

Lecture 15

# COMPILER DESIGN

# Announcements

- **HW4:** OAT v. 1.0
  - Parsing & basic code generation

# Why Inference Rules?

- Allow a compact, precise way of specifying language properties
  - About 20 pages for full Java, versus
  - Hundreds of pages of prose Java Language Specification
- Correspond closely to recursive AST traversal for implementing them
- Type checking/inference tries to prove a different judgment  $C \vdash e : t$ 
  - By searching backward through the rules
- Compiling is “*interpreting*” the type checking rules  $[[C \vdash e : t]]$ 
  - Compilation follows the type checking judgment
- Strong mathematical foundations, e.g., “Curry-Howard correspondence”
  - Programming Language ~ Logic
  - Program ~ Proof
  - Type ~ Proposition



# OPTIMIZING CONTROL

# Standard Evaluation

- Consider compiling the following program fragment

```
if (x & !y | !w)
    z = 3;
else
    z = 4;
return z;
```



```
%tmp1 = icmp Eq [[y]], 0      ; !y
%tmp2 = and [[x]] %tmp1
%tmp3 = icmp Eq [[w]], 0
%tmp4 = or %tmp2, %tmp3
%tmp5 = icmp Eq %tmp4, 0
br %tmp5, label %else, label %then
```

```
then:
    store [[z]], 3
    br %merge
```

```
else:
    store [[z]], 4
    br %merge
```

```
merge:
    %tmp6 = load [[z]]
    ret %tmp6
```

# Observation

- Usually, we want the translation  $\llbracket e \rrbracket$  to produce a value
  - $\llbracket C \vdash e : t \rrbracket = (\text{ty}, \text{operand}, \text{stream})$
  - e.g.  $\llbracket C \vdash e_1 + e_2 : \text{int} \rrbracket = (\text{i64}, \%tmp, [\%tmp = \text{add } \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket])$
- But, when the compiled expression appears in a test
  - The program jumps to one label or another after the comparison
  - Otherwise, it never uses the value
- In many cases, we can avoid “materializing” the value (i.e. storing it in a temporary) and thus produce better code
  - This idea also lets us implement different functionality too:  
e.g. short-circuiting Boolean expressions

# Idea: Use a different translation for tests

- Usual Expression translation:  $\llbracket C \vdash e : t \rrbracket = (\text{ty}, \text{operand}, \text{stream})$
- Conditional branch translation of Booleans, without materializing value  
 $\llbracket C \vdash e : \text{bool}@ \rrbracket \text{ ltrue lfalse} = \text{stream}$

$\llbracket C, \text{rt} \vdash \text{if}(e) \text{ then } s_1 \text{ else } s_2 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket,$

```
    insns3
  then:
    [[s1]]
    br %merge
  else:
    [[s2]]
    br %merge
  merge:
```

## Notes

- Two extra arguments
  - “true” branch label
  - “false” branch label
- Doesn’t “return a value”
- Aside: this is a form of continuation-passing translation

where

$\llbracket C, \text{rt} \vdash s_1 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket, \text{ insns}_1$

$\llbracket C, \text{rt} \vdash s_2 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket, \text{ insns}_2$

$\llbracket C \vdash e : \text{bool}@ \rrbracket \text{ then else} = \text{insns}_3$

# Short Circuit Compilation: Expressions

- $\llbracket C \vdash e : \text{bool}@ \rrbracket \text{ ltrue } \text{ lfalse} = \text{ insns}$

$$\frac{}{\llbracket C \vdash \text{false} : \text{bool}@ \rrbracket \text{ ltrue } \text{ lfalse} = [\text{br } \% \text{ lfalse}]} \text{ FALSE}$$
$$\frac{}{\llbracket C \vdash \text{true} : \text{bool}@ \rrbracket \text{ ltrue } \text{ lfalse} = [\text{br } \% \text{ ltrue}]} \text{ TRUE}$$
$$\frac{\llbracket C \vdash e : \text{bool}@ \rrbracket \text{ lfalse } \text{ ltrue} = \text{ insns}}{\llbracket C \vdash !e : \text{bool}@ \rrbracket \text{ ltrue } \text{ lfalse} = \text{ insns}} \text{ NOT}$$



# Short Circuit Evaluation

Idea: build the logic into the translation

$$\llbracket C \vdash e1 : \text{bool}@ \rrbracket \text{ ltrue right} = \text{insns}_1 \quad \llbracket C \vdash e2 : \text{bool}@ \rrbracket \text{ ltrue lfalse} = \text{insns}_2$$
$$\llbracket C \vdash e1 \mid e2 : \text{bool}@ \rrbracket \text{ ltrue lfalse} =$$

```
insns1
right:
insns2
```

$$\llbracket C \vdash e1 : \text{bool}@ \rrbracket \text{ right lfalse} = \text{insns}_1 \quad \llbracket C \vdash e2 : \text{bool}@ \rrbracket \text{ ltrue lfalse} = \text{insns}_2$$
$$\llbracket C \vdash e1 \& e2 : \text{bool}@ \rrbracket \text{ ltrue lfalse} =$$

```
insns1
right:
insns2
```

where `right` is a fresh label

# Short-Circuit Evaluation

- Consider compiling the following program fragment

```
if (x & !y | !w)
    z = 3;
else
    z = 4;
return z;
```



```
%tmp1 = icmp Eq [[x]], 0
br %tmp1, label %right2, label %right1

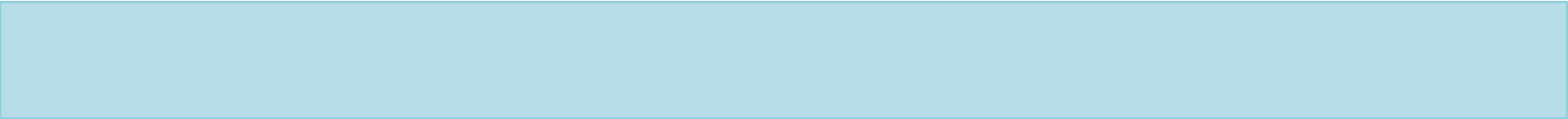
right1:
    %tmp2 = icmp Eq [[y]], 0
    br %tmp2, label %then, label %right2

right2:
    %tmp3 = icmp Eq [[w]], 0
    br %tmp3, label %then, label %else

then:
    store [[z]], 3
    br %merge

else:
    store [[z]], 4
    br %merge

merge:
    %tmp5 = load [[z]]
    ret %tmp5
```



Compiling lambda calculus to straight-line code.  
Representing evaluation environments at runtime.

# CLOSURE CONVERSION

# “Functional” languages

- Languages like ML, Haskell, Scheme, Python, C#, Java 8, Swift
  - Functions can be passed as arguments (e.g. map or fold)
  - Functions can be returned as values (e.g. compose)
  - **Function nest:**
    - Inner function can refer to variables bound in the outer function

```
let add = fun x -> fun y -> x + y
```

```
let inc = add 1
```

```
let dec = add -1
```

```
let compose = fun f -> fun g -> fun x -> f (g x)
```

```
let id = compose inc dec
```

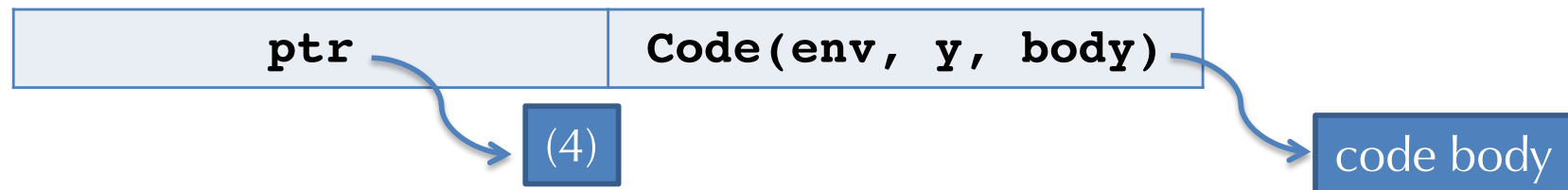
- How do we implement such functions?
  - In an interpreter? In a compiled language?

# Compiling First-class Functions

- To implement first-class functions on a processor, there are 2 problems
  - Must implement substitution of free variables
  - Must separate “code” from “data”
- Reify the substitution
  - Move substitution from the meta language to the object language by making the data structure & lookup operation explicit
  - The environment-based interpreter is one step in this direction
- Closure conversion
  - Eliminates free variables by packaging up the needed environment in the data structure
- Hoisting
  - Separates code from data, pulling closed code to the top level

# Example of Closure Creation

- Recall the “add” function  
`let add = fun x -> fun y -> x + y`
- Consider the inner function: `fun y -> x + y`
- When run the function application: `add 4`  
the program builds a closure and returns it
  - The closure is a pair of the environment and a code pointer



- The code pointer takes a pair of parameters: env and y
  - The function code is (essentially)  
`fun (env, y) -> let x = nth env 0 in x + y`

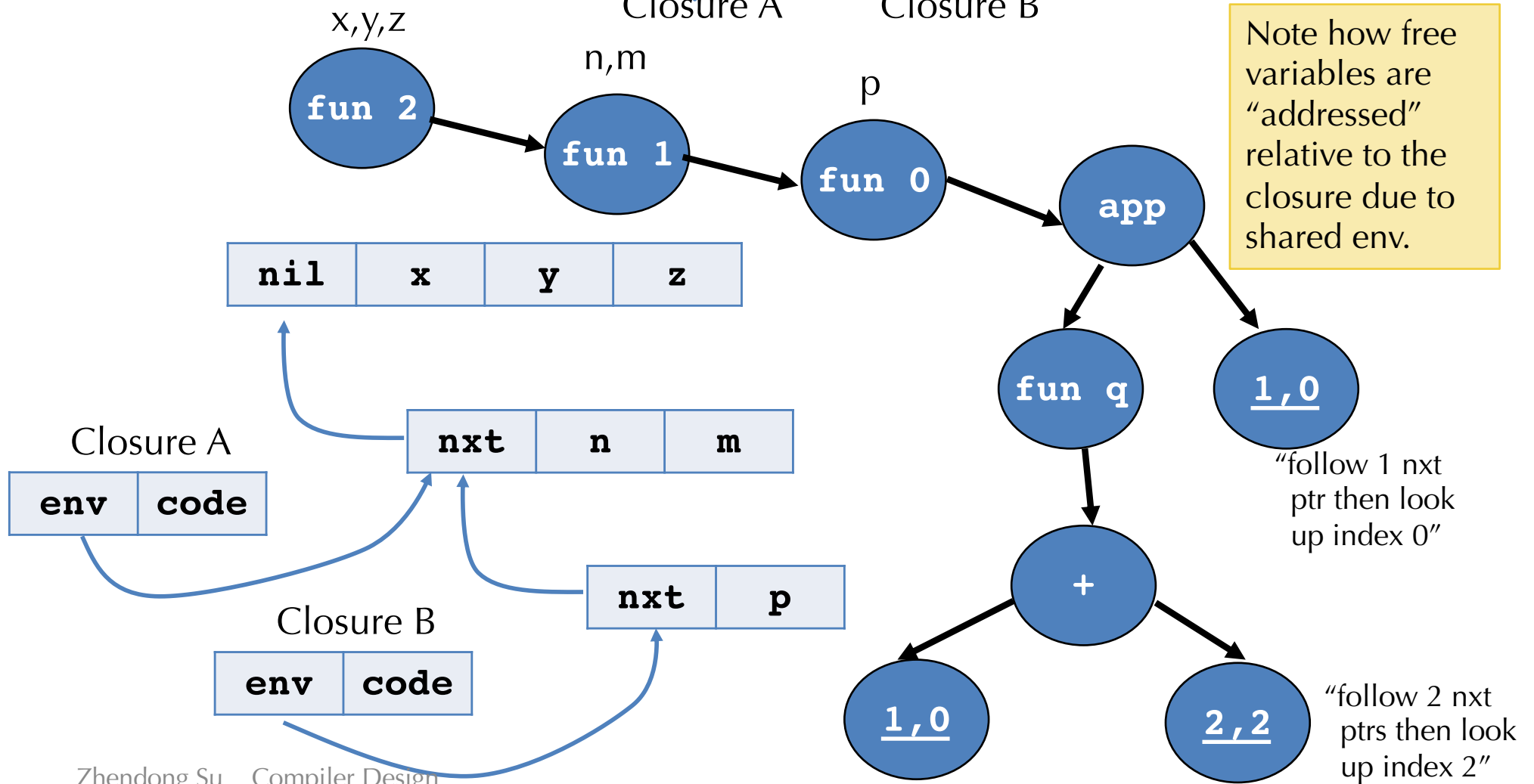
# Representing Closures

- The simple closure conversion doesn't generate very efficient code
  - It stores all the values for variables in the environment, even if they aren't needed by the function body
  - It copies the environment values each time a nested closure is created
  - It uses a linked-list data structure for tuples
- There are many options
  - Store only the values for free variables in the body of the closure
  - Share subcomponents of the environment to avoid copying
  - Use vectors or arrays rather than linked structures

# Array-based Closures with N-ary Functions

```
(fun (x y z) ->
  (fun (n m) -> (fun p -> (fun q -> n + z) x)
```

Closure A
Closure B







Scope, Types, and Context

# STATIC ANALYSIS

# Adding Integers to Lambda Calculus

exp ::=  
| ...  
| n *constant integers*  
| exp<sub>1</sub> + exp<sub>2</sub> *binary arithmetic operation*

val ::=  
| fun x -> exp *functions are values*  
| n *integers are values*

n{v/x} = n *constants have no free vars*  
(e<sub>1</sub> + e<sub>2</sub>){v/x} = (e<sub>1</sub>{v/x} + e<sub>2</sub>{v/x}) *substitute everywhere*

$$\frac{\text{exp}_1 \Downarrow n_1 \quad \text{exp}_2 \Downarrow n_2}{\text{exp}_1 + \text{exp}_2 \Downarrow (n_1 \llbracket + \rrbracket n_2)}$$

object-level '+'      meta-level '+'

**NOTE:** there are no rules for the case where exp<sub>1</sub> or exp<sub>2</sub> evaluate to functions! The semantics is *undefined* in those cases.

# Variable Scoping

- Problem: How to determine whether a declared variable is in scope?
- Issues
  - Q1: Which variables are available at a given program location?
  - Q2: Can the same identifier be reused (i.e., shadowing), or it is an error?
- Example: code below is syntactically correct, but not well-formed
  - `y` and `q` are used without being defined anywhere

```
int fact(int x) {  
    var acc = 1;  
    while (x > 0) {  
        acc = acc * y;  
        x = q - 1;  
    }  
    return acc;  
}
```

Q: Can we solve this problem by changing the parser to rule out such programs?

# Type Checking / Static Analysis

The interpreter from the `Eval3` module (`fun.ml`, Lec 13)

```
let rec eval env e =
  match e with
  | ...
  | Add (e1, e2) ->
    (match (eval env e1, eval env e2) with
     | (IntV i1, IntV i2) -> IntV (i1 + i2)
     | _ -> failwith "tried to add non-integers")
  | ...
```

- The interpreter might fail at runtime
  - Not all operations are defined for all values (e.g. `3/0`, `3 + true`, ...)
- A compiler can't generate sensible code for this case
  - A naïve implementation might "add" an integer and a function pointer



See `tc.ml` (`lec13.zip`)

# STATICALLY RULING OUT PARTIALITY: TYPE CHECKING

# Notes about this Typechecker

- The interpreter only evaluates the body of a function when it's applied
- Typechecker always check function's body (even if it's never applied)
  - Assume the input has some type (say  $t_1$ )
  - Reflect this information in the type of the function ( $t_1 \rightarrow t_2$ )
- Dually, at a call site ( $e_1 e_2$ ), we don't know what *closure* we'll get, but
  - Can calculate  $e_1$ 's type
  - Check  $e_2$  is an argument of the right type
  - Determine what type  $e_1$  will return

(Q1) Why is this an approximation?

(Q2) What if `well_typed` always returns `false`?

# Contexts and Inference Rules

- Need to keep track of contextual information
  - What variables are in scope?
  - What are their types?
  - What information do we have about each syntactic construct?
- What relationships are there among the syntactic objects?
  - e.g. is one type a subtype of another?
- How do we describe this information?
  - In the compiler, there's a mapping from variables to information we know about them – the "context"
  - The compiler has a collection of (mutually recursive) functions that follow the structure of the syntax

# Type Judgments

- In the judgment:  $E \vdash e : t$ 
  - $E$  is a *typing environment* or a *type context*
  - $E$  maps variables to types and is simply a set of bindings of the form:  
 $x_1 : t_1, x_2 : t_2, \dots, x_n : t_n$
- For example:  $x : \text{int}, b : \text{bool} \vdash \text{if } (b) \ 3 \ \text{else } x : \text{int}$
- What do we need to know to decide whether “if (b) 3 else x” has type int in the environment  $x : \text{int}, b : \text{bool}$ ?
  - $b$  must be a bool                    i.e.         $x : \text{int}, b : \text{bool} \vdash b : \text{bool}$
  - $3$  must be an int                    i.e.         $x : \text{int}, b : \text{bool} \vdash 3 : \text{int}$
  - $x$  must be an int                    i.e.         $x : \text{int}, b : \text{bool} \vdash x : \text{int}$



# Simply-typed Lambda Calculus

For the language in “tc.ml”, we have five inference rules

INT

$$\frac{}{E \vdash i : \text{int}}$$

VAR

$$\frac{x : T \in E}{E \vdash x : T}$$

ADD

$$\frac{E \vdash e_1 : \text{int} \quad E \vdash e_2 : \text{int}}{E \vdash e_1 + e_2 : \text{int}}$$

FUN

$$\frac{E, x : T \vdash e : S}{E \vdash \text{fun } (x:T) \rightarrow e : T \rightarrow S}$$

APP

$$\frac{E \vdash e_1 : T \rightarrow S \quad E \vdash e_2 : T}{E \vdash e_1 e_2 : S}$$

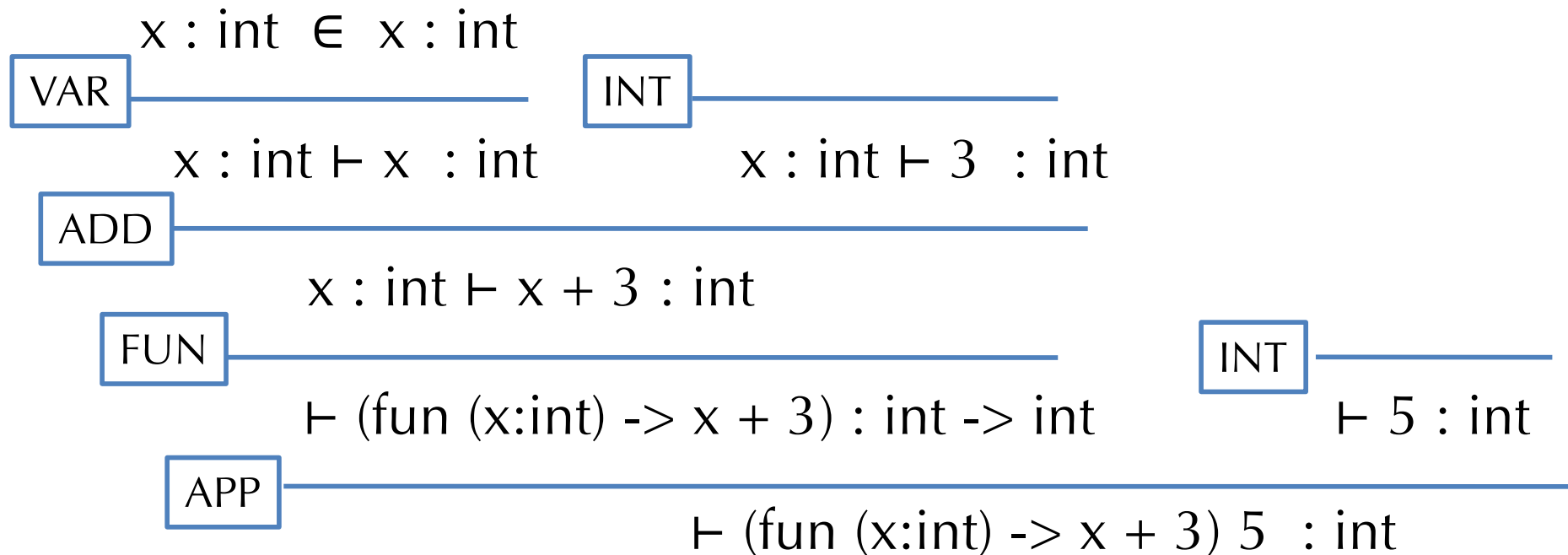
Note how these rules correspond to the code

# Type Checking Derivations

- *Derivation* or *proof tree*
  - Nodes: judgments
  - Edges: connect premises to a conclusion (according to inference rules)
- Leaves of the tree are *axioms* (i.e., rules with no premises)
  - Example: the INT rule is an axiom
- **Goal of the type checker:** verify that such a tree exists
  
- Ex: Find a tree for the following code using the given inference rules

```
⊢ (fun (x:int) -> x + 3) 5 : int
```

# Example Derivation Tree



Note

- The OCaml function `typecheck` verifies the existence of this tree
- Recursive calls for running `typecheck` follow the same shape as this tree
- `x : int ∈ E` is implemented by the function `lookup`

# Type Safety

*"Well typed programs do not go wrong."*

– Robin Milner, 1978

**Theorem:** (simply typed lambda calculus with integers)

If  $\vdash e : t$ , then there exists a value  $v$  such that  $e \Downarrow v$

- Note: This is a *very* strong property
  - Well-typed programs never executes undefined code like `3 + (fun x -> 2)`
  - Simply-typed lambda calculus terminates (i.e., not Turing complete)

# Type Safety For General Languages

## Theorem: (Type Safety)

If  $\vdash P : t$  is a well-typed program, then either

- (a) the program terminates in a well-defined way, or
- (b) the program continues computing forever

- Well-defined termination could include
  - halting with a return value
  - raising an exception
- Type safety rules out undefined behavior
  - abusing "unsafe" casts: converting pointers to integers, etc.
  - treating non-code values as code (and vice-versa)
  - breaking the type abstractions of the language
- What is "defined" depends on the language semantics

# Arrays

- Array constructs are not hard
- First: add a new type constructor:  $T[]$

$$\boxed{\text{NEW}} \quad \frac{E \vdash e_1 : \text{int} \quad E \vdash e_2 : T}{E \vdash \text{new } T[e_1](e_2) : T[]}$$

$e_1$ : size of newly alloc. array  
 $e_2$ : initializes the array

$$\boxed{\text{INDEX}} \quad \frac{E \vdash e_1 : T[] \quad E \vdash e_2 : \text{int}}{E \vdash e_1[e_2] : T}$$

$$\boxed{\text{UPDATE}} \quad \frac{E \vdash e_1 : T[] \quad E \vdash e_2 : \text{int} \quad E \vdash e_3 : T}{E \vdash e_1[e_2] = e_3 \text{ ok}}$$

Note: These rules don't ensure array indices are within bounds, which should be checked *dynamically*

# Tuples

- ML-style tuples with statically known number of products
- First: add a new type constructor:  $T_1 * \dots * T_n$

TUPLE

$$E \vdash e_1 : T_1 \quad \dots \quad E \vdash e_n : T_n$$

---

$$E \vdash (e_1, \dots, e_n) : T_1 * \dots * T_n$$

PROJ

$$E \vdash e : T_1 * \dots * T_n \quad 1 \leq i \leq n$$

---

$$E \vdash \#i e : T_i$$

# References

- ML-style references (note that ML uses only expressions)
- First, add a new type constructor: **T ref**

REF

$$\frac{E \vdash e : T}{E \vdash \text{ref } e : T \text{ ref}}$$

DEREF

$$\frac{E \vdash e : T \text{ ref}}{E \vdash !e : T}$$

ASSIGN

$$\frac{E \vdash e_1 : T \text{ ref} \quad E \vdash e_2 : T}{E \vdash e_1 := e_2 : \text{unit}}$$

Note the similarity with the rules for arrays





Beyond describing “structure”... describing “properties”

Types as sets

Subsumption

# TYPES, MORE GENERALLY

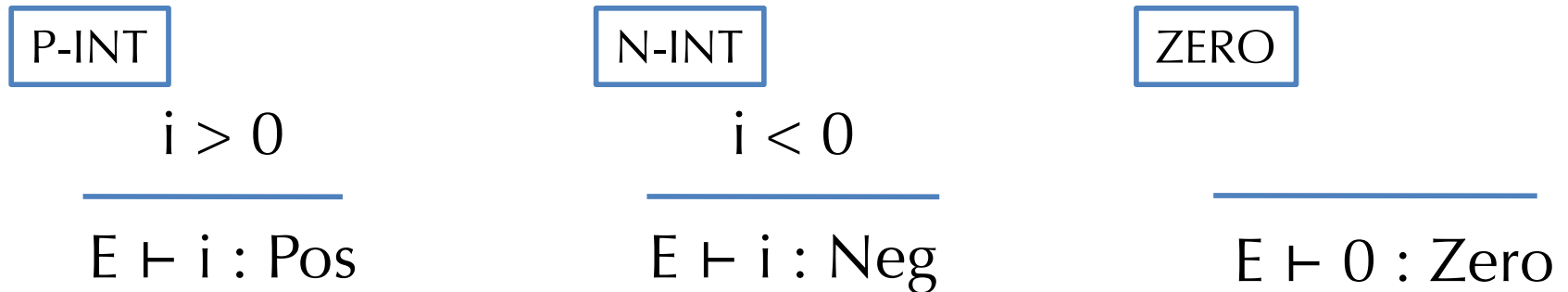
# What are types, anyway?

- A *type* is just a predicate on the set of values in a system
  - E.g., the type “int” can be thought of as a boolean function that returns “true” on integers and “false” otherwise
  - Equivalently, we can think of a type as just a *subset* of all values
- For efficiency and tractability, the predicates are usually very simple
  - Types are an *abstraction* mechanism
- We can easily add new types that distinguish different subsets of values

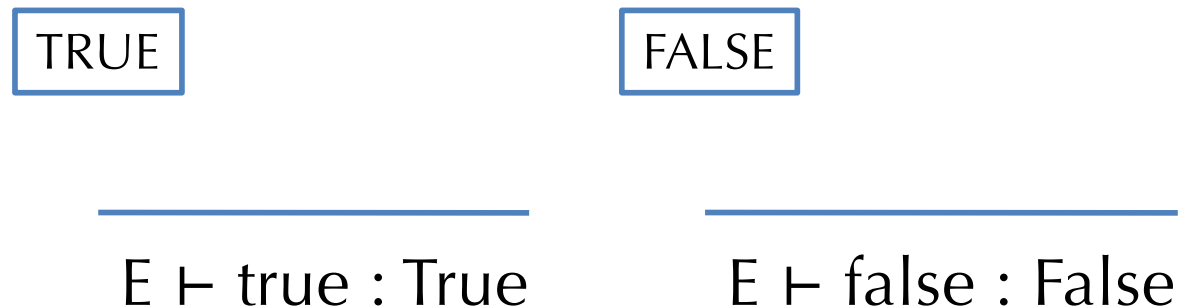
```
type tp =  
  | IntT          (* type of integers *)  
  | PostT | NegT | ZeroT (* refinements of ints *)  
  | BoolT        (* type of booleans *)  
  | TrueT | FalseT (* subsets of booleans *)  
  | AnyT         (* any value *)
```

# Modifying the typing rules

- We need to refine the typing rules too
- Some easy cases
  - Just split up the integers into their more refined cases



- Same for booleans



# What about “if”?

- Two cases are easy

$$\begin{array}{c} \boxed{\text{IF-T}} \quad E \vdash e_1 : \text{True} \quad E \vdash e_2 : T \\ \hline E \vdash \text{if } (e_1) e_2 \text{ else } e_3 : T \end{array} \qquad \begin{array}{c} \boxed{\text{IF-F}} \quad E \vdash e_1 : \text{False} \quad E \vdash e_3 : T \\ \hline E \vdash \text{if } (e_1) e_2 \text{ else } e_3 : T \end{array}$$

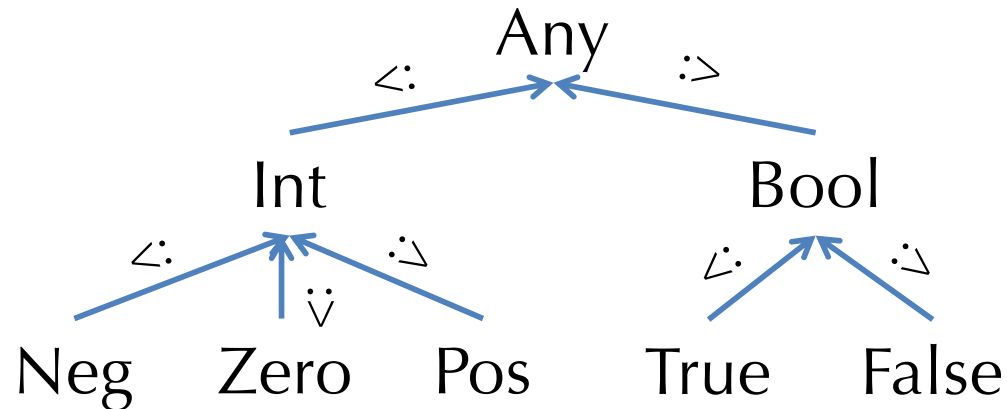
- What if we don't know statically which branch will be taken?
- Consider the typechecking problem

$x:\text{bool} \vdash \text{if } (x) 3 \text{ else } -1 : ?$

- The true branch has type Pos, while the false branch has type Neg
  - What should be the result type of the whole if?

# Subtyping and Upper Bounds

- If we view types as sets of values, there is a natural inclusion relation  
 $\text{Pos} \subseteq \text{Int}$
- This subset relation gives rise to a *subtype* relation:  $\text{Pos} <: \text{Int}$
- Such inclusions give rise to a *subtyping hierarchy*



- For types  $T_1, T_2$ , define their *least upper bound* (LUB) wrt the hierarchy
  - Example:  $\text{LUB}(\text{True}, \text{False}) = \text{Bool}$ ,  $\text{LUB}(\text{Int}, \text{Bool}) = \text{Any}$
  - Note: may want to add types for “NonZero”, “NonNegative”, “NonPositive”, so that set union on values corresponds to taking LUBs on types

# “If” Typing Rule Revisited

- For statically unknown conditionals, we want the return value to be the LUB of the types of the branches

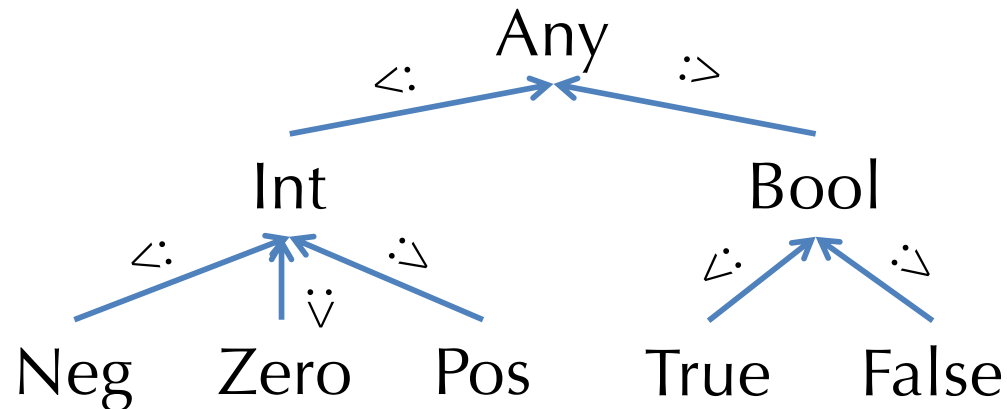
IF-BOOL

$$\frac{E \vdash e_1 : \text{bool} \quad E \vdash e_2 : T_1 \quad E \vdash e_3 : T_2}{E \vdash \text{if } (e_1) e_2 \text{ else } e_3 : \text{LUB}(T_1, T_2)}$$

- Note  $\text{LUB}(T_1, T_2)$  is the most precise type (wrt the hierarchy) that is able to describe any value that has either type  $T_1$  or type  $T_2$
- In math notation,  $\text{LUB}(T_1, T_2)$  is sometimes written  $T_1 \vee T_2$
- LUB is also called the *join* operation

# Subtyping Hierarchy

- A *subtyping hierarchy*



- The subtyping relation is a *partial order*
  - **Reflexive**:  $T \langle: T$  for any type  $T$
  - **Transitive**:  $T_1 \langle: T_2$  and  $T_2 \langle: T_3$  then  $T_1 \langle: T_3$
  - **Antisymmetric**: If  $T_1 \langle: T_2$  and  $T_2 \langle: T_1$  then  $T_1 = T_2$

# Soundness of Subtyping Relations

- We don't have to treat every subset of the integers as a type
  - e.g., we left out the type NonNeg
- A subtyping relation  $T_1 <: T_2$  is *sound* if it approximates the underlying semantic subset relation
- Formally: write  $\llbracket T \rrbracket$  for the subset of (closed) values of type  $T$ 
  - i.e.  $\llbracket T \rrbracket = \{v \mid \vdash v : T\}$
  - e.g.  $\llbracket \text{Zero} \rrbracket = \{0\}$ ,  $\llbracket \text{Pos} \rrbracket = \{1, 2, 3, \dots\}$
- If  $T_1 <: T_2$  implies  $\llbracket T_1 \rrbracket \subseteq \llbracket T_2 \rrbracket$ , then  $T_1 <: T_2$  is sound
  - Ex 1:  $\text{Pos} <: \text{Int}$  is sound, since  $\{1, 2, 3, \dots\} \subseteq \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$
  - Ex 2:  $\text{Int} <: \text{Pos}$  is not sound, since it is *not* the case that  $\{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\} \subseteq \{1, 2, 3, \dots\}$



# Soundness of LUBs

- Whenever we have a sound subtyping relation, it follows that
$$\llbracket \text{LUB}(T_1, T_2) \rrbracket \supseteq \llbracket T_1 \rrbracket \cup \llbracket T_2 \rrbracket$$
  - Note that the LUB is an over approximation of the “semantic union”
  - Example:  $\llbracket \text{LUB}(\text{Zero}, \text{Pos}) \rrbracket = \llbracket \text{Int} \rrbracket = \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\} \supseteq \{0, 1, 2, 3, \dots\} = \{0\} \cup \{1, 2, 3, \dots\} = \llbracket \text{Zero} \rrbracket \cup \llbracket \text{Pos} \rrbracket$
- Using LUBs in the typing rules yields sound approximations of the program behavior (as if the IF-B rule)
- It just so happens that LUBs on types  $<: \text{Int}$  correspond to +

ADD

$$E \vdash e_1 : T_1 \quad E \vdash e_2 : T_2 \quad T_1 <: \text{Int} \quad T_2 <: \text{Int}$$

---

$$E \vdash e_1 + e_2 : T_1 \vee T_2$$

# Subsumption Rule

- When we add subtyping judgments of the form  $T <: S$  we can uniformly integrate it into the type system generically

SUBSUMPTION

$$E \vdash e : T \quad T <: S$$

---

$$E \vdash e : S$$

- Subsumption allows any value of type  $T$  to be treated as an  $S$  whenever  $T <: S$
- Adding this rule makes the search for typing derivations more difficult
  - this rule can be applied anywhere, since  $T <: T$
  - But careful engineering of the typing system can incorporate the subsumption rule into a deterministic algorithm
  - See for example the OAT type system

# Downcasting

- What happens if we have an Int, but need something of type Pos?
  - At compile time, we don't know whether the Int is greater than zero
  - At run time, we do
- Add a “checked downcast”

$$\frac{E \vdash e_1 : \text{Int} \quad E, x : \text{Pos} \vdash e_2 : T_2 \quad E \vdash e_3 : T_3}{E \vdash \text{ifPos } (x = e_1) e_2 \text{ else } e_3 : T_2 \vee T_3}$$

- At runtime, ifPos checks whether  $e_1$  is  $> 0$ . If so, branches to  $e_2$  and otherwise branches to  $e_3$
- Inside the expression  $e_2$ ,  $x$  is the name for  $e_1$ 's value, which is known to be strictly positive because of the dynamic check
- Note such rules force the programmer to add the appropriate checks
  - We could give integer division the type:  $\text{Int} \rightarrow \text{NonZero} \rightarrow \text{Int}$