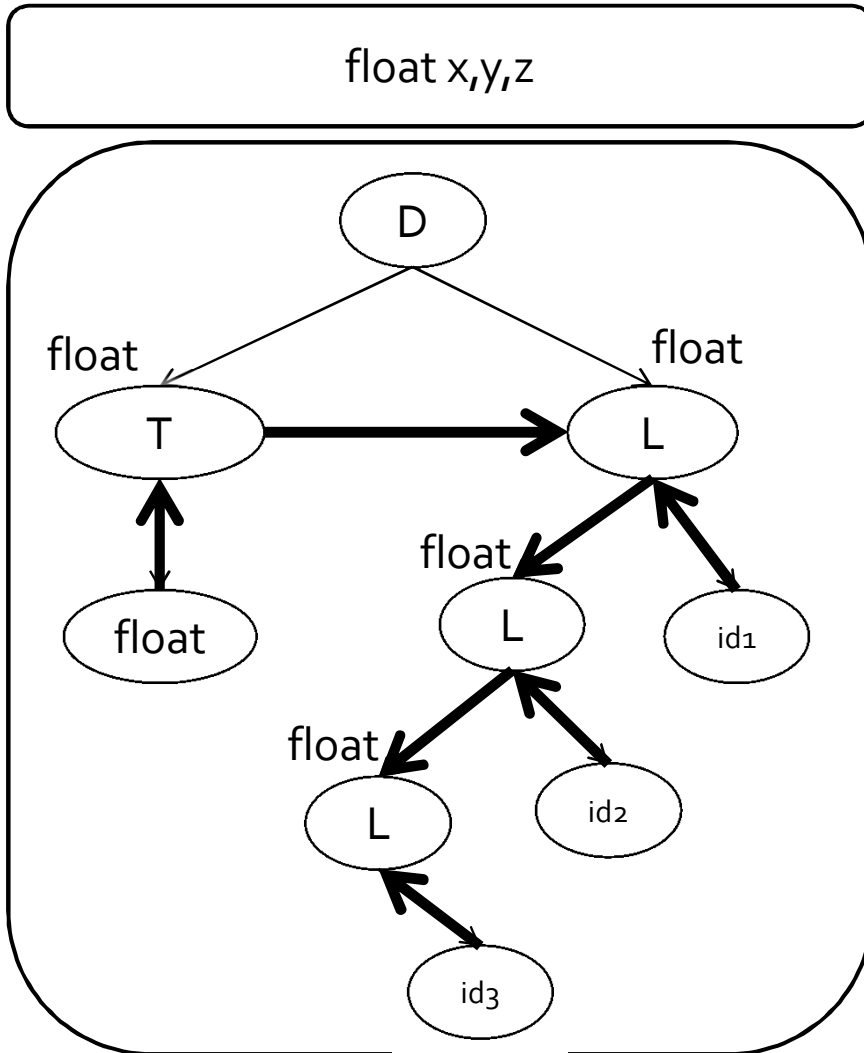
But what about cycles?

- For a given attribute grammar hard to detect if it has cyclic dependencies
 - Exponential cost
- Special classes of attribute grammars
 - Our “usual trick”
 - sacrifice generality for predictable performance

Inherited vs. Synthesized Attributes

- Synthesized attributes
 - Computed from children of a node
- Inherited attributes
 - Computed from parents and siblings of a node
- Attributes of tokens are technically considered as synthesized attributes

example



Production	Semantic Rule
$D \rightarrow T L$	$L.in = T.type$
$T \rightarrow int$	$T.type = integer$
$T \rightarrow float$	$T.type = float$
$L \rightarrow L_1, id$	$L_1.in = L.in$ $addType(id.entry, L.in)$
$L \rightarrow id$	$addType(id.entry, L.in)$

inherited

synthesized

S-attributed Grammars

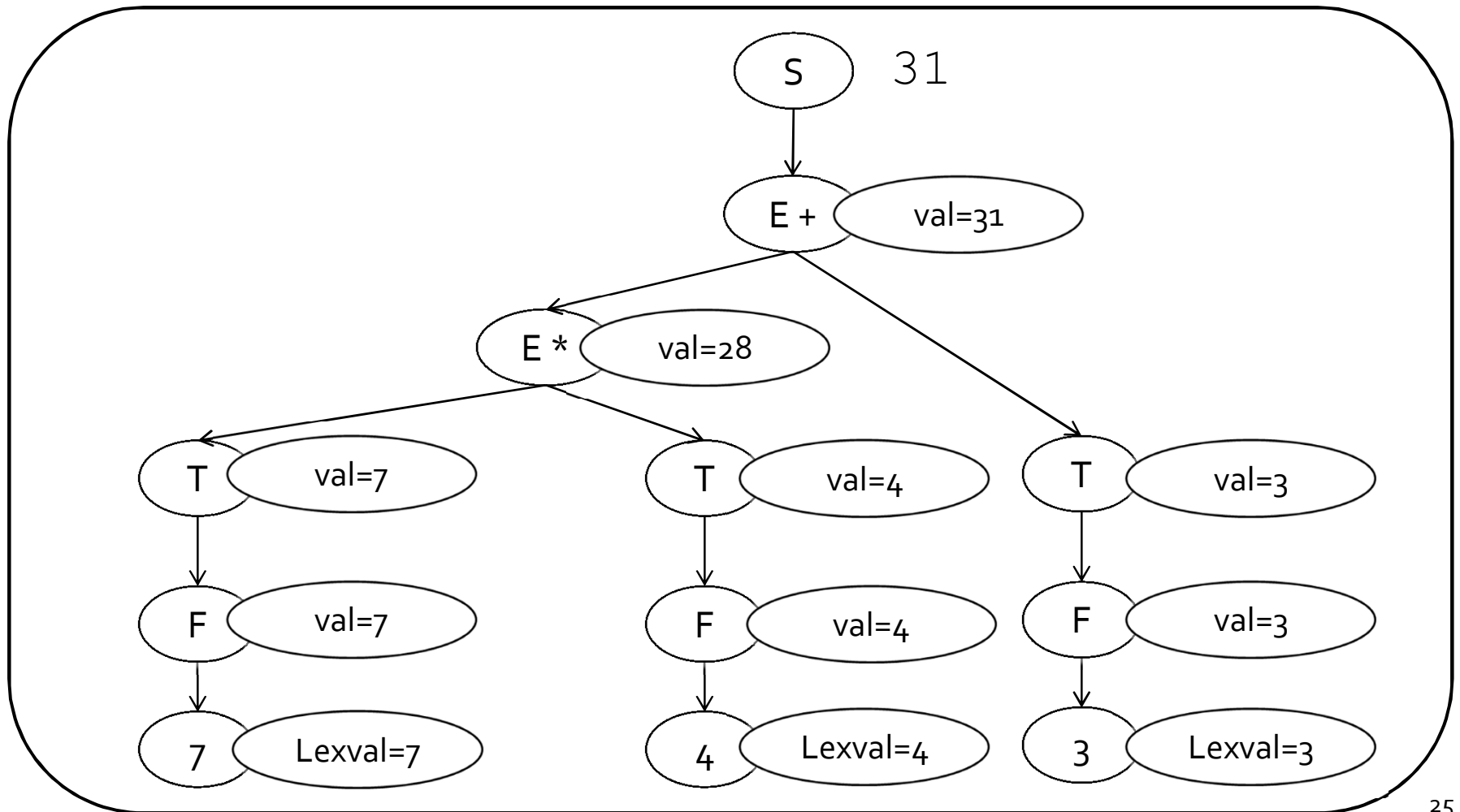
- Special class of attribute grammars
- Only uses synthesized attributes (S-attributed)
- No use of inherited attributes

- Can be computed by any bottom-up parser **during parsing**
- Attributes can be stored on the parsing stack
- Reduce operation computes the (synthesized) attribute from attributes of children

S-attributed Grammar: example

Production	Semantic Rule
$S \rightarrow E ;$	$\text{print}(E.\text{val})$
$E \rightarrow E_1 + T$	$E.\text{val} = E_1.\text{val} + T.\text{val}$
$E \rightarrow T$	$E.\text{val} = T.\text{val}$
$T \rightarrow T_1 * F$	$T.\text{val} = T_1.\text{val} * F.\text{val}$
$T \rightarrow F$	$T.\text{val} = F.\text{val}$
$F \rightarrow (E)$	$F.\text{val} = E.\text{val}$
$F \rightarrow \text{digit}$	$F.\text{val} = \text{digit}.\text{lexval}$

example



L-attributed grammars

- L-attributed attribute grammar when every attribute in a production $A \rightarrow X_1 \dots X_n$ is
 - A synthesized attribute, or
 - An inherited attribute of X_j , $1 \leq j \leq n$ that only depends on
 - Attributes of $X_1 \dots X_{j-1}$ to the left of X_j , or
 - Inherited attributes of A

Example: typesetting



- Vertical geometry
 - pointsize (ps) – size of letters in a box. Subscript text has smaller point size of 0.7p.
 - baseline
 - height (ht) – distance from top of the box to the baseline
 - depth (dp) – distance from baseline to the bottom of the box.

Example: typesetting

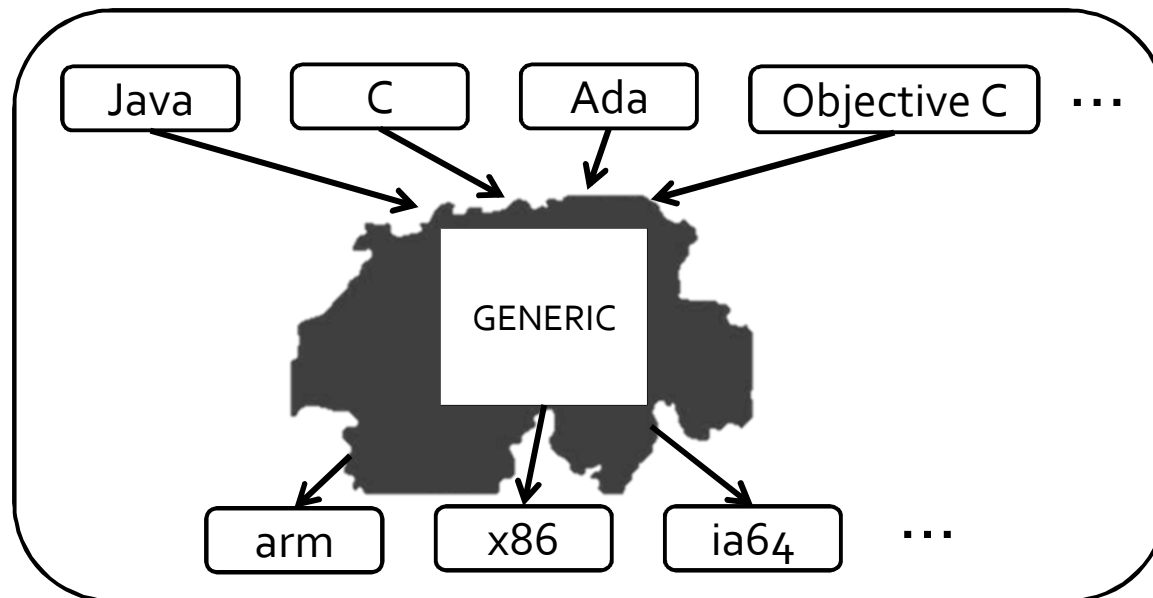
production	semantic rules
$S \rightarrow B$	$B.ps = 10$
$B \rightarrow B_1 B_2$	$B_1.ps = B.ps$ $B_2.ps = B.ps$ $B.ht = \max(B_1.ht, B_2.ht)$ $B.dp = \max(B_1.dp, B_2.dp)$
$B \rightarrow B_1 \text{ sub } B_2$	$B_1.ps = B.ps$ $B_2.ps = 0.7 * B.ps$ $B.ht = \max(B_1.ht, B_2.ht - 0.25 * B.ps)$ $B.dp = \max(B_1.dp, B_2.dp - 0.25 * B.ps)$
$B \rightarrow \text{text}$	$B.ht = \text{getHt}(B.ps, \text{text.lexval})$ $B.dp = \text{getDp}(B.ps, \text{text.lexval})$

Attribute grammars: summary

- Contextual analysis can move information between nodes in the AST
 - Even when they are not “local”
- Attribute grammars
 - Attach attributes and semantic actions to grammar
- Attribute evaluation
 - Build dependency graph, topological sort, evaluate
- Special classes with pre-determined evaluation order: S-attributed, L-attributed

Intermediate Representation

- “neutral” representation between the front-end and the back-end
 - Abstracts away details of the source language
 - Abstract away details of the target language
- A compiler may have multiple intermediate representations and move between them
- In practice, the IR may be biased toward a certain language (e.g., GENERIC in gcc)



Intermediate Representation(s)

- Annotated abstract syntax tree
- Three address code
- ...

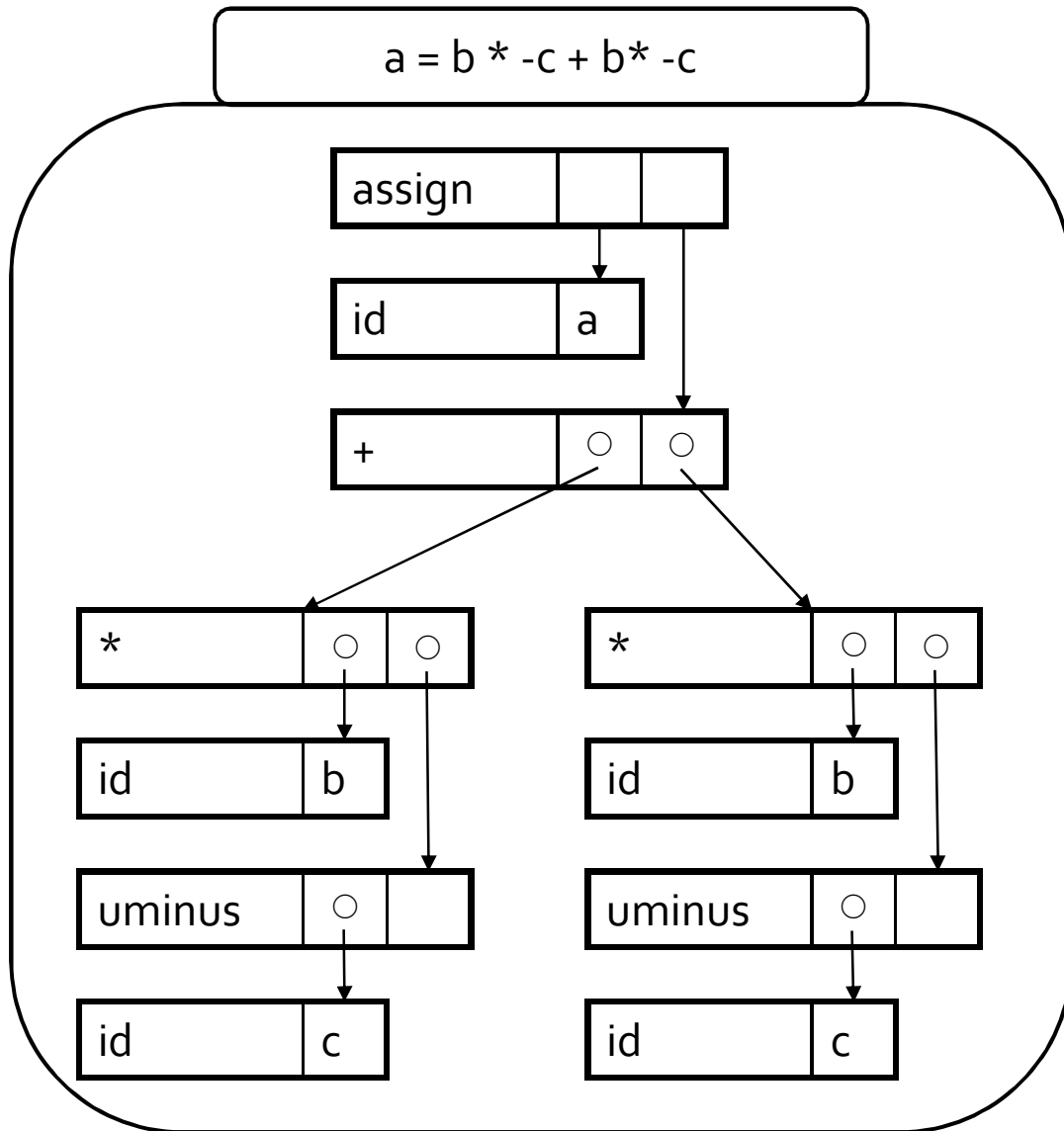
Example: Annotated AST

production	semantic rule
$S \rightarrow id := E$	$S.nptr = \text{makeNode}(\text{'assign'}, \text{makeLeaf}(id, id.place), E.nptr)$
$E \rightarrow E_1 + E_2$	$E.nptr = \text{makeNode}(\text{'+'}, E_1.nptr, E_2.nptr)$
$E \rightarrow E_1 * E_2$	$E.nptr = \text{makeNode}(\text{'*'}, E_1.nptr, E_2.nptr)$
$E \rightarrow -E_1$	$E.nptr = \text{makeNode}(\text{'uminus'}, E_1.nptr)$
$E \rightarrow (E_1)$	$E.nptr = E_1.nptr$
$E \rightarrow id$	$E.nptr = \text{makeLeaf}(id, id.place)$

- `makeNode` – creates new node for unary/binary operator
- `makeLeaf` – creates a leaf
- `id.place` – pointer to symbol table

Example

$a = b * -c + b * -c$

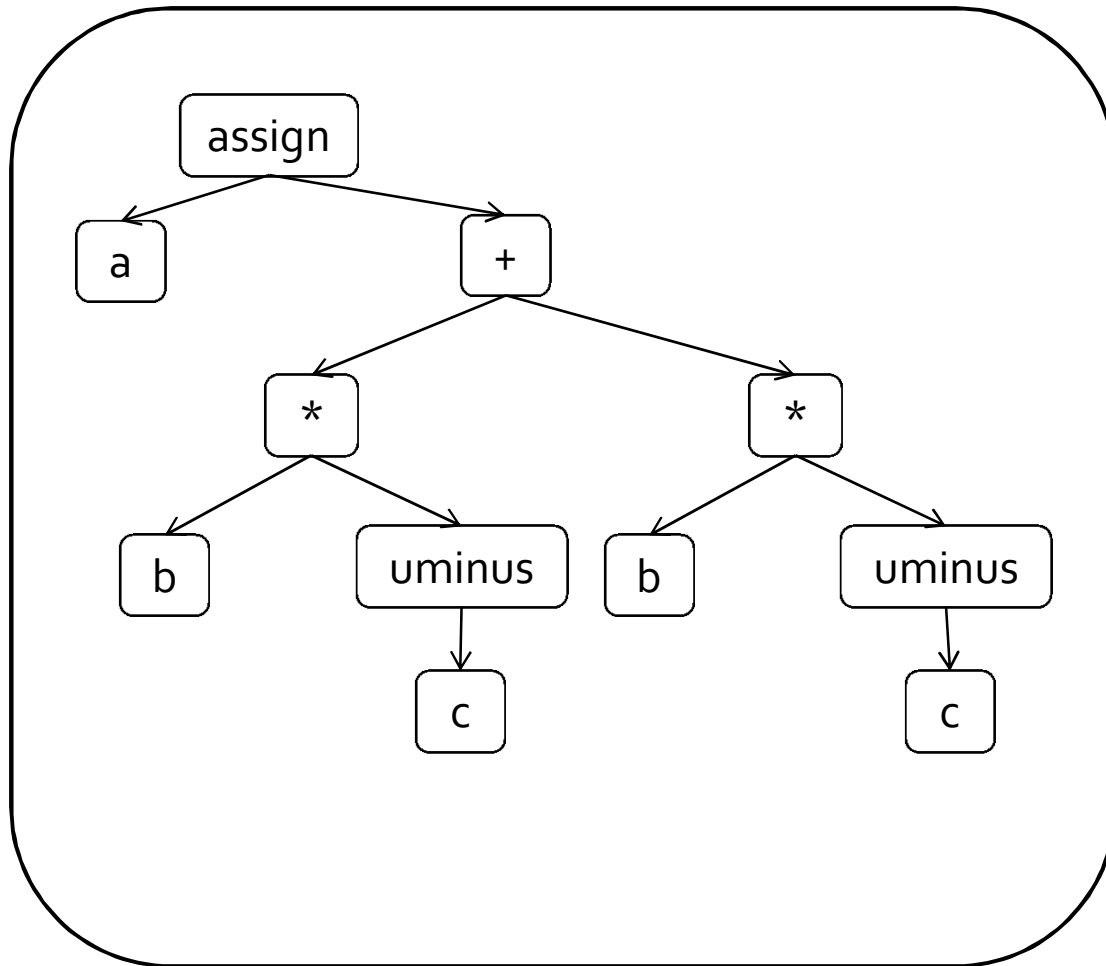


0	id	b	
1	id	c	
2	uminus	1	
3	*	0	2
4	id	b	
5	id	c	
6	uminus	5	
7	*	4	6
8	+	3	7
9	id	a	
10	assign	9	8
11	...		

Three Address Code (3AC)

- Every instruction operates on three addresses
 - $\text{result} = \text{operand}_1 \text{ operator } \text{operand}_2$
- Close to low-level operations in the machine language
 - Operator is a basic operation
- Statements in the source language may be mapped to multiple instructions in three address code

Three address code: example



```
t1 := - c  
t2 := b * t1  
t3 := - c  
t4 := b * t3  
t5 := t2 + t4  
a := t5
```

Three address code: example instructions

instruction	meaning
$x := y \text{ op } z$	assignment with binary operator
$x := \text{op } y$	assignment unary operator
$x := y$	assignment
$x := \&y$	assign address of y
$x := *y$	assignment from deref y
$*x := y$	assignment to deref x

instruction	meaning
goto L	unconditional jump
if x relop y goto L	conditional jump

Array operations

- Are these 3AC operations?

`x := y[i]`

```
t1 := &y      ; t1 = address-of y
t2 := t1 + i  ; t2 = address of y[i]
x   := *t2    ; value stored at y[i]
```

`x[i] := y`

```
t1 := &x      ; t1 = address-of x
t2 := t1 + i  ; t2 = address of x[i]
*t2 := y     ; store through pointer
```

Three address code: example

```
int main(void) {  
    int i;  
    int b[10];  
    for (i = 0; i < 10; ++i)  
        b[i] = i*i;  
}
```

```
i := 0 ; assignment  
L1: if i >= 10 goto L2 ; conditional jump  
    t0 := i*i  
    t1 := &b ; address-of operation  
    t2 := t1 + i ; t2 holds the address of b[i]  
    *t2 := t0 ; store through pointer  
    i := i + 1  
    goto L1  
L2:
```

Three address code

- Choice of instructions and operators affects code generation and optimization
- Small set of instructions
 - Easy to generate machine code
 - Harder to optimize
- Large set of instructions
 - Harder to generate machine code
- Typically prefer small set and smart optimizer

Creating 3AC

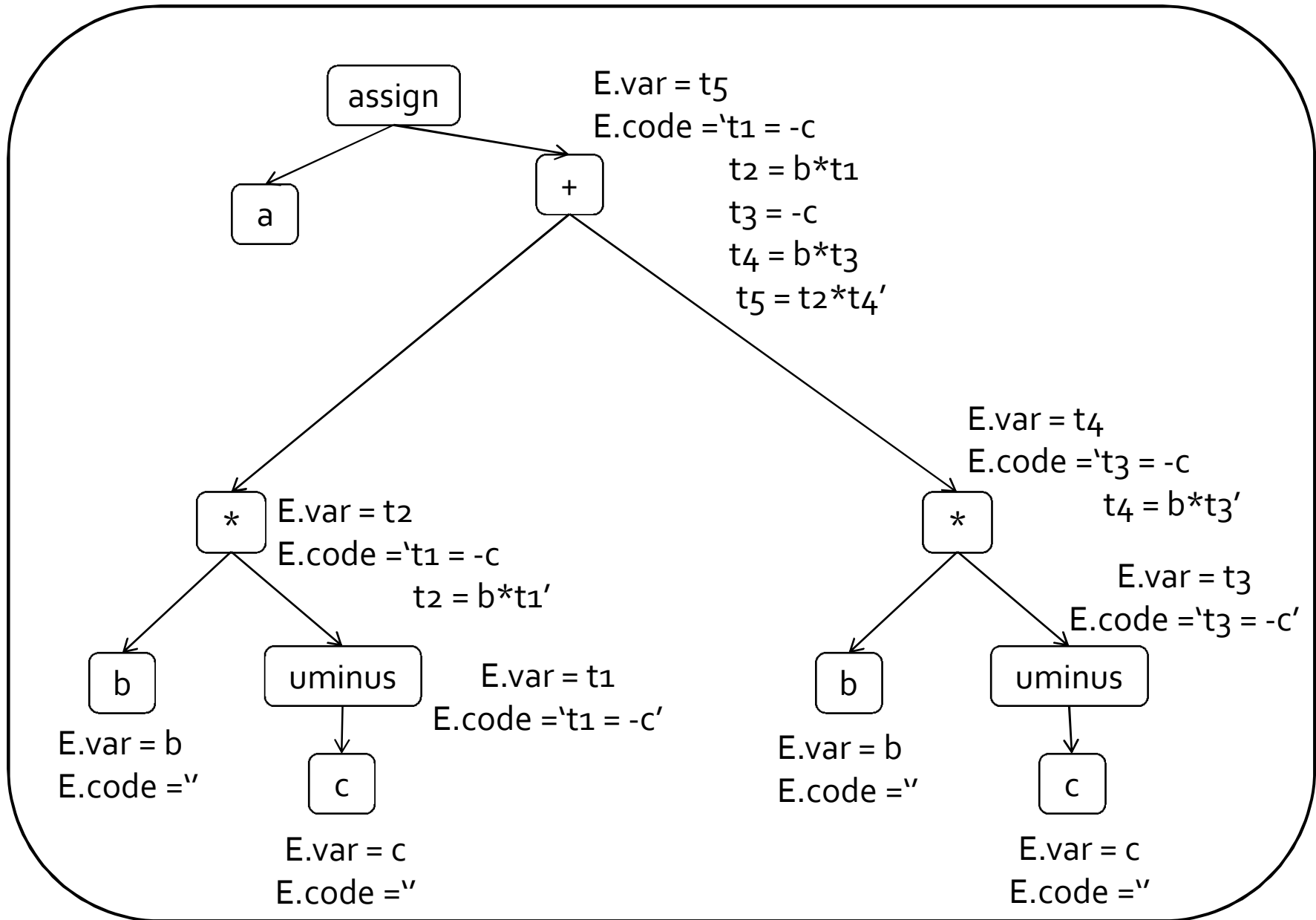
- Assume bottom up parser
 - Why?
- Creating 3AC via syntax directed translation
- Attributes
 - code – code generated for a nonterminal
 - var – name of variable that stores result of nonterminal
- freshVar – helper function that returns the name of a fresh variable

Creating 3AC: expressions

production	semantic rule
$S \rightarrow id := E$	$S.code := E.code \parallel gen(id.var := E.var)$
$E \rightarrow E_1 + E_2$	$E.var := freshVar();$ $E.code = E_1.code \parallel E_2.code \parallel gen(E.var := E_1.var '+' E_2.var)$
$E \rightarrow E_1 * E_2$	$E.var := freshVar();$ $E.code = E_1.code \parallel E_2.code \parallel gen(E.var := E_1.var '*' E_2.var)$
$E \rightarrow - E_1$	$E.var := freshVar();$ $E.code = E_1.code \parallel gen(E.var := 'uminu' E_1.var)$
$E \rightarrow (E_1)$	$E.var := E_1.var$ $E.code = '(' \parallel E_1.code \parallel ')'$
$E \rightarrow id$	$E.var := id.var; E.code = ''$

(we use \parallel to denote concatenation of intermediate code fragments)

example

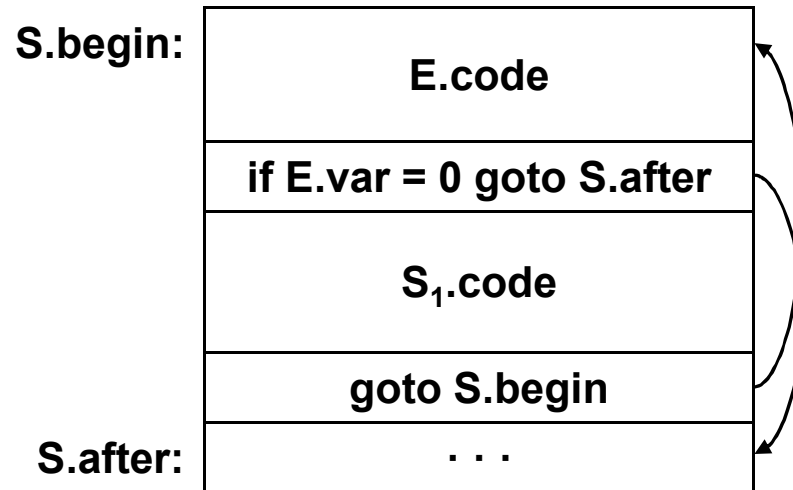


Creating 3AC: control statements

- 3AC only supports conditional/unconditional jumps
- Add labels
- Attributes
 - begin – label marks beginning of code
 - after – label marks end of code
- Helper function `freshLabel()` allocates a new fresh label

Creating 3AC: control statements

$S \rightarrow \text{while } E \text{ do } S_1$



production	semantic rule
$S \rightarrow \text{while } E \text{ do } S_1$	$S.\text{begin} := \text{freshLabel}();$ $S.\text{after} = \text{freshLabel}();$ $S.\text{code} :=$ $\text{gen}(S.\text{begin} \text{ ':'}) \parallel E.\text{code} \parallel$ $\text{gen}(\text{'if' } E.\text{var} \text{ '=' '0' 'goto' } S.\text{after}) \parallel$ $S_1.\text{code} \parallel \text{gen}(\text{'goto' } S.\text{begin}) \parallel \text{gen}(S.\text{after} \text{ ':'})$

Representing 3AC

- Quadruple (op,arg1,arg2,result)
- Result of every instruction is written into a new temporary variable
- Generates many variable names
- Can move code fragments without complicated renaming
- Alternative representations may be more compact

```
t1 = - c
t2 = b * t1
t3 = - c
t4 = b * t3
t5 = t2 * t4
a = t5
```

	op	arg 1	arg 2	result
(0)	uminus	c		t ₁
(1)	*	b	t ₁	t ₂
(2)	uminus	c		t ₃
(3)	*	b	t ₃	t ₄
(4)	+	t ₂	t ₄	t ₅
(5)	:=	t ₅		a

Allocating Memory

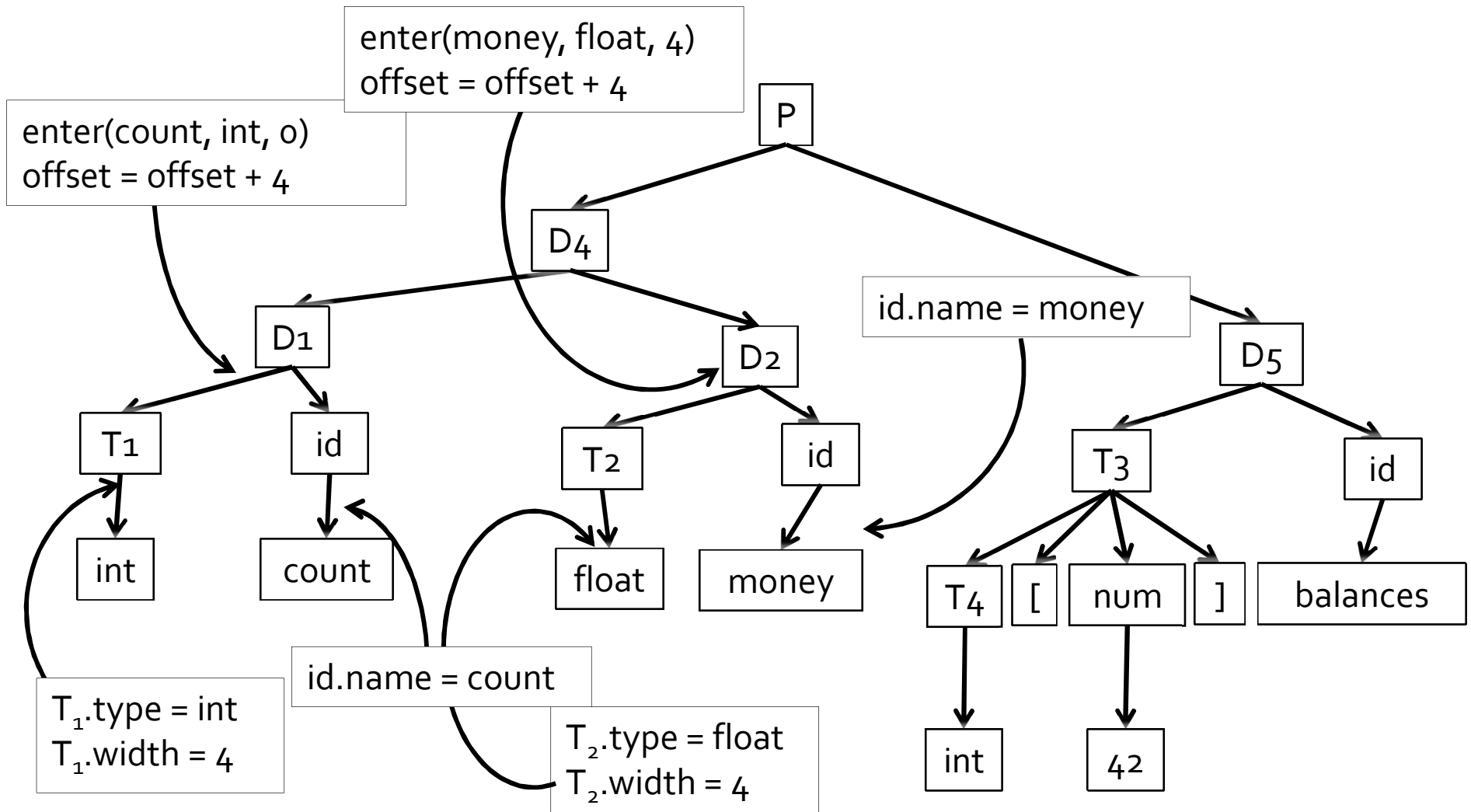
- Type checking helped us guarantee correctness
- Also tells us
 - How much memory allocate on the heap/stack for variables
 - Where to find variables (based on offsets)
 - Compute address of an element inside array (size of stride based on type of element)

Allocating Memory

- Global variable “offset” with memory allocated so far

production	semantic rule
$P \rightarrow D$	{ offset := 0 }
$D \rightarrow D D$	
$D \rightarrow T \text{ id};$	{ enter(id.name, T.type, offset); offset += T.width }
$T \rightarrow \text{integer}$	{ T.type := int; T.width = 4 }
$T \rightarrow \text{float}$	{ T.type := float; T.width = 8 }
$T \rightarrow T_1[\text{num}]$	{ T.type = array (num.val, T ₁ .Type); T.width = num.val * T ₁ .width; }
$T \rightarrow *T_1$	{ T.type := pointer(T ₁ .type); T.width = 4 }

Allocating Memory



Adjusting to bottom-up

production	semantic rule
$P \rightarrow M D$	
$M \rightarrow \varepsilon$	{ offset := o }
$D \rightarrow D D$	
$D \rightarrow T \text{ id};$	{ enter(id.name, T.type, offset); offset += T.width }
$T \rightarrow \text{integer}$	{ T.type := int; T.width = 4 }
$T \rightarrow \text{float}$	{ T.type := float; T.width = 8 }
$T \rightarrow T_1[\text{num}]$	{ T.type = array (num.val, T1.Type); T.width = num.val * T1.width; }
$T \rightarrow *T_1$	{ T.type := pointer(T1.type); T.width = 4 }

Generating IR code

- Option 1
accumulate code in AST attributes
- Option 2
emit IR code to a file during compilation
 - If for every production the code of the left-hand-side is constructed from a concatenation of the code of the RHS in some fixed order

Expressions and assignments

production	semantic action
$S \rightarrow id := E$	{ p:= lookup(id.name); if p \neq null then emit (p $\text{'} := \text{'}$ E.var) else error }
$E \rightarrow E_1 \text{ op } E_2$	{ E.var := freshVar(); emit (E.var $\text{'} := \text{'}$ E1.var op E2.var) }
$E \rightarrow - E_1$	{ E.var := freshVar(); emit (E.var $\text{'} := \text{'}$ 'uminus' E1.var) }
$E \rightarrow (E_1)$	{ E.var := E1.var }
$E \rightarrow id$	{ p:= lookup(id.name); if p \neq null then E.var :=p else error }

Boolean Expressions

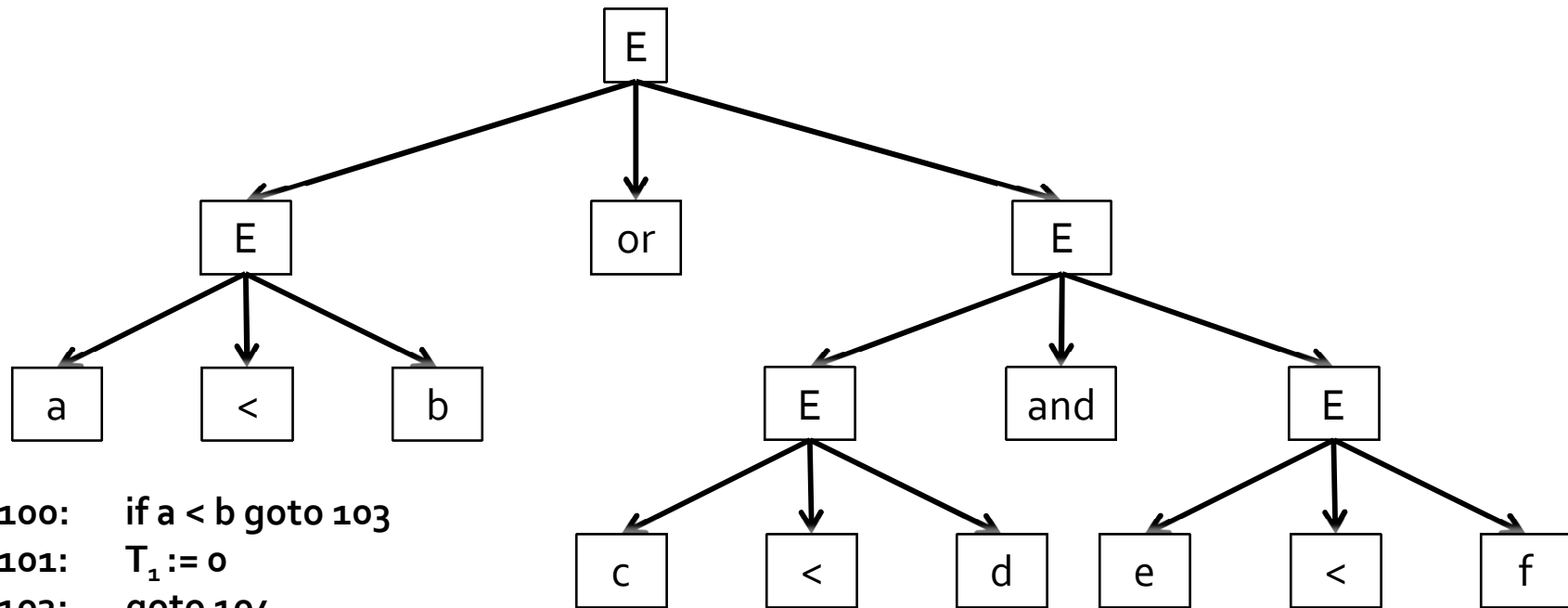
production	semantic action
$E \rightarrow E_1 \text{ op } E_2$	{ E.var := freshVar(); emit(E.var := ' E1.var op E2.var) }
$E \rightarrow \text{not } E_1$	{ E.var := freshVar(); emit(E.var := 'not' E1.var) }
$E \rightarrow (E_1)$	{ E.var := E1.var }
$E \rightarrow \text{true}$	{ E.var := freshVar(); emit(E.var := '1') }
$E \rightarrow \text{false}$	{ E.var := freshVar(); emit(E.var := '0') }

- Represent true as 1, false as 0
- Wasteful representation, creating variables for true/false

Boolean expressions via jumps

production	semantic action
$E \rightarrow id_1 \text{ op } id_2$	<pre>{ E.var := freshVar(); emit('if' id1.var relop id2.var 'goto' nextStmt+2); emit(E.var := '0'); emit('goto ' nextStmt + 1); emit(E.var := '1') }</pre>

Example



100: if a < b goto 103
101: $T_1 := 0$
102: goto 104
103: $T_1 := 1$

104: if c < d goto 107
105: $T_2 := 0$
106: goto 108
107: $T_2 := 1$

108: if e < f goto 111
109: $T_3 := 0$
110: goto 112
111: $T_3 := 1$
112: $T_4 := T_2 \text{ and } T_3$
113: $T_5 := T_1 \text{ or } T_4$

Short circuit evaluation

- Second argument of a boolean operator is only evaluated if the first argument does not already determine the outcome
- $(x \text{ and } y)$ is equivalent to `if x then y else false;`
- $(x \text{ or } y)$ is equivalent to `if x then true else y`

example

$a < b$ or $(c < d$ and $e < f)$

```
100: if a < b goto 103
101: T1 := 0
102: goto 104
103: T1 := 1
104: if c < d goto 107
105: T2 := 0
106: goto 108
107: T2 := 1
108: if e < f goto 111
109: T3 := 0
110: goto 112
111: T3 := 1
112: T4 := T2 and T3
113: T5 := T1 and T4
```

naive

```
100: if a < b goto 105
101: if !(c < d) goto 103
102: if e < f goto 105
103: T := 0
104: goto 106
105: T := 1
106:
```

Short circuit evaluation

More examples

```
int denom = 0;
if (denom && nom/denom) {
    oops_i_just_divided_by_zero();
}
```

```
int x=0;
if (++x>0 && x==1) {
    hmmm();
}
```

Summary

- Three address code (3AC)
- Generating 3AC
- Boolean expressions
- Short circuit evaluation

Next time

- Generating IR for control structures
 - While, for, if
- backpatching

The End