

5. Semantic Analysis

Mircea Lungu

Oscar Nierstrasz

Thanks to Jens Palsberg and Tony Hosking for their kind permission to reuse and adapt the CS132 and CS502 lecture notes.

<http://www.cs.ucla.edu/~palsberg/>

<http://www.cs.purdue.edu/homes/hosking/>

**William Tell is a folk hero of Switzerland;
she was an exceptional marksman.**

Conference in Vienna in 1964 best summarized by T. B. Steel:

“I don’t fully know myself how to describe the semantics of a language. I daresay nobody does or we wouldn’t be here”

The Genesis of Attribute Grammars

Donald E. Knuth

Roadmap



- > Context-sensitive analysis
- > Strategies for semantic analysis
- > Attribute grammars
- > Symbol tables and type-checking

See, *Modern compiler implementation in Java* (Second edition), Chapter 5.

Roadmap



- > **Context-sensitive analysis**
- > Strategies for semantic analysis
- > Attribute grammars
- > Symbol tables and type-checking

```
{  
...  
  x := y [1];  
  z := x + y;  
...  
}
```

Semantic Analysis

Compilation is driven by the syntactic structure of the program as discovered by the parser

Semantic routines:

- interpret meaning of the program based on its syntactic structure
- two purposes:
 - *finish **analysis** by deriving context-sensitive information*
 - *begin **synthesis** by generating the IR or target code*
- associated with individual productions of a context free grammar or sub-trees of a syntax tree

Context-sensitive analysis

What context-sensitive questions might the compiler ask?

1. Is x scalar, an array, or a function?
2. Is x declared before it is used?
3. Are any names declared but not used?
4. Which declaration of x is being referenced?
5. Is an expression type-consistent?
6. Does the dimension of a reference match the declaration?
7. Where can x be stored? (heap, stack, ...)
8. Does $*p$ reference the result of a `malloc()`?
9. Is x defined before it is used?
10. Is an array reference in bounds?
11. Does function `foo` produce a constant value?
12. Can p be implemented as a memo-function?

These questions cannot be answered with a context-free grammar

Context-sensitive analysis

- > What are the challenges?
 - questions and answers involve non-local information
 - answers depend on values, not syntax
 - answers may involve computation
- > Several approaches:
 - symbol tables*: central store for facts; express checking code
 - attribute grammars*: specify non-local computations; automatic evaluators
 - language design*: simplify language; avoid problems

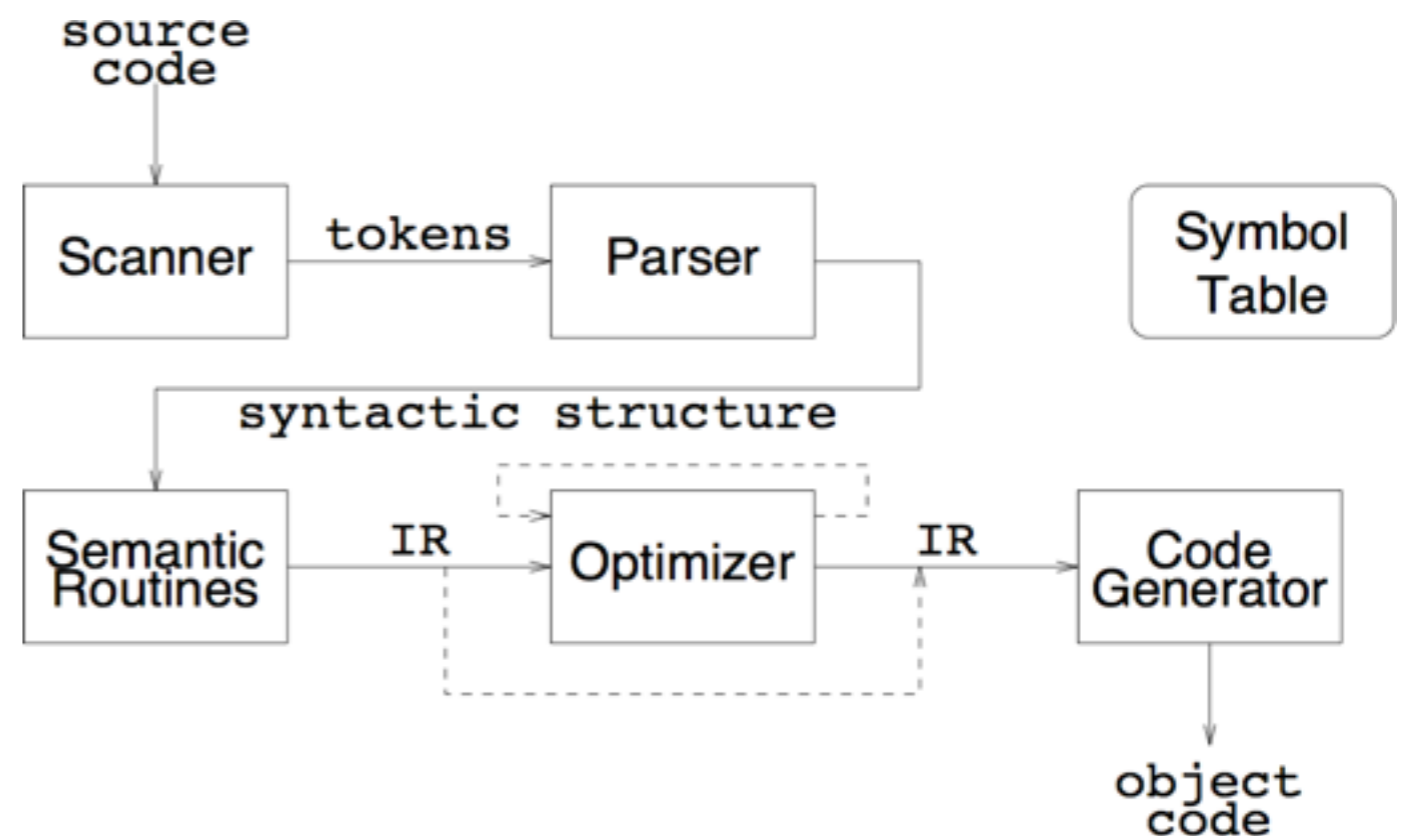
Roadmap



- > Context-sensitive analysis
- > **Strategies for semantic analysis**
- > Attribute grammars
- > Symbol tables and type-checking

Alternatives for semantic processing

- One-pass compiler and synthesis
- Two-pass compiler + peephole compiler & IR synthesis + code generation pass
- Multi-pass analysis synthesis



One-pass compilers

- > Interleave scanning, parsing and translation
 - no explicit IR
 - generate target code directly
 - *emit short sequences of instructions on each parser action*
 - *little or no optimization possible (minimal context)*
- > Can add *peephole optimization pass*
 - extra pass over generated code through small window (“peephole”) of instructions
 - smooths out “rough edges” between code emitted by subsequent calls to code generator

Two-pass: analysis & IR synthesis + code generation

- > Generate explicit IR as interface to code generator
 - linear (e.g., tuples)
 - can emit multiple tuples at a time for better code context
- > Advantages
 - easier retargeting (IR must be expressive enough for different machines!)
 - can add optimization pass later (multi-pass synthesis)

Multi-pass analysis

> Several passes, read/write intermediate files

1. scan source file, generate tokens
 - *place identifiers and constants in symbol table*
2. parse token file
 - *generate semantic actions or linearized parse tree*
3. process declarations to symbol table
4. semantic checking with IR synthesis

> Motivations:

- Historical: constrained address spaces
- Language: *e.g.*, declaration after use
- Multiple analyses over IR tree

Multi-pass synthesis

- > Passes operate on linear or tree-structured IR
- > Options:
 - code generation and peephole optimization
 - multi-pass IR transformation
 - *machine-independent then dependent optimizations*
 - high-level to low-level IR transformation before code generation
 - *e.g., in **gcc** high-level trees drive generation of low-level Register Transfer Language for machine-independent optimization*
 - language-independent front ends
 - retargetable back ends

Roadmap



- > Context-sensitive analysis
- > Strategies for semantic analysis
- > **Attribute grammars**
- > Symbol tables and type-checking

Attribute grammars

- > Add attributes to the syntax tree or PEG:
 - can add attributes (fields) to each node
 - specify equations to define values
 - propagate values up (synthesis) or down (inheritance)
- > **Example:** ensuring that constants are immutable
 - add *type* and *class* attributes to expression nodes
 - add rules to production for `:=`
 1. *check that LHS.class is variable (not constant)*
 2. *check that LHS.type and RHS.type are compatible*

Attribute grammar actions

PRODUCTION	SEMANTIC RULES
$D \rightarrow T L$	$L.in := T.type$
$T \rightarrow \mathbf{int}$	$T.type := \text{integer}$
$T \rightarrow \mathbf{real}$	$T.type := \text{real}$
$L \rightarrow L_1 , \mathbf{id}$	$L_1.in := L.in$ $\text{addtype}(\mathbf{id}.entry, L.in)$
$L \rightarrow \mathbf{id}$	$\text{addtype}(\mathbf{id}.entry, L.in)$

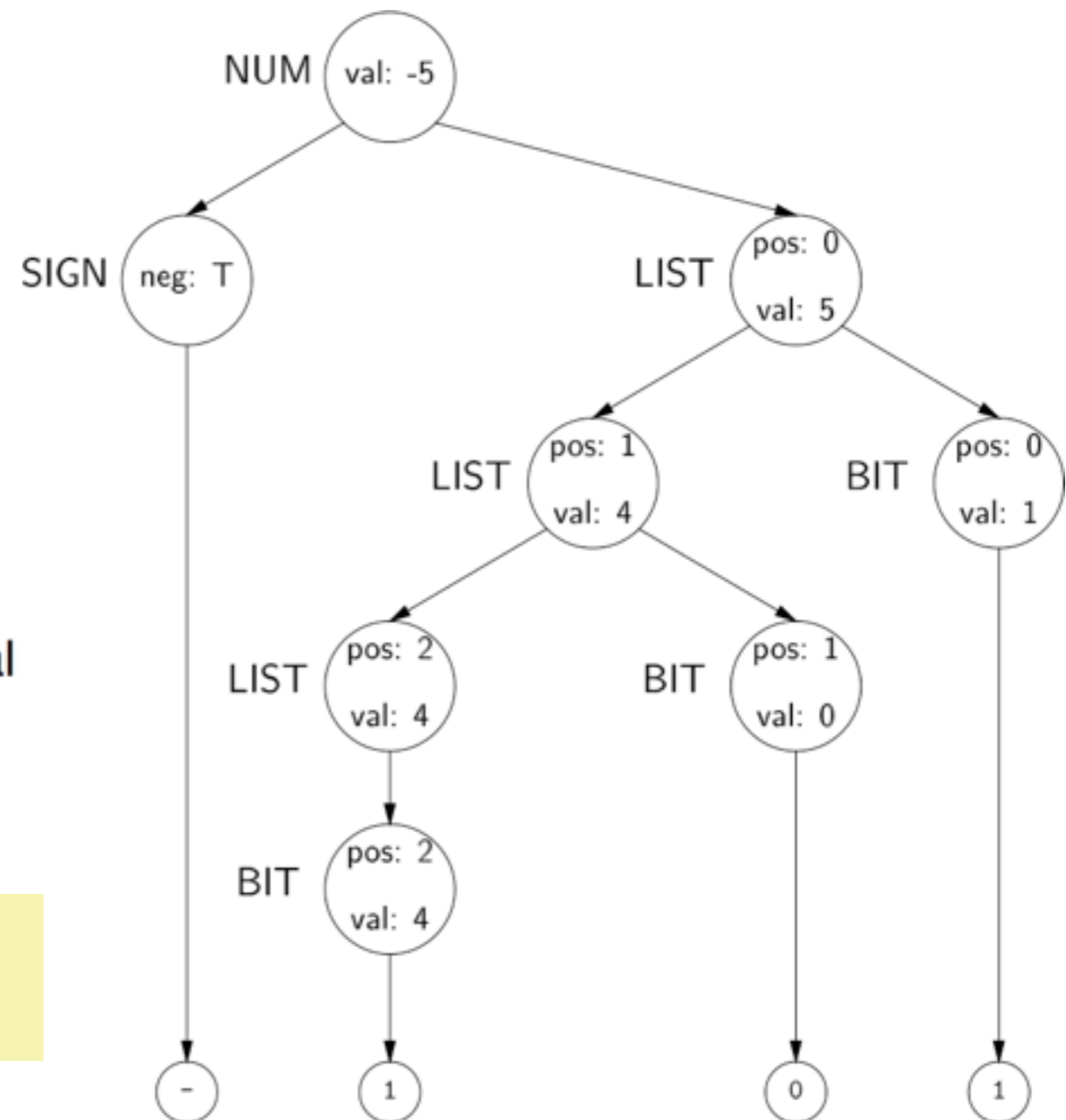
- > tree attributes specified by grammar
- > productions associated with attribute assignments
- > each attribute defined uniquely and locally
- > identical terms are labeled uniquely

Example: evaluating signed binary numbers

PRODUCTION	SEMANTIC RULES
NUM → SIGN LIST	LIST.pos := 0 if SIGN.neg NUM.val := -LIST.val else NUM.val := LIST.val
SIGN → +	SIGN.neg := false
SIGN → -	SIGN.neg := true
LIST → BIT	BIT.pos := LIST.pos LIST.val := BIT.val
LIST → LIST ₁ BIT	LIST ₁ .pos := LIST.pos + 1 BIT.pos := LIST.pos LIST.val := LIST ₁ .val + BIT.val
BIT → 0	BIT.val := 0
BIT → 1	BIT.val := 2 ^{BIT.pos}

- *val* and *neg* are synthetic attributes
- *pos* is an inherited attribute

Attributed parse tree for -101

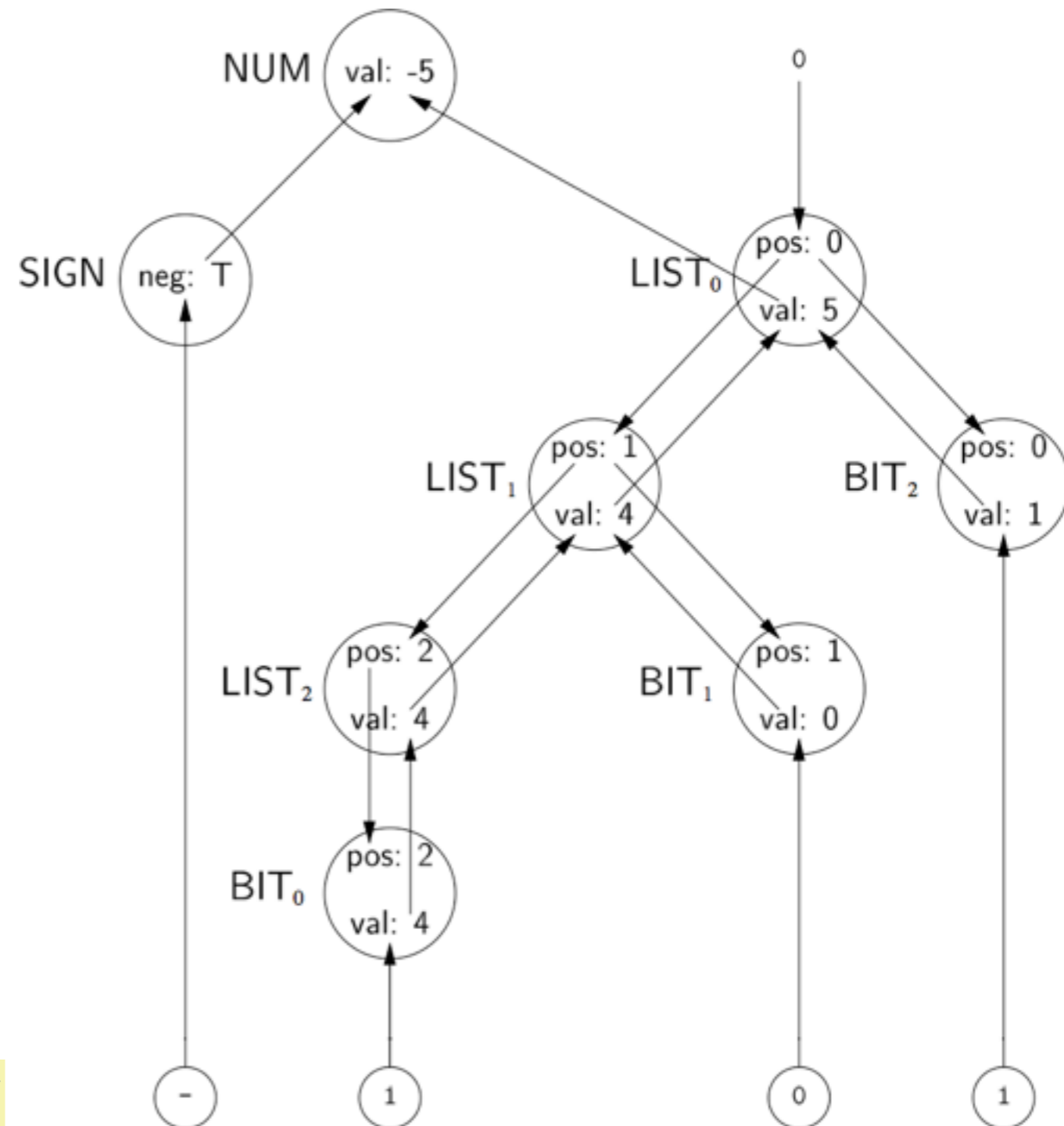


Attribute dependency graph

- *nodes* represent *attributes*
- *edges* represent *flow of values*
- graph must be *acyclic*
- *topologically sort* to order attributes
 - use this order to evaluate rules
 - order depends on both grammar and input string!

- | | |
|---------------------------|----------------------------|
| 1. SIGN.neg | 8. BIT ₀ .val |
| 2. LIST ₀ .pos | 9. LIST ₂ .val |
| 3. LIST ₁ .pos | 10. BIT ₁ .val |
| 4. LIST ₂ .pos | 11. LIST ₁ .val |
| 5. BIT ₀ .pos | 12. BIT ₂ .val |
| 6. BIT ₁ .pos | 13. LIST ₀ .val |
| 7. BIT ₂ .pos | 14. NUM.val |

Evaluating in this order yields NUM.val = -5



Evaluation strategies

> ***Parse-tree methods***

1. build the parse tree
2. build the dependency graph
3. topologically sort the graph
4. evaluate it

> ***Rule-based methods***

1. analyse semantic rules at compiler-construction time
2. determine static ordering for each production's attributes
3. evaluate its attributes in that order at compile time

> ***Oblivious methods***

1. ignore the parse tree and the grammar
2. choose a convenient order (e.g., left-to-right traversal) and use it
3. repeat traversal until no more attribute values can be generated

Attribute grammars in practice

> *Advantages*

- clean formalism
- automatic generation of evaluator
- high-level specification

> *Disadvantages*

- evaluation strategy determines efficiency
- increase space requirements
- parse tree evaluators need dependency graph
- results distributed over tree
- circularity testing

Historically, attribute grammars have been judged too large and expensive for industrial-strength compilers.

Roadmap



- > Context-sensitive analysis
- > Strategies for semantic analysis
- > Attribute grammars
- > **Symbol tables and type-checking**

Symbol tables

- > For compile-time efficiency, compilers often use a symbol table:
 - associates lexical *names* (symbols) with their *attributes*
- > What items should be entered?
 - variable names
 - constants
 - procedure and function names
 - literal constants and strings
 - compiler-generated temporaries (we'll get there)
- > Separate table of structure layouts for types (field offsets and lengths)

A symbol table is a compile-time structure

Symbol table information

- > What kind of information might the compiler need?
 - textual name
 - data type
 - dimension information (*for aggregates*)
 - declaring procedure
 - lexical level** of declaration
 - storage class (*heap, stack, text ...*)
 - offset in storage
 - if record, pointer to structure table
 - if parameter, by-reference or by-value?
 - can it be aliased? to what other names?
 - number and type of arguments to functions

Lexical Scoping

```
class C {  
  int x;  
  void m(int y) {  
    int z;  
    if (y>x) {  
      int w=z+y;  
      return w;  
    }  
    return y;  
  }  
}
```

With lexical scoping the definition of a name is determined by its *static scope*. A stack suffices to track the current definitions.

scope of y and z

scope of w

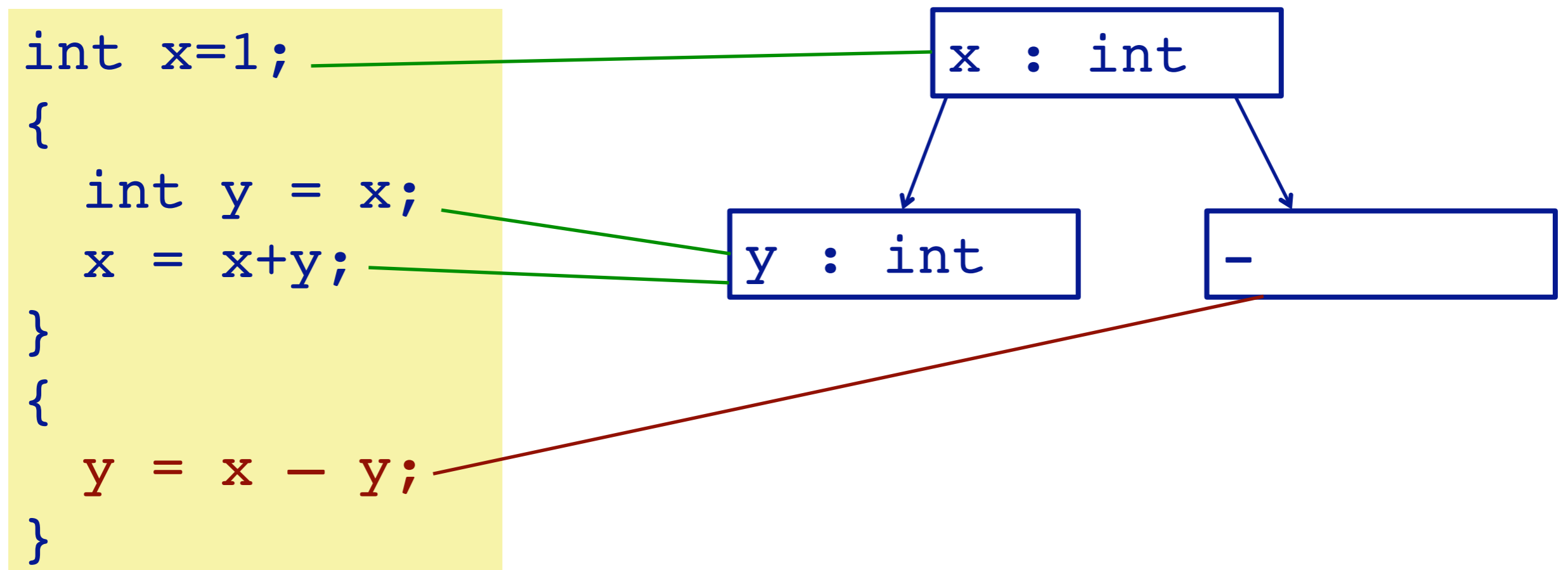
scope of x

Nested scopes: block-structured symbol tables

- > What information is needed?
 - when we ask about a name, we want the *most recent declaration*
 - the declaration may be from the current scope or some enclosing scope
 - innermost scope overrides declarations from outer scopes
- > Key point: *new declarations (usually) occur only in current scope*
- > What operations do we need?
 - `void put(Symbol key, Object value)` – bind key to value
 - `Object get(Symbol key)` – return value bound to key
 - `void beginScope()` – remember current state of table
 - `void endScope()` – restore table to state at most recent scope that has not been ended

May need to preserve list of locals for the debugger

Checking variable declarations in a hierarchical symbol table



Efficient Implementation of Symbol Tables

Implementation options

1. functional
2. imperative

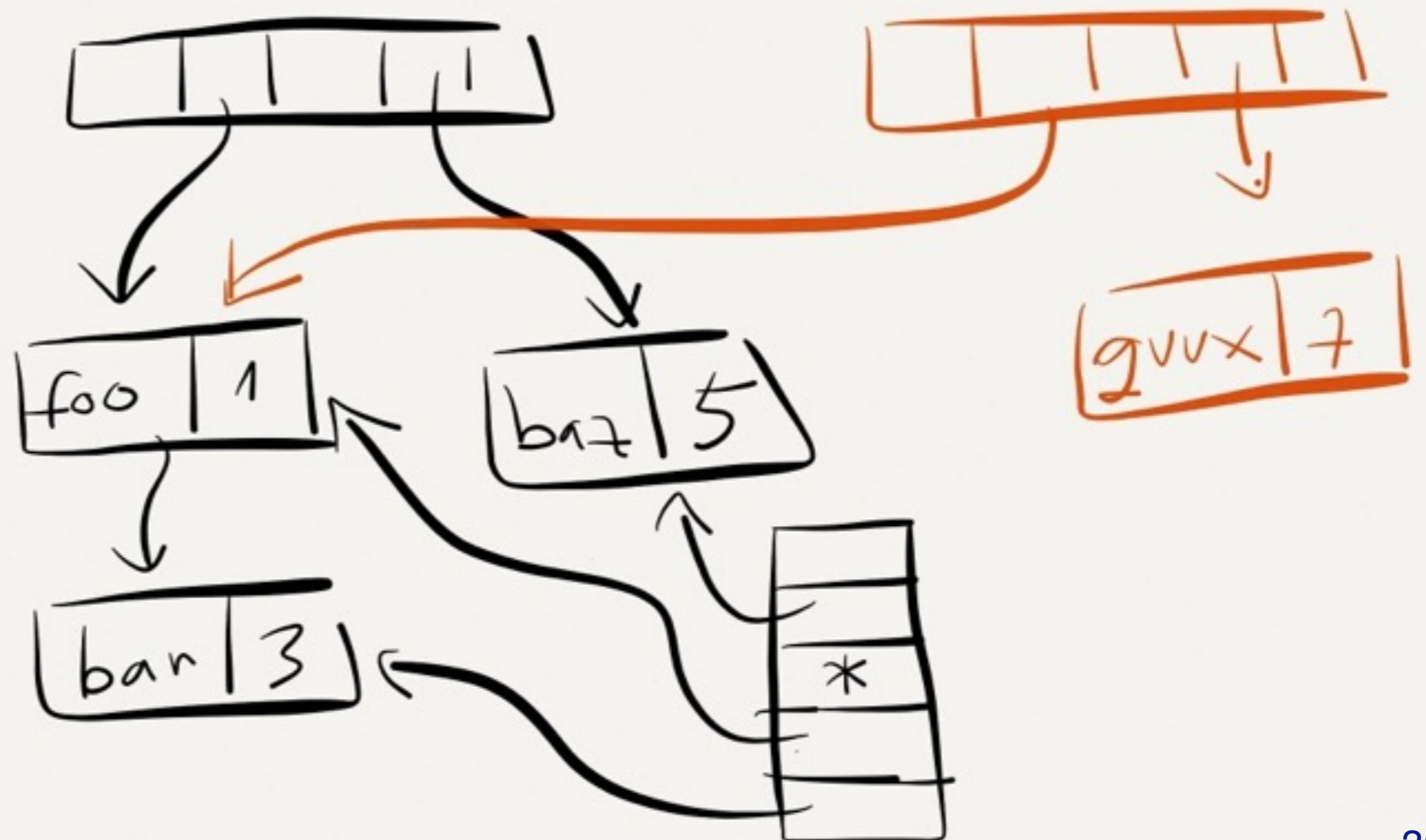
How to ensure efficiency, with thousands of distinct identifiers in a large program?

Efficient Implementation of Symbol Tables

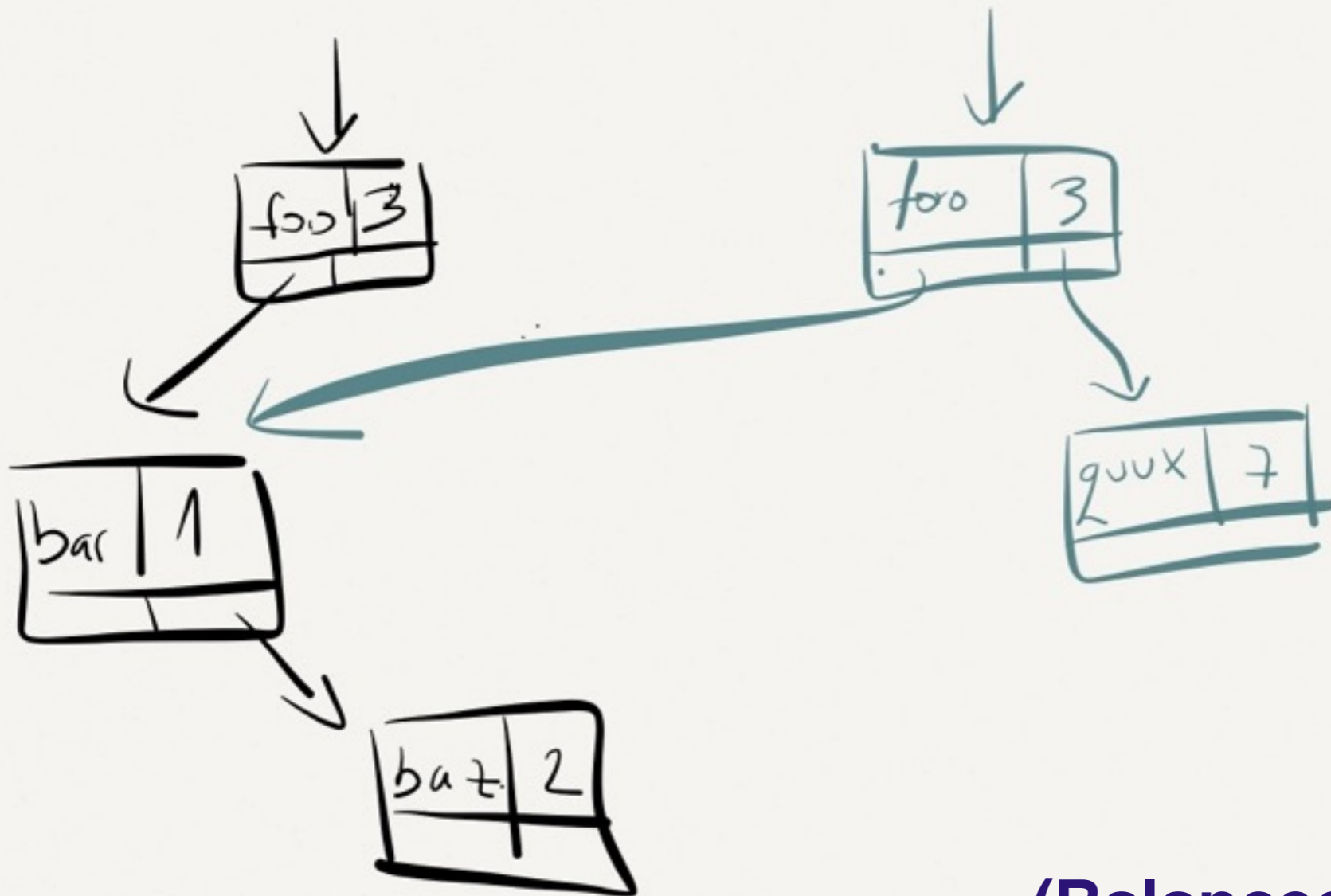
```
int foo, bar;
foo = ++bar;
if (bar > 10) then
{
    boolean baz;
    baz = true;
}

// and assume
hash(foo) = hash(bar)
hash(baz) = hash(quux)
```

Hash tables support an **imperative** (destructive) implementation



Efficient Implementation of Symbol Tables (2)



(Balanced) binary trees support a functional (non-destructive) implementation. A persistence data structure.

Attribute information

- > Attributes are internal representations of declarations
- > Symbol table associates names with attributes

- > Names may have different attributes depending on their meaning:
 - *variables*: type, procedure level, frame offset
 - *types*: type descriptor, data size/alignment
 - *constants*: type, value
 - *procedures*: formals (names/types), result type, block information (local decls.), frame size

Static and Dynamic Typing

A language is statically typed if it is always possible to *determine the (static) type* of an expression *based on the program text alone*.

A language is dynamically typed if *only values have fixed type*. Variables and parameters may take on different types at run-time, and must be checked immediately before they are used.

A language is “strongly typed” if it is impossible to perform an operation on the wrong kind of object.

Type consistency may be assured by

- I. compile-time type-checking,
- II. type inference, or
- III. dynamic type-checking.

See: Programming Languages course

Type expressions

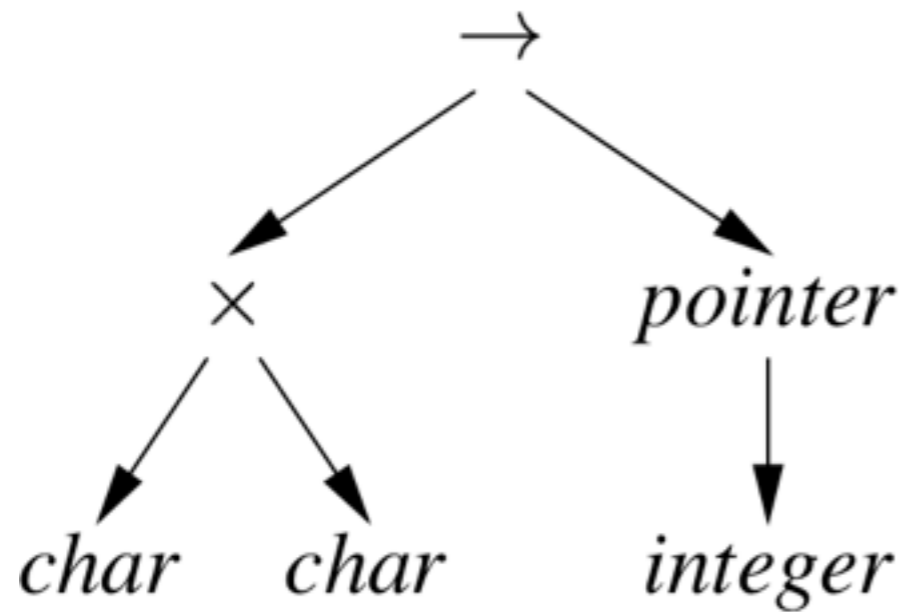
Type expressions are a textual representation for types:

1. basic types: *boolean, char, integer, real*, etc.
2. type names
3. constructed types (constructors applied to type expressions):
 - a) $\text{array}(I, T)$ denotes *array* of elements type T , index type I
e.g., $\text{array}(1 \dots 10, \text{integer})$
 - b) $T_1 \times T_2$ denotes *Cartesian product* of type expressions T_1 and T_2
 - c) $\text{record}(\dots)$ denotes *record* with named fields
e.g., $\text{record}((a \times \text{integer}), (b \times \text{real}))$
 - d) $\text{pointer}(T)$ denotes the type “*pointer* to object of type T ”
 - e) $D \rightarrow R$ denotes type of *function* mapping domain D to range R
e.g., $\text{integer} \times \text{integer} \rightarrow \text{integer}$

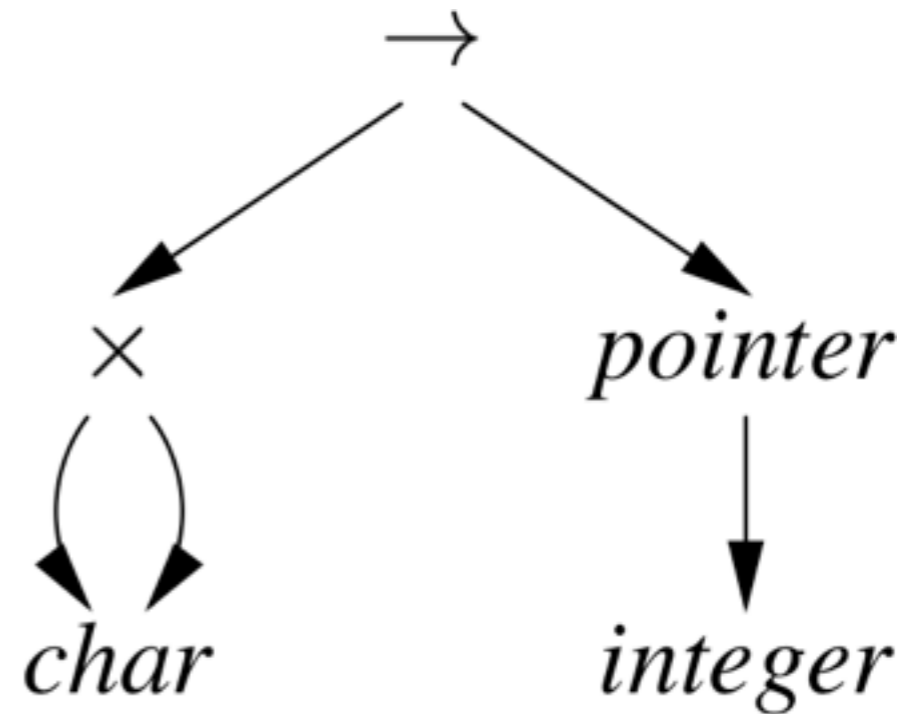
Type descriptors

Type descriptors are **compile-time structures** representing type expressions

e.g., $char \times char \rightarrow pointer(integer)$



or



Type compatibility

Type checking needs to determine *type equivalence*

Two approaches:

- > *Name equivalence*: each type name is a distinct type
- > *Structural equivalence*: two types are equivalent iff they have **the same structure** (after substituting type expressions for type names)
 - $s \equiv t$ iff s and t are the same basic types
 - $\text{array}(s_1, s_2) \equiv \text{array}(t_1, t_2)$ iff $s_1 \equiv t_1$ and $s_2 \equiv t_2$
 - $s_1 \times s_2 \equiv t_1 \times t_2$ iff $s_1 \equiv t_1$ and $s_2 \equiv t_2$
 - $\text{pointer}(s) \equiv \text{pointer}(t)$ iff $s \equiv t$
 - $s_1 \rightarrow s_2 \equiv t_1 \rightarrow t_2$ iff $s_1 \equiv t_1$ and $s_2 \equiv t_2$

Type compatibility: example

Consider:

```
type link = ^cell
var next : link;
var last : link;
var p : ^cell;
var q, r : ^cell;
```

Under name equivalence:

- next and last have the same type
- p, q and r have the same type
- p and next have different type

Under structural equivalence all variables have the same type

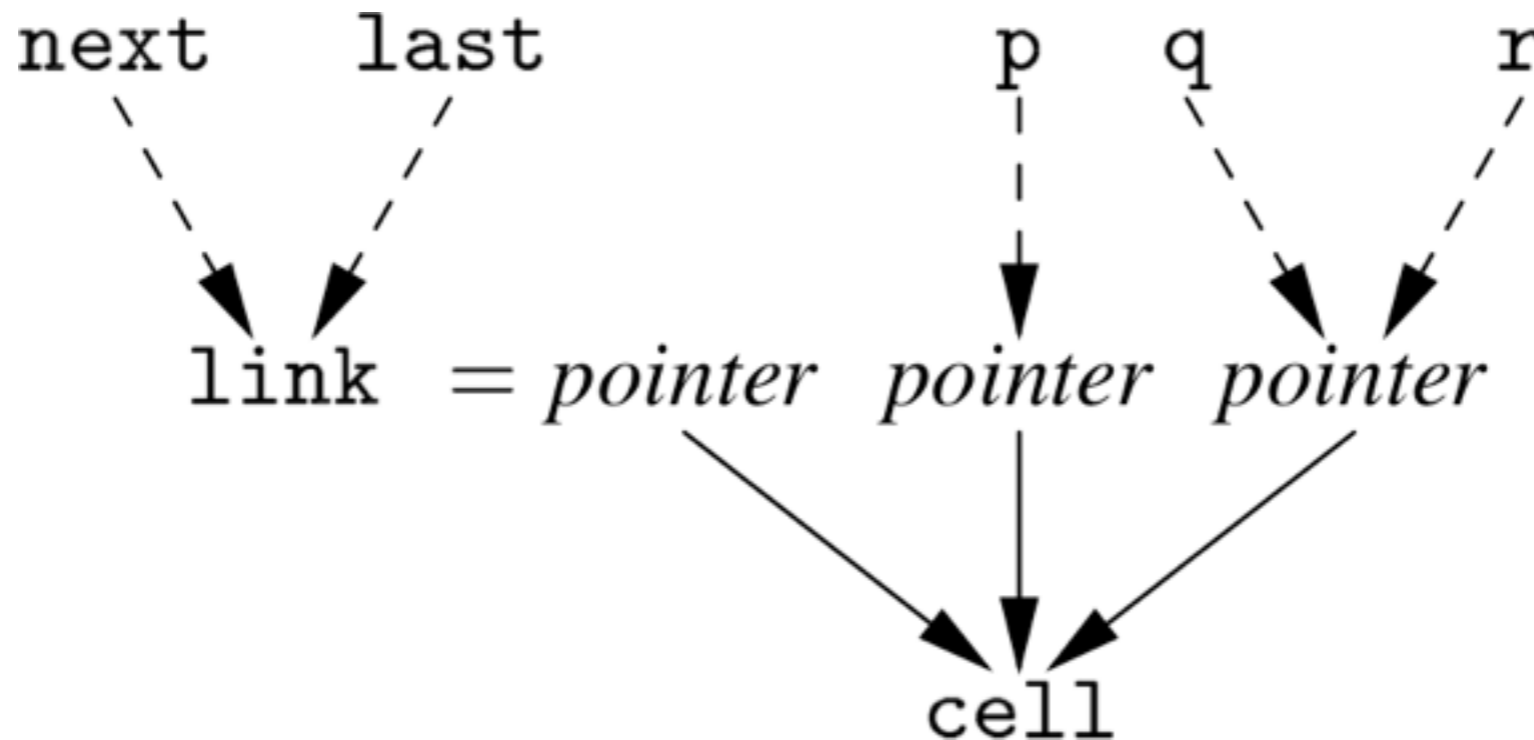
Ada/Pascal/Modula-2 are somewhat confusing: they treat distinct type definitions as distinct types, so

- p has different type from q and r (!)

Type compatibility: Pascal-style name equivalence

Build compile-time structure called a ***type graph***:

- each constructor or basic type creates a node
- each name creates a leaf (associated with the type's descriptor)



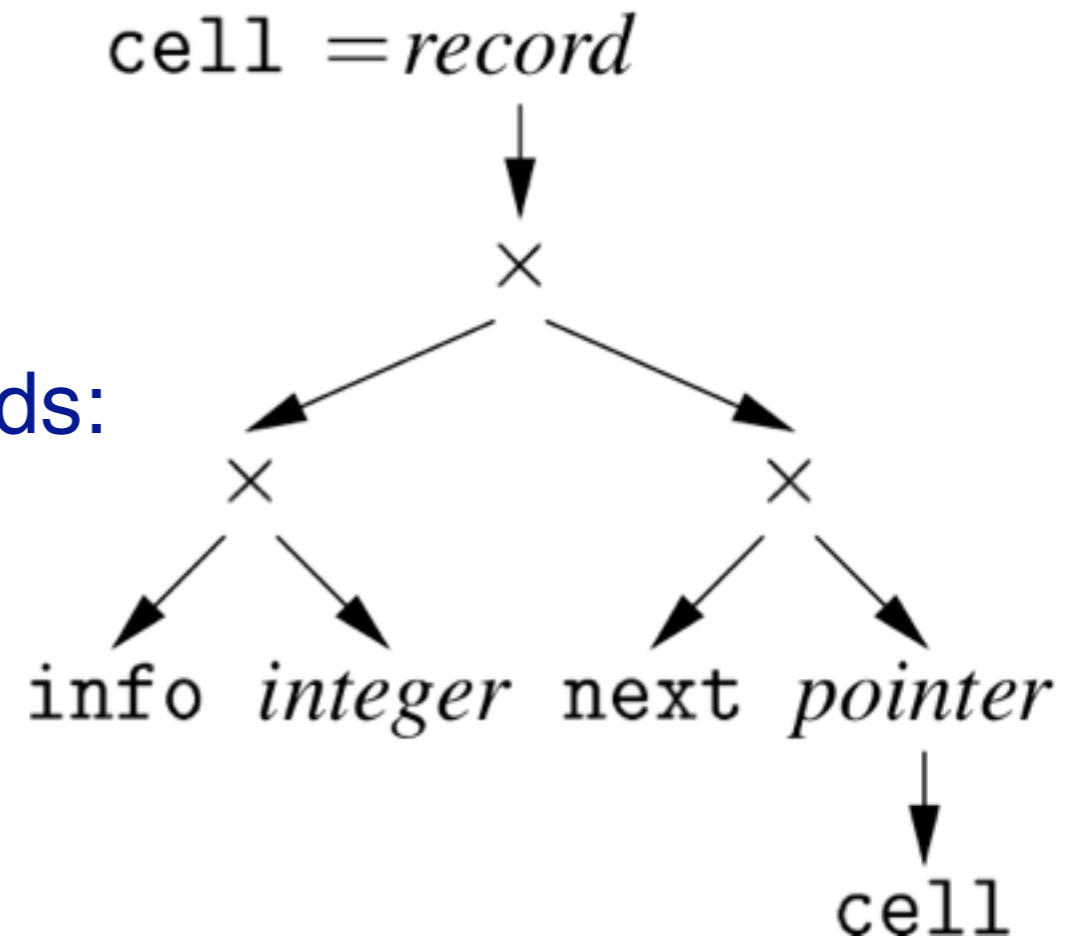
Type expressions are equivalent if they are represented by the same node in the graph

Type compatibility: recursive types

Consider:

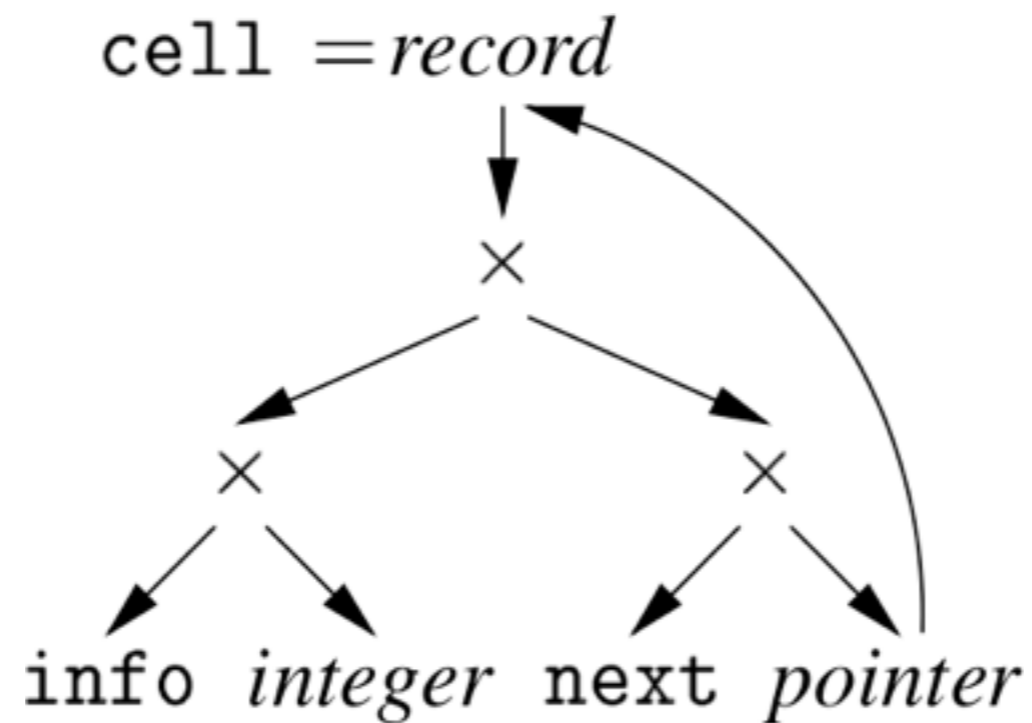
```
type link = ^cell
var cell = record
    info : integer;
    next : link;
end
```

Expanding `link` in the type graph yields:



Type compatibility: recursive types

Allowing cycles in the type graph eliminates `cell`:



Type rules

Type-checking rules can be formalized to prove soundness and correctness.

$$\frac{f : A \rightarrow B, x : A}{f(x) : B}$$

If f is a function from A to B , and x is of type A , then $f(x)$ is a value of type B .






Example: Featherweight Java

Syntax:	Expression typing:	
$\text{CL} ::= \text{class } C \text{ extends } C \{ \bar{C} \bar{f}; K \bar{M} \}$	$\Gamma \vdash x \in \Gamma(x) \quad (\text{T-VAR})$	
$K ::= C(\bar{C} \bar{f}) \{ \text{super}(\bar{f}); \text{this}.\bar{f} = \bar{f}; \}$	$\frac{\Gamma \vdash e_0 \in C_0 \quad \text{fields}(C_0) = \bar{C} \bar{f}}{\Gamma \vdash e_0.f_i \in C_i} \quad (\text{T-FIELD})$	
$M ::= C m(\bar{C} \bar{x}) \{ \text{return } e; \}$	$\frac{\Gamma \vdash e_0 \in C_0 \quad \text{mtype}(m, C_0) = \bar{D} \rightarrow C \quad \Gamma \vdash \bar{e} \in \bar{C} \quad \bar{C} \triangleleft \bar{D}}{\Gamma \vdash e_0.m(\bar{e}) \in C} \quad (\text{T-INVK})$	
$e ::= \begin{array}{l} x \\ e.f \\ e.m(\bar{e}) \\ \text{new } C(\bar{e}) \\ (C)e \end{array}$	$\frac{\text{fields}(C) = \bar{D} \bar{f} \quad \Gamma \vdash \bar{e} \in \bar{C} \quad \bar{C} \triangleleft \bar{D}}{\Gamma \vdash \text{new } C(\bar{e}) \in C} \quad (\text{T-NEW})$	
Subtyping:	$\frac{\Gamma \vdash e_0 \in D \quad D \triangleleft C \quad C \neq D}{\Gamma \vdash (C)e_0 \in C} \quad (\text{T-UCAST})$	
$C \triangleleft C$	$\frac{\Gamma \vdash e_0 \in D \quad C \not\triangleleft D \quad D \not\triangleleft C \quad \text{stupid warning}}{\Gamma \vdash (C)e_0 \in C} \quad (\text{T-SCAST})$	
$\frac{C \triangleleft D \quad D \triangleleft E}{C \triangleleft E}$	<th>Method typing:</th>	Method typing:
$\frac{CT(C) = \text{class } C \text{ extends } D \{ \dots \}}{C \triangleleft D}$	$\frac{\bar{x} : \bar{C}, \text{this} : C \vdash e_0 \in E_0 \quad E_0 \triangleleft C_0 \quad CT(C) = \text{class } C \text{ extends } D \{ \dots \} \quad \text{override}(m, D, \bar{C} \rightarrow C_0)}{C_0 m(\bar{C} \bar{x}) \{ \text{return } e_0; \} \text{ OK IN } C}$	
Computation:	Class typing:	
$\frac{\text{fields}(C) = \bar{C} \bar{f}}{(\text{new } C(\bar{e})) . f_i \rightarrow e_i} \quad (\text{R-FIELD})$	$K = C(\bar{D} \bar{g}, \bar{C} \bar{f}) \{ \text{super}(\bar{g}); \text{this}.\bar{f} = \bar{f}; \}$	
$\frac{\text{mbody}(m, C) = (\bar{x}, e_0)}{(\text{new } C(\bar{e})) . m(\bar{d}) \rightarrow [\bar{d}/\bar{x}, \text{new } C(\bar{e})/\text{this}]e_0} \quad (\text{R-INVK})$	$\frac{\text{fields}(D) = \bar{D} \bar{g} \quad M \text{ OK IN } C}{\text{class } C \text{ extends } D \{ \bar{C} \bar{f}; K \bar{M} \} \text{ OK}}$	
$\frac{C \triangleleft D}{(D)(\text{new } C(\bar{e})) \rightarrow \text{new } C(\bar{e})} \quad (\text{R-CAST})$		





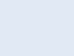

Used to prove that generics could be added to Java without breaking the type system.

Igarashi, Pierce and Wadler, "Featherweight Java: a minimal core calculus for Java and GJ", OOPSLA '99
doi.acm.org/10.1145/320384.320395

Can you answer these questions?

-  Why can semantic analysis be performed by the parser?*
-  What are the pros and cons of introducing an IR?*
-  Why must an attribute dependency graph be acyclic?*
-  Why would be the use of a symbol table at run-time?*
-  Why does Java adopt nominal (name-based) rather than structural type rules?*

What you should know!

-  Why is semantic analysis mostly context-sensitive?*
-  What is “peephole optimization”?*
-  Why was multi-pass semantic analysis introduced?*
-  What is an attribute grammar? How can it be used to support semantic analysis?*
-  What kind of information is stored in a symbol table?*
-  How is type-checking performed?*



Attribution-ShareAlike 4.0 International (CC BY-SA 4.0)

You are free to:

Share — copy and redistribute the material in any medium or format

Adapt — remix, transform, and build upon the material for any purpose, even commercially.

The licensor cannot revoke these freedoms as long as you follow the license terms.

Under the following terms:



Attribution — You must give appropriate credit, provide a link to the license, and indicate if changes were made. You may do so in any reasonable manner, but not in any way that suggests the licensor endorses you or your use.



ShareAlike — If you remix, transform, or build upon the material, you must distribute your contributions under the same license as the original.

No additional restrictions — You may not apply legal terms or technological measures that legally restrict others from doing anything the license permits.

<http://creativecommons.org/licenses/by-sa/4.0/>