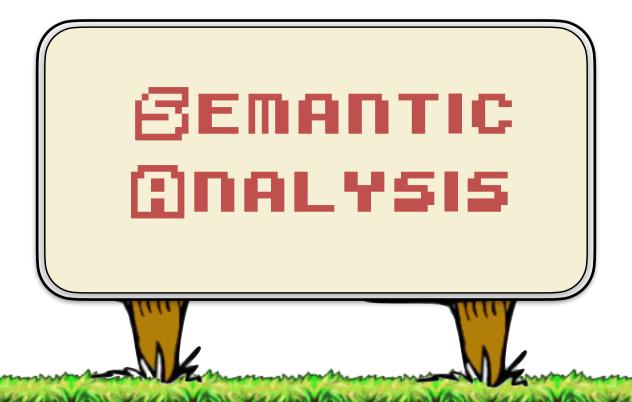
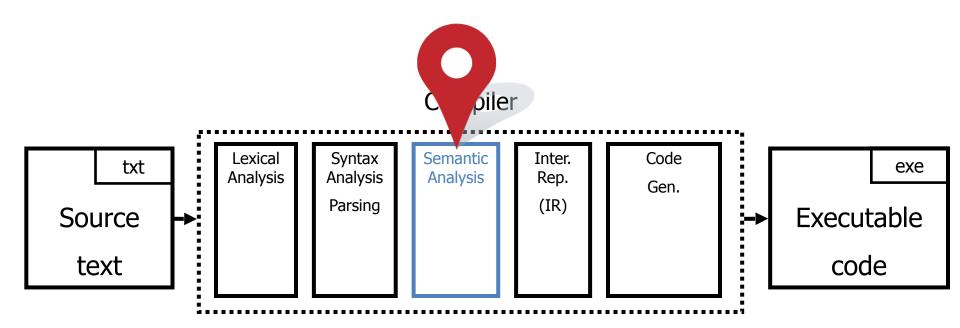
THEORY OF GOMPILATION

LECTURE 04

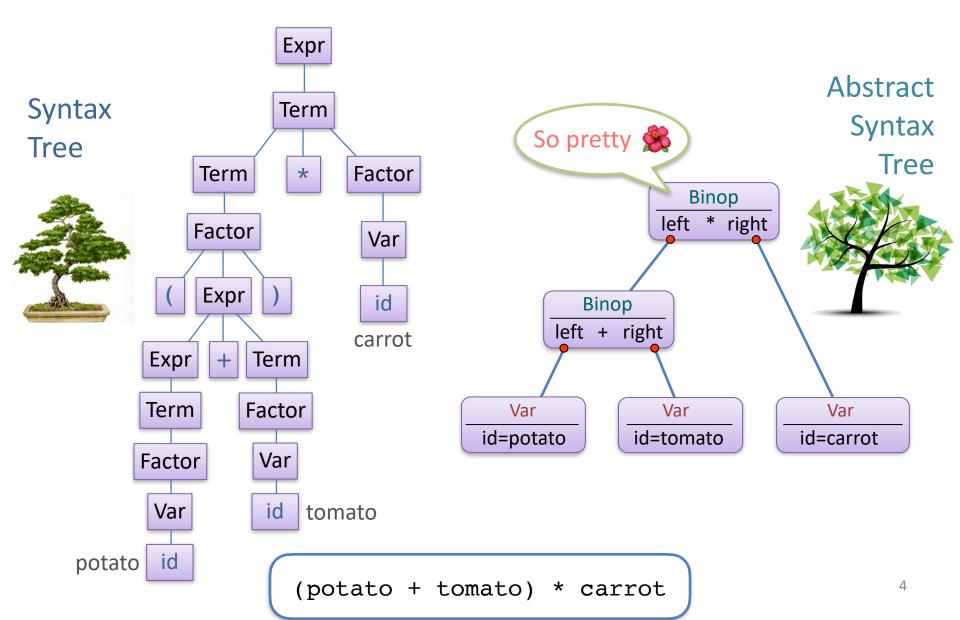


You are here

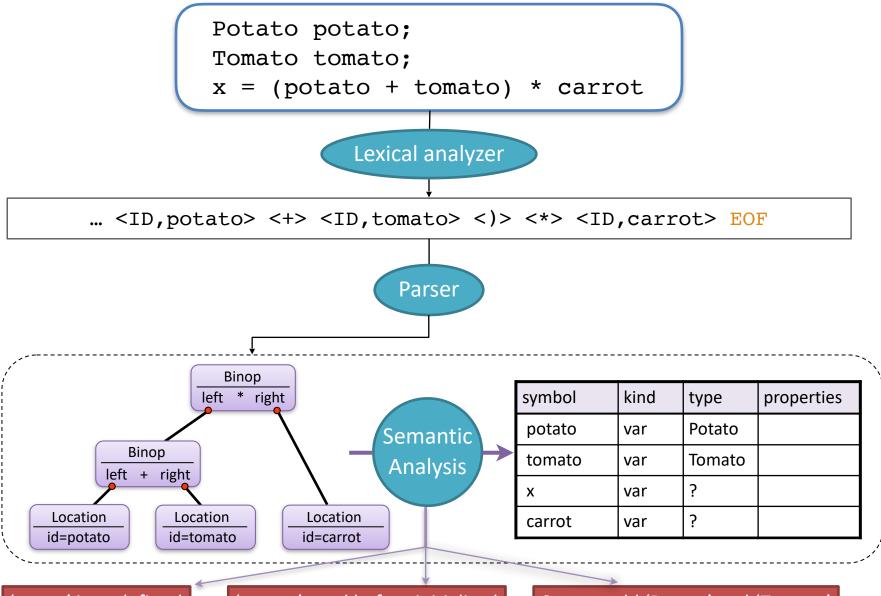


You are here txt Lexical Process Semantic Syntax tokens characters AST Analysis text input Analysis **Analysis** Source text exe Executable code

Building the Abstract Syntax Tree



What We Want



Semantic Analysis

- Often called "Contextual analysis"
 - ▶ As opposed to our syntax analysis which was "context free"
- Properties that cannot be formula
 - ▶ Declare before use
 - ▶ Type checking
 - ▶ Initialization
 - **)** ...
- Properties that are clumsy to fo
 - "break" only appears inside a loo
 - **)** ...

```
stmts \rightarrow \epsilon | stmt; stmts
stmt \rightarrow if expr \{ stmts \} s else
          while expr { stmts \ell }
          let id : type = expr
          | id = expr
s else \rightarrow \varepsilon | else { stmts }
stmts^{\ell} \rightarrow \epsilon \mid stmt^{\ell} ; stmts^{\ell}
stmt^{\ell} \rightarrow if expr \{ stmts^{\ell} \} s else^{\ell}
          | while expr { stmts<sup>\ell</sup> }
          let id: type = expr
          | id = expr | break
s_{else} \rightarrow \epsilon \mid else \{ stmts^{\ell} \}
```

oter 11

Semantic Analysis

Identification

- Gather information about each named item in the program
- ▶ e.g., what is the declaration for each usage

Context checking

- Type checking
- ▶ e.g., the condition in an if-statement is a Boolean

Identification

```
month : integer RANGE 1..12;
month := 1;
while (month <= 12) {
   print(month_name[month]);
   month := month + 1;
}</pre>
```

Symbol table

```
month : integer RANGE 1..12;
month_name : string RANGE 1..12;
month := 1;
while (month <= 12) {
    print(month_name[month]);
    month := month + 1;</pre>
month := nonth + 1;
```

- A table containing information about identifiers in the program
- Single entry for each named item

Not so fast...

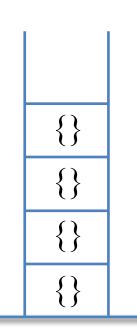
```
A struct field named "i"
struct one int {
   int i;
                                   A struct variable named "i"
} i;
                            Assignment to the "i" field of struct "i"
main() {
 i.i = 42;
 int t = i.i;
                                  Reading the "i" field of struct "i"
 printf("%d",t);
```

Not so fast...

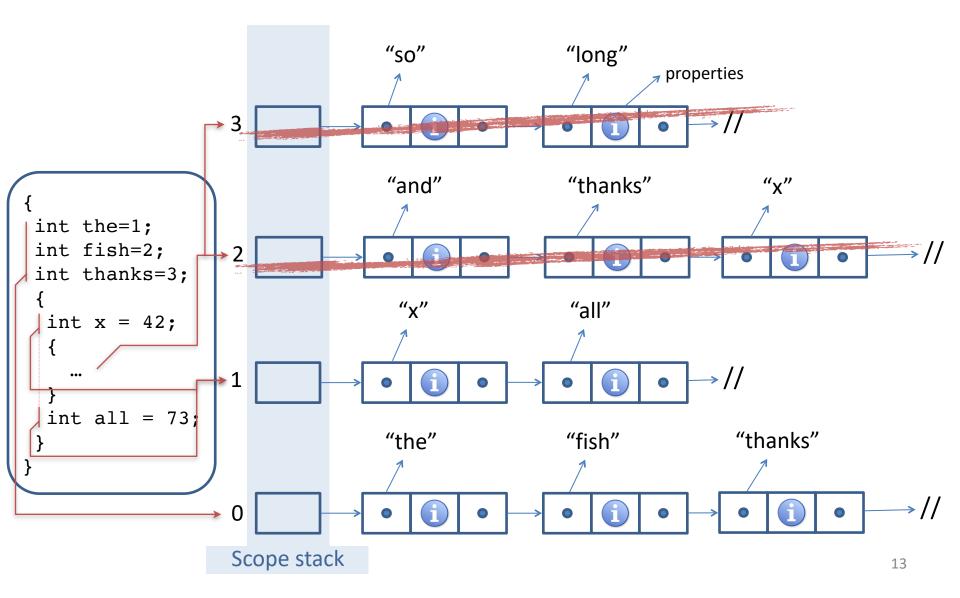
```
A struct field named "i"
struct one int {
   int i;
                                  A struct variable named "i"
} i;
main() {
                            Assignment to the "i" field of struct "i"
 i.i = 42;
 int t = i.i;
                                  Reading the "i" field of struct "i"
 printf("%d",t);
                                          another variable
   int i = 73;
                                            named "i"?!
  printf("%d",i);
                                    Now what?!
                                                           11
```

Scopes

- Typically: stack structured scopes
- Scope entry
 - push new empty scope element
- Scope exit
 - pop scope element and discard its content
- Identifier declaration
 - identifier created inside (current) top scope
- Identifier Lookup
 - Search for identifier top-down in scope stack



Scope-structured Symbol Table

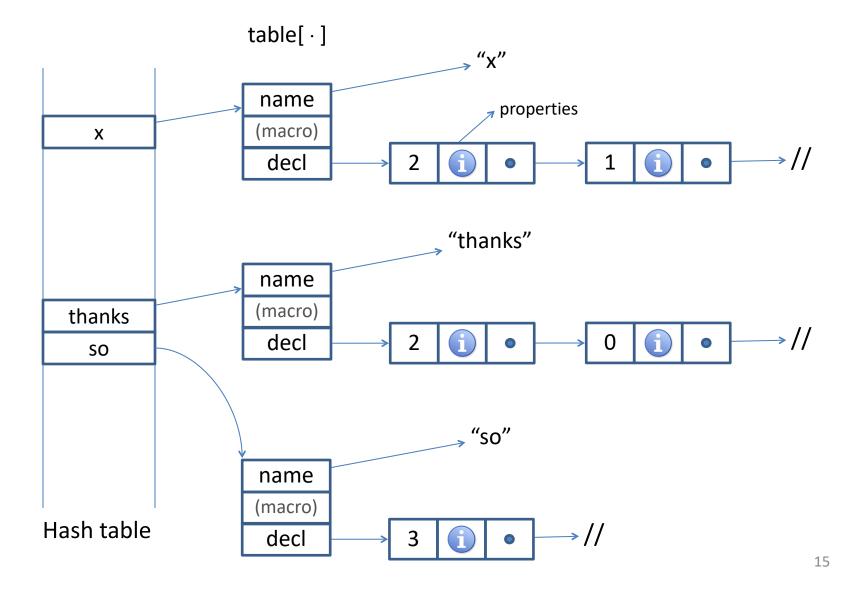


Scope and Symbol Table

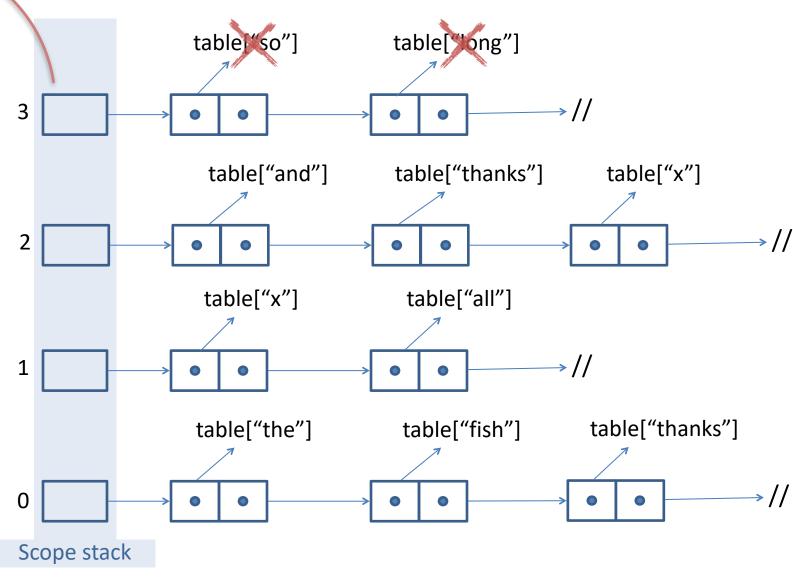
- Scope × Identifier → properties
 - Expensive lookup: when a reference to x is encountered, the entire table needs to be searched

- A better solution
 - Hash table over identifiers
 - For each key: list of scopes where it has been defined (innermost first)

Hash Table-based Symbol Table



Scope info



Remember Lexing+Parsing?

- How did we know to always map an identifier to the same token?
 - ▶ We didn't! Each occurrence is a new token (but with the same lexical attribute, "name").

Now is the time to match them up.

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

- Semantic attributes
 - Attributes attached to grammar symbols
- Semantic actions
 - Define how to update the attributes

Attribute grammars

Attribute Grammars

- Semantic attributes
 - Every grammar symbol has attached attributes
 - Example: N.name
- Semantic actions
 - Every production rule can define how to assign
 values to attributes

 Create semantic attribute

```
T → int { T.name = "int"; } Use lexical attribute

T → struct id { T.name = id.value; }

Use semantic attribute

D → T id ;

D.type = T.name; of child node

D.varname = id.value; }

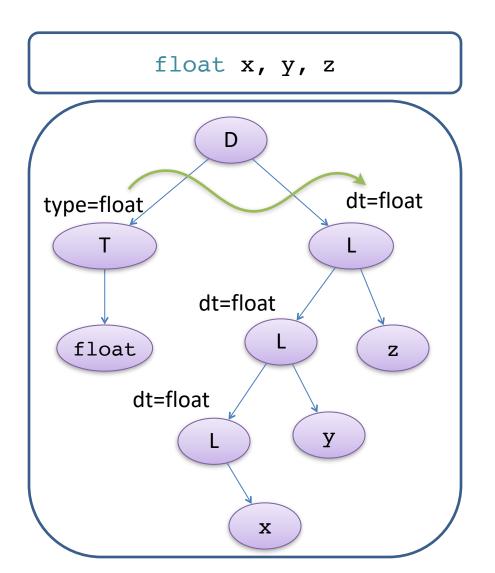
D.varname = id.value; }
```

Indexed symbols

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

```
\begin{array}{c} \text{Decl} \rightarrow \text{Decl} \text{, id} \\ \text{Decl} \rightarrow \text{Decl}_1 \text{, id} \\ \text{Decl}_1 \text{.dtype} = \text{Decl.dtype;} \end{array} \}
```

Example — Analyzing Declarations



Production	Semantic Rule	
$D \rightarrow T L$	L.dt = T.type	
$T \rightarrow \text{int}$	T.type = int	
$T \rightarrow float$	T.type = float	
$L \rightarrow L_1$, id	L ₁ .dt = L.dt addType(id.entry, L.dt)	
$L \rightarrow id$	addType(id.entry, L.dt)	

Attribute Evaluation

- Build the AST
- Fill attributes of terminals with their lexical values
- Execute semantic actions 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 it's available
 - Each attribute will get a value only once

Attribute Dependencies

All semantic actions take the form

$$a_1 = f_1(b_{1-1}, b_{1-2}, ...)$$

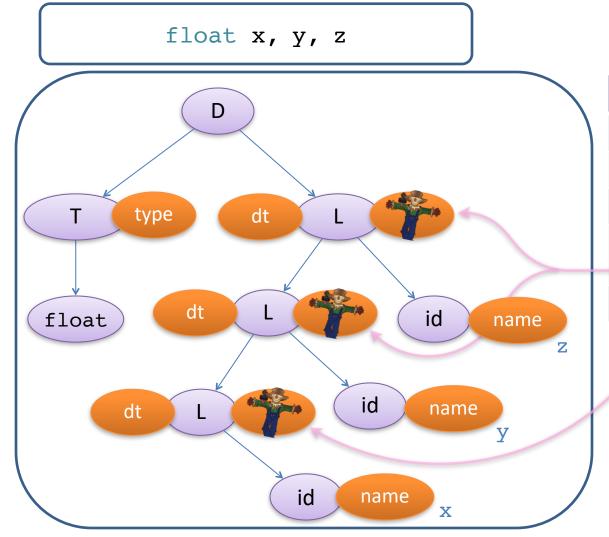
 $a_2 = f_2(b_{2-1}, b_{2-2}, ...)$

▶ For actions with side effects, e.g. print(b.name): we introduce a dummy attribute as their result.



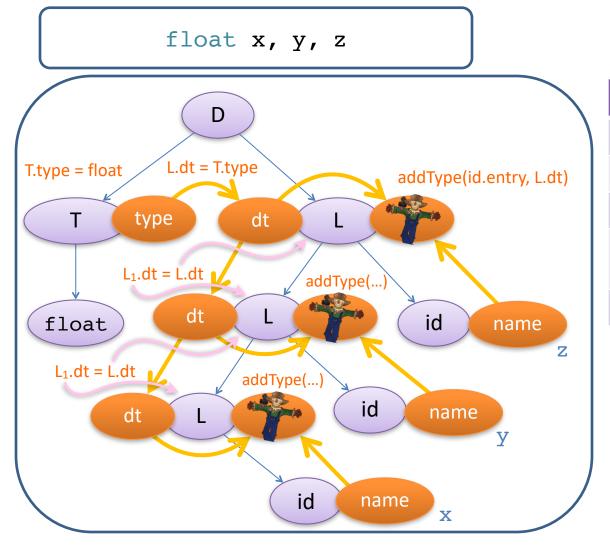
- Build a directed dependency graph G
 - ▶ For every attribute a of a node u in the AST, create a node u.a
 - ▶ For every dependency of the form $u.a_i = \cdots v.b_j \cdots$ create an edge $v.b_i \rightarrow u.a_i$

Example



Prod.	Semantic Rule
$D\toTL$	L.dt = T.type
$T \to int$	T.type = int
$T \to \mathtt{float}$	T.type = float
$L \rightarrow L_1$, id	L ₁ .dt = L.dt addType(id.name, L.dt)
$L \rightarrow id$	addType(id.name, L.dt)

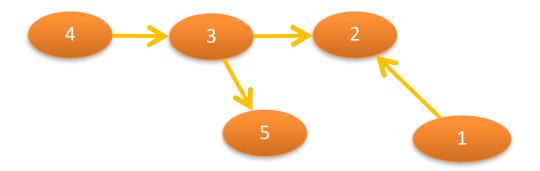
Example



Prod.	Semantic Rule
$D\toTL$	L.dt = T.type
$T \to int$	T.type = int
$T \to \mathtt{float}$	T.type = float
$L \rightarrow L_1$, id	L ₁ .dt = L.dt addType(id.name, L.dt)
$L \rightarrow id$	addType(id.name, L.dt)

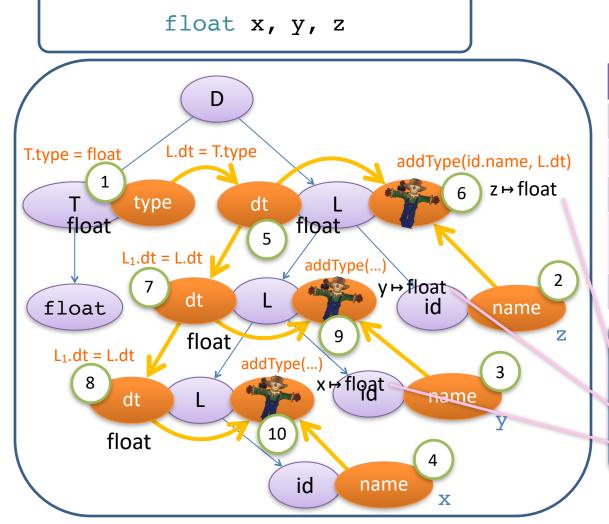
Topological Order

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



Example topological orderings: $\langle 1 4 3 2 5 \rangle$, $\langle 4 3 5 1 2 \rangle$

Example



Prod.	Semantic Rule
$D\toTL$	L.dt = T.type
$T \to int$	T.type = int
$T \to \mathtt{float}$	T.type = float
$L \rightarrow L_1$, id	L ₁ .dt = L.dt addType(id.name, L.dt)
$L \rightarrow id$	addType(id.name, L.dt)

symbol	kind	type	properties
⇒ z	var	float	
⇒ y	var	float	
→ X	var	float	

Symbol Table

But what about cycles?

- For a given attribute grammar hard to detect if it has cyclic dependencies
 - Exponential cost

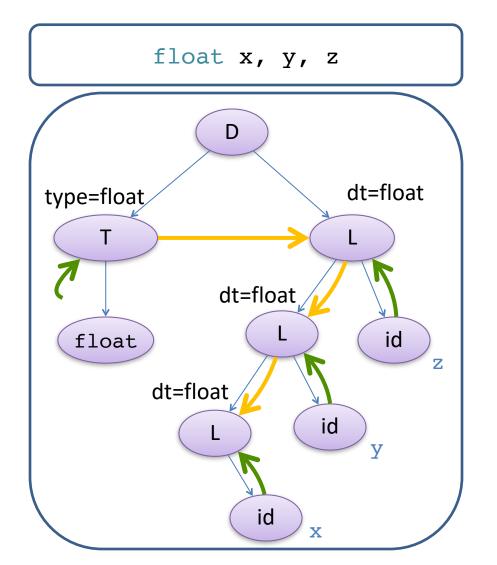
- Special classes of attribute grammars
 - Our "usual trick":
 sacrifice generality for predictable performance

Synthesized vs. Inherited Attributes

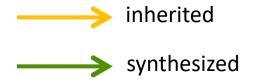
- Synthesized attributes
 - Attributes whose values at a given node depend only on the attributes of its children (and itself)
- Inherited attributes
 - Attributes whose values at a given node depend
 only on the attributes of its parent and siblings

Attributes that don't depend on anything (e.g. those of a token) — by convention, are classified as *synthesized* attributes.

Synthesized vs. Inherited Attributes



Prod.	Semantic Rule
$D \to T \; L$	L.dt = T.type
$T \to int$	T.type = integer
$T \rightarrow float$	T.type = float
$L \to L_1, id$	L ₁ .dt = L.dt addType(id.name, L.dt)
$L \rightarrow id$	addType(id.name, L.dt)



S-attributed Grammars

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

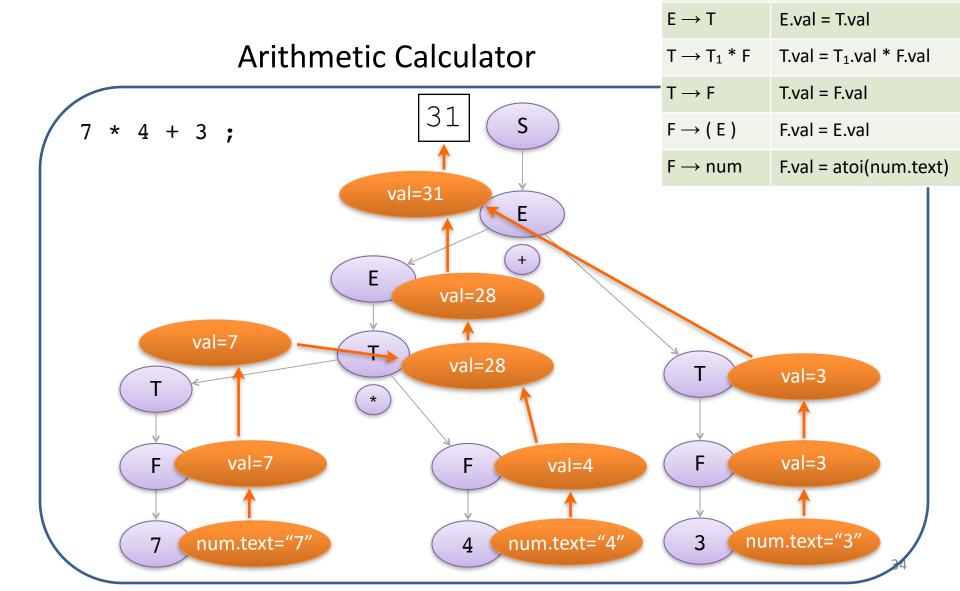
- Can be computed by any bottom-up parser during parsing — no need to construct dependency graph
 - Attributes can be stored on the parsing stack
 - Reduce operation computes the (synthesized) attribute from attributes of children

S-attributed Grammars

Arithmetic Calculator

Production	Semantic Rule
$S \rightarrow E$;	print(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
$F \rightarrow num$	F.val = atoi(num.text)

S-attributed Grammars



Production

 $E \rightarrow E_1 + T$

 $S \rightarrow E$;

Semantic Rule

 $E.val = E_1.val + T.val$

print(E.val)

L-attributed Grammars

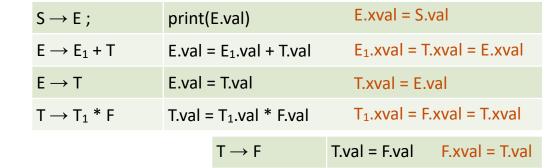
- L-attributed attribute grammar: when every attribute in a production $A \rightarrow X_1...X_n$ is either
 - ▶ A synthesized attribute, or
 - ▶ An inherited attribute of X_j , $1 \le j \le n$ that only depends on
 - Attributes (synthesized or inherited) of $X_1...X_{j-1}$ to the **left** of X_j
 - Inherited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

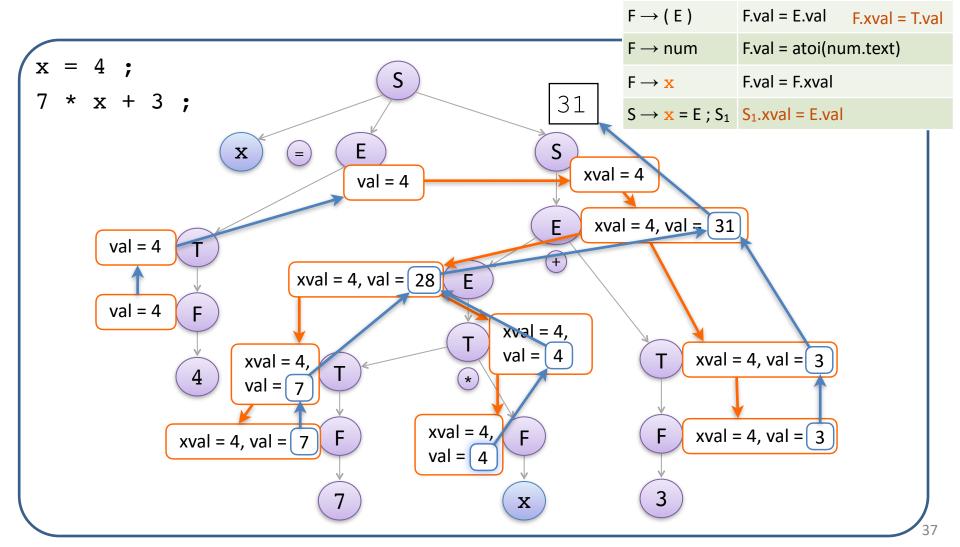
L-attributed Grammars

Arithmetic Calculator with Variables

$$e.g.,$$
 $x = 4 ; 7 * x + 3 ;$

Production	Semantic Rule	
$S \rightarrow E$;	print(E.val)	E.xval = S.xval
$E \rightarrow E_1 + T$	$E.val = E_1.val + T.val$	$E_1.xval = T.xval$
$E \rightarrow T$	E.val = T.val	E.xval = T.xval
$T \to T_1 * F$	$T.val = T_1.val * F.val$	T ₁ .xval = F.xval
$T \rightarrow F$	T.val = F.val	F.xval = T.xval
$F \rightarrow (E)$	F.val = E.val	E.xval = F.xval
$F \rightarrow num$	F.val = atoi(num.text)	
$F \rightarrow x$	F.val = F.xval	
$S \rightarrow \mathbf{x} = E ; S_1$	$S_1.xval = E.val$	





- In recursive-descent parsers:
 - Pass inherited attributes down as arguments

```
D() {
   T();
   L();
}
L() {
   match(ID);
   A();
}
A() {
   if (current == COMMA) {
      match(COMMA);
      L();
   } else if (current == EOF) {
   } else error();
```

```
Prod.Semantic RuleD \rightarrow TLL.dt = T.typeT \rightarrow intT.type = integerT \rightarrow floatT.type = floatL \rightarrow id AA.dt = L.dt<br/>addType(id.entry, L.dt)A \rightarrow LL.dt = A.dtA \rightarrow E
```

```
T() {
    if (current == INT) {
        match(INT);
    } else if (current == FLOAT) {
        match(float);
    } else error();
}
```

- In recursive-descent parsers:
 - Pass inherited attributes down as arguments

```
D() {
   Attrs t = T();
   L({dt rtipe});
} Pass synthesized
   attributes up as
   return values

L(Attrs ih) {
   Token id = match(ID);
   addType(id.entry, ih.dt);
   A({dt rib.dt});
```

if (current == COMMA) {

} else if (current == EOF) {

match(COMMA);

} else error();

 $L(\{dt \mapsto ih.dt\});$

A(Attrs ih) {

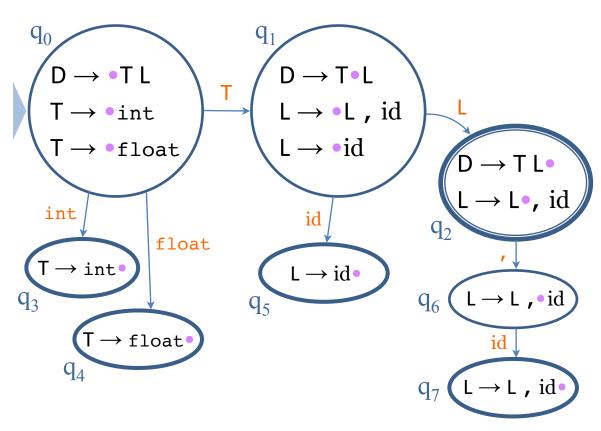
```
Prod.Semantic RuleD \rightarrow T LL.dt = T.typeT \rightarrow intT.type = integerT \rightarrow floatT.type = floatL \rightarrow id AA.dt = L.dt<br/>addType(id.entry, L.dt)A \rightarrow LL.dt = A.dtA \rightarrow LL.dt = A.dt
```

```
T() {
    if (current == INT) {
        match(INT);
        return {type \( \to \) int};
    } else if (current == FLOAT) {
        match(float);
        return {type \( \to \) float};
    } else error();
}
```

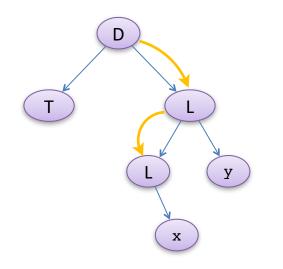
• In shift-reduce parsers:

Semantic actions are performed during reduce.

We have a problem — tree is built bottom-up

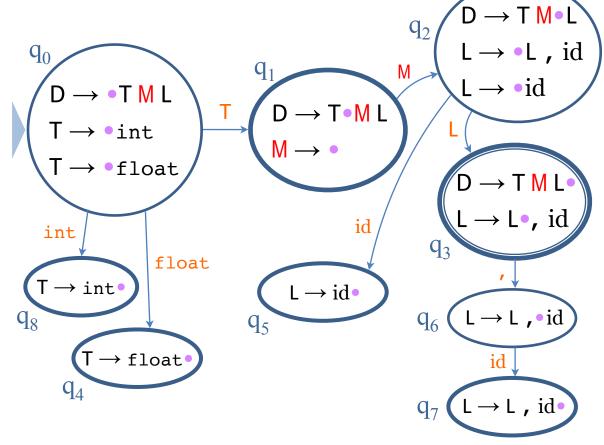


Prod.	Semantic Rule
$D\toTL$	L.dt = T.type
$T \to int$	T.type = integer
$T \rightarrow float$	T.type = float
$L \rightarrow L_1$, id	L ₁ .dt = L.dt addType(id.entry, L.dt)
$L \rightarrow id$	addType(id.entry, L.dt)

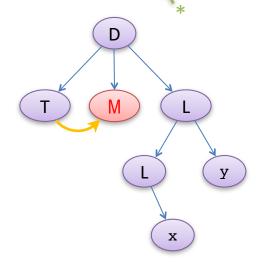


• In shift-reduce parsers:





Prod.	Semantic Rule	
$D\toTML$	dtype = null	
$T \to int$	T.type = integer	
$T \rightarrow float$	T.type = float	
$L \to L_1 \text{, id}$	addType(id.entry, dtype)	
$L \rightarrow id$	addType(id.entry, dtype)	
$M \to \epsilon$	dtype = T.type	



Marker Variables

 Since a marker only appears in one production rule, it is commonly abbreviated:

$D \rightarrow TMI$	→ T M L action ₁		
D / I IVI L	action	$D \rightarrow T \{action_2\} L$	action ₁
$M \rightarrow \epsilon$	action ₂	D / I (action2) L	action
$ V \rightarrow \varepsilon$ action2	actions		

action₂ is called a mid-rule action.

It is important to remember that adding mid-rule actions inherently changes the grammar. The marker variable is there even if it is not explicitly visible.

(A) STABILOBOS

STABILO BOS

Marker Variables

(№ STABILO BOS

In particular, marker variables and the associated ε-productions can violate your grammar's LR(0)/SLR/LALR/LR(1)-ness



Prod.	Semantic Rule
$S \rightarrow D$	
$D\toTL$	
$D \rightarrow TL[]$	
$T \to int$	T.type = integer
$T \rightarrow float$	T.type = float
$L \rightarrow id$	addType(id.entry, dtype)

This grammar is LR(0)

Prod.	Semantic Rule
$S \rightarrow D$	
$D\toT{}^{\mathop{M}}L$	
$D \rightarrow T N L []$	
$T \rightarrow int$	T.type = integer
$T \rightarrow float$	T.type = float
$L \rightarrow id$	addType(id.entry, dtype)
$M \rightarrow \epsilon$	dtype = T.type
$N \rightarrow \epsilon$	<pre>dtype = array(T.type)</pre>

This grammar is not even LR(1)

Semantic Checks

- Scope rules
 - Use symbol table to check:
 - No multiple definition of same identifier
 - Identifiers defined before used

- Type checking
 - Check that types in the program are consistent
 - o How?

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

- Typing rules specify
 - which types can be combined with certain operator
 - Assignment of expression to variable
 - Formal and actual parameters of a method call

Examples

```
string string
"drive" + "drink"
    string
    int string
    42 + "the answer"
    ERROR
```

Type System

- A type system of a programming language is a way to define how "good" programs behave
 - Good programs = well-typed programs
 - Bad programs = programs containing type errors

Static Typing most checking done at compile time

Dynamic Typing most checking done at runtime

Static Typing vs. Dynamic Typing

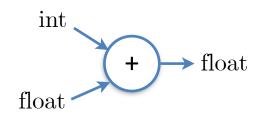
- Static type checking is conservative
 - Any program that is determined to be well-typed is free from certain kinds of errors
 - May reject programs that cannot be statically determined to be safe
 if (some-function(x))
 - ▶ Why?
- Dynamic type checking
 - May accept more programs as valid (runtime info)
 - Errors not caught at compile time
 - Runtime overhead

y = "the answer"

print(y * y);

Type Checking Rules

- Specify for each operator
 - Types of operands
 - Type of result



- Basic Types
 - Building blocks for the type system (type rules)
 - e.g., int, boolean, (sometimes) string
- Type Expressions
 - Array types
 - Function types
 - Record types / Classes

Typing Rules

If E_1 has type int and E_2 has type int, then $E_1 + E_2$ has type int

 E_1 : int E_2 : int

 $E_1 + E_2 : int$

Similarly,

 E_1 : int E_2 : int

 $E_1 * E_2 : long$

If E_1 has type int and E_2 has type float, then $E_1 + E_2$ has type float

 E_1 : int E_2 : float

 $E_1 + E_2$: float

 E_1 : float E_2 : int

 $E_1 + E_2$: float

...and the other way around

Basic Type System

true: boolean

false: boolean

Constants

int-literal: int

string-literal: string

 E_1 : int

 E_2 : int

 $E_1 \diamond E_2$: int

Binary Arithmetic Operators

 E_1 : int

 E_2 : int

 $E_1 \oplus E_2$: boolean

Comparison Operators

 $E_1:T$

 $E_2:T$

 $E_1 \oplus E_2$: boolean

 $riangleright \in \{ ==, != \}$

Basic Type System

E₁: boolean

E₂: boolean

E₁ * E₂ : boolean

⊛ ∈ { &&, || }

Binary Logical Operators

 E_1 : int

 $-E_1$: int

E₁: boolean

! E₁: boolean

Unary **Operators**

Array Operations

 $\mathsf{E}_1:\mathsf{T}[]$

 E_1 .length: int

 $E_1:T[]$

 E_2 : int

 $E_1[E_2]:T$

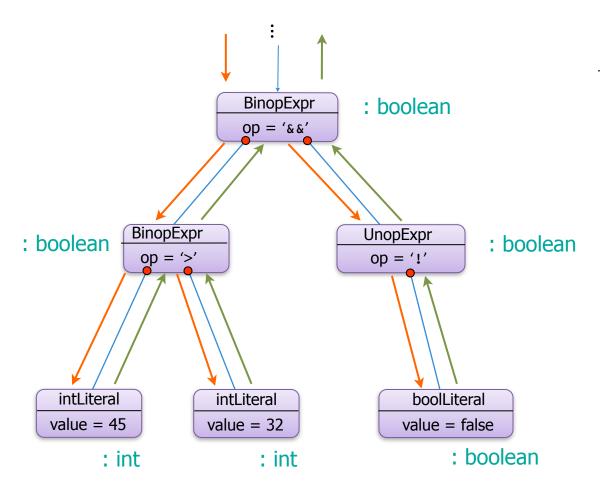
 E_1 : int

new $T[E_1]$: T[]

- Traverse AST and assign types for AST nodes
 - Use typing rules to compute node types

- Alternative: type-check during parsing
 - (Slightly) more complicated
 - But naturally also more efficient

- Use semantic actions
 - ▶ In derivations corresponding to expressions, check that typing rules are satisfied



 E_1 : boolean E_2 : boolean

E₁ ⊛ E₂ : boolean

for $\circledast \in \{ \&\&, | | \}$

E₁: boolean

! E₁ : boolean

 E_1 : int E_2 : int

 $E_1 \oplus E_2$: boolean

false: boolean

int-literal: int

Strongly Typed vs. Weakly Typed

Output: 73

warning: initialization makes integer from pointer without a cast

- Coercion
- Strongly typed
 - C, C++, Java
- Weakly typed
 - JavaScript, Perl, PHP

(Not everybody agrees on this classification)

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

Perl

```
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

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

Type Declarations

- So far, we had a fixed set of types
- But types can also be user-defined

Type Declarations

series: array (1..42) of Real;



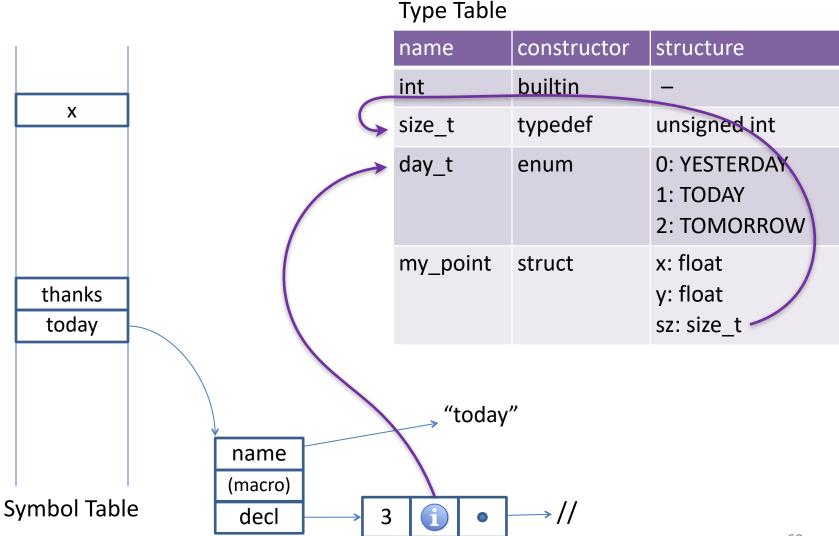
```
type type01@line_73 is array (1..42) of Real;
series : type01@line_73;
```

Type Table

 All types in a compilation unit are collected in a type table

- For each type, its table entry contains
 - ▶ Type constructor: basic, record, array, pointer,...
 - Size and alignment requirements
 - to be used later in code generation
 - Types of components (if applicable)
 - e.g., type of array element, types of record fields

Type Table + Symbol Table



Forward References

```
typedef struct List_Entry *Ptr_List_Entry;
::
struct List_Entry {
    int element;
    Ptr_List_Entry next;
};
```

- Forward references must be resolved at a later point
- A forward reference is added to the type table (as "unresolved"), and later updated when the type declaration is met
- At the end of scope, check that all forward refs have been resolved
- Must add check for circular dependencies between types

Nominal Type System

Type Equivalence = Name Equality

```
Type t1 = ARRAY[Integer] OF Integer;

Type t2 = ARRAY[Integer] OF Integer;
```

t1 not (nominally) equivalent to t2

```
Type t3 = ARRAY[Integer] OF Integer;
Type t4 = t3
```

t3 equivalent to t4

Structural Type System

Type Equivalence = Structure Isomorphism

```
Type t5 = RECORD c: Integer; p: POINTER TO t5; END RECORD;
Type t6 = RECORD c: Integer; p: POINTER TO t6; END RECORD;
Type t7 =
 RECORD
   c: Integer;
   p: POINTER TO
    RECORD
      c: Integer;
      p: POINTER TO t5;
     END RECORD:
END RECORD;
```

t5, t6, t7 are all (structurally) equivalent

In Practice

- Almost all modern languages use a nominal type system
 - ▶ Why?
- One exception to this is type aliases (C's typedef)

▶ Some languages (e.g. Ada) have two kinds for this purpose:

```
type listptr is pointer to Node pointer to Node ≡ listptr

type listptr is new pointer to Node pointer to Node ≡ listptr
```

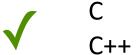
Coercions

 If we expect a value of type T₁ at some point in the program, and find a value of type T₂, is that acceptable?

float
$$x = 3.141$$
;
int $y = x$;



Java C# Scala Rust Ada



Coercions

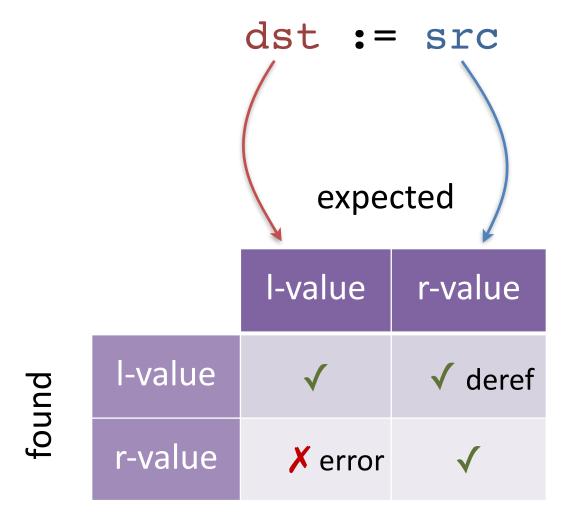
 If we expect a value of type T₁ at some point in the program, and find a value of type T₂, is that acceptable?

I-values and r-values

- A "location" that can be assigned to is called an l-value — it evaluates to an address
- A "datum" that can be read and copied is called an r-value — it evaluates to the actual value

I-values and r-values (example)

I-values and r-values



I-values and r-values (example)

$$x := A[1]$$
 $x := A[A[1]]$
 $A[A[1]] := x + 1$
 $x + 1 := A[1]$
 $A[A[1]] := x + not ok$

Reference and const modifiers

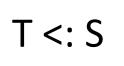
 These fall neatly into the I-value/r-value categories:

* What about const int& y?

- A basic concept in Object-oriented Programming
 - Every class has (can have) a superclass

};

- Subtyping relation:
 - "T is a sub-type of S"

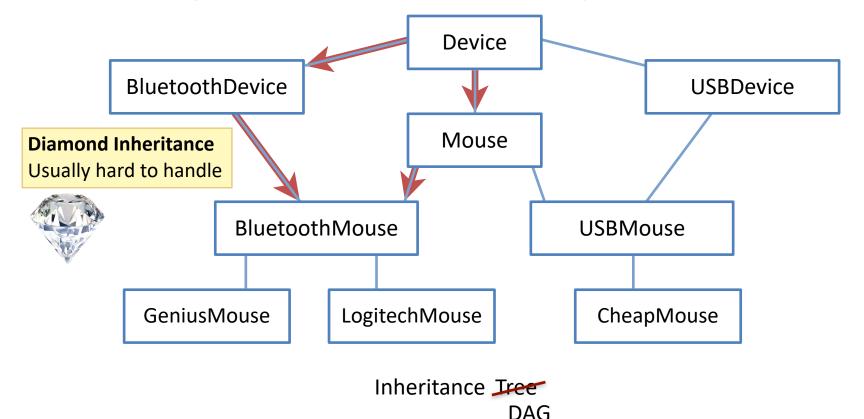


Handled with corresponding typing rule

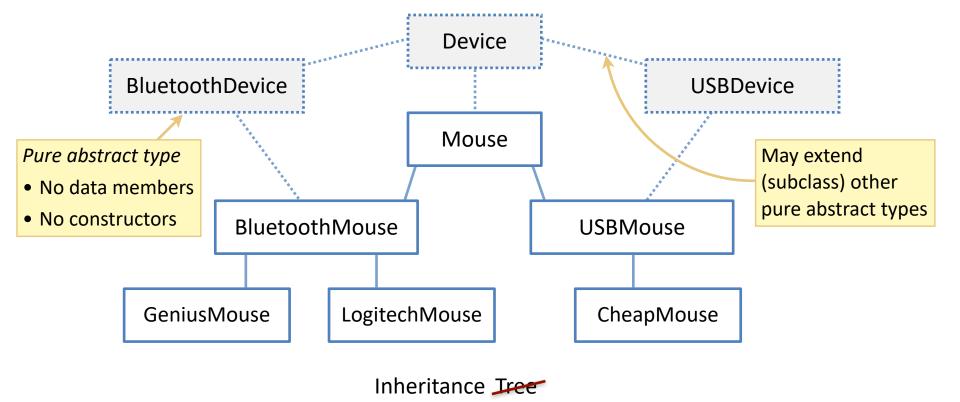
$$\frac{\mathsf{E}_1:\mathsf{T}\quad \mathsf{T}<:\mathsf{S}}{\mathsf{E}_1:\mathsf{S}}$$

- As a consequence:
 - Each term has more than one type
 - Need to find the appropriate type for each context

- A basic concept in Object-oriented Programming
 - Every class has (can have) a superclass



 Marking some of the types as pure abstract helps alleviate some of the difficulty in multiple inheritance



DAG

- Upcast
 - ➤ The compiler will insert an implicit conversion from subtype to supertype (similar to coercions)

```
Mouse m = new GeniusMouse();
registerDevice(m);
```





Summary

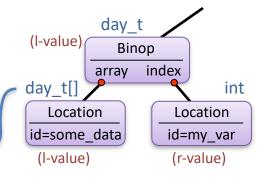
Static semantic checks

- Identification matches applied occurrences of identifier to its defining occurrence
 - The symbol table maintains this information

- Type checking checks which type combinations are legal
 - Each node in the AST of an expression represents either an I-value (location) or an r-value (value) of some type

name	pos	type
month	1	int 112
month_name	2	string[112]
day	19	day_t

name	constructor	structure
int	builtin	-
size_t	typedef	unsigned int
day_t	enum	0: YESTERDAY 1: TODAY 2: TOMORROW



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
 - with a few tricks, attributes can be computed while building the AST

