

Lecture 18

# EECS 483: COMPILER CONSTRUCTION

# Announcements

- HW4: OAT v. 1.0
  - Due tomorrow evening.
- HW5: OAT v. 2.0
  - records, function pointers, type checking, array-bounds checks, etc.
  - Due: Tuesday, April 9<sup>th</sup>
  - Available late tomorrow, review during Wednesday's lecture.
  - **As usual, Start Early!**
- Next two weeks
  - Professor New will be on family starting soon (at the latest next Monday)
    - Lectures by Eric or guest professors
    - Only remote office hours for Prof. New, no appointments
  - April 8: Total solar eclipse (only one this century in the USA!)
    - During lecture time, so we will cancel.

# Type Judgments

- In the judgment:  $E \vdash e : t$ 
  - $E$  is a *typing environment* or a *type context*
  - $E$  maps variables to types. It is just 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

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

ADD

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

FUN

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

APP

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

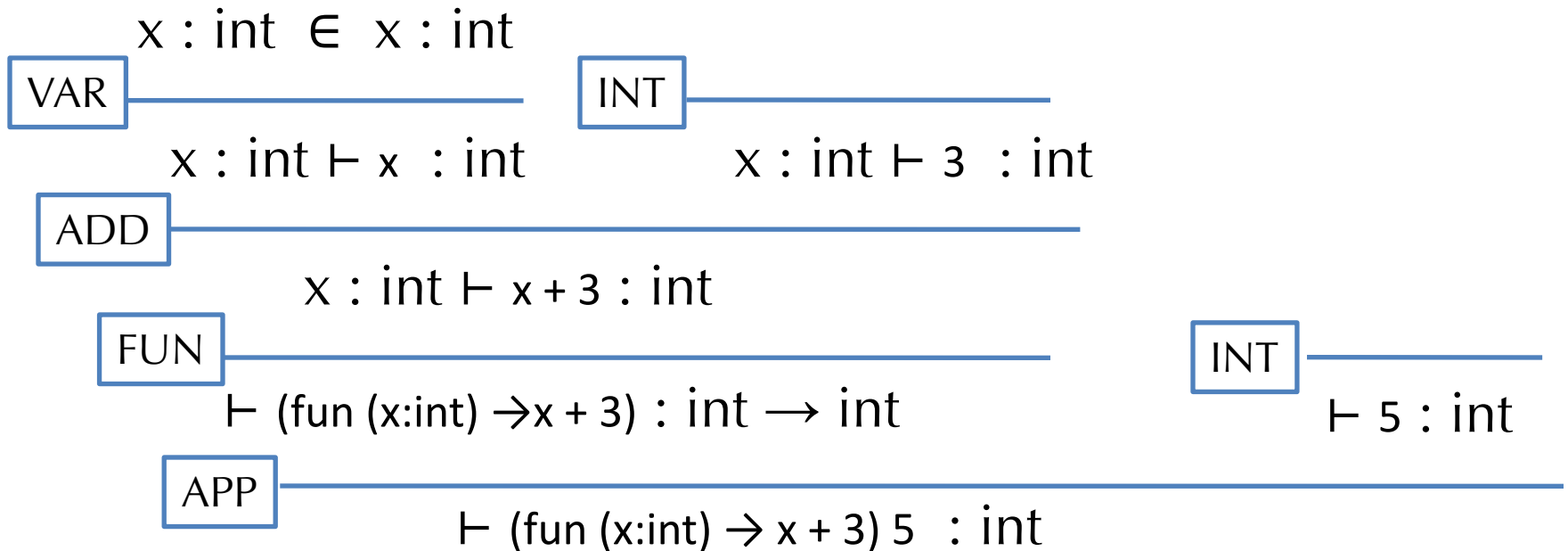
- Note how these rules correspond to the code.

# Type Checking Derivations

- A *derivation* or *proof tree* has (instances of) judgments as its nodes and edges that connect premises to a conclusion according to an inference rule.
- Leaves of the tree are *axioms* (i.e. rules with no premises)
  - Example: the INT rule is an axiom
- Goal of the typechecker: verify that such a tree exists.
- Example: Find a tree for the following program using the inference rules on the previous slide:

$\vdash (\text{fun } (x:\text{int}) \rightarrow x + 3) \ 5 \ : \text{int}$

# Example Derivation Tree



- Note: the OCaml function `typecheck` verifies the existence of this tree. The structure of the recursive calls when running `typecheck` is the same shape as this tree!
- Note that  $x : \text{int} \in E$  is implemented by the function `lookup`

# Ill-typed Programs

- Programs without derivations are ill-typed

Example: There is no type  $T$  such that

$$\vdash (\text{fun } (x:\text{int}) \rightarrow x \ 3) \ 5 : T$$

$x : \text{int} \rightarrow T \notin x : \text{int}$

VAR

$x : \text{int} \vdash x : \text{int} \rightarrow T$

$x : \text{int} \vdash 3 : \text{int}$

APP

$x : \text{int} \vdash x \ 3 : T$

FUN

$\vdash (\text{fun } (x:\text{int}) \rightarrow x \ 3) : \text{int} \rightarrow T$

$\vdash 5 : \text{int}$

APP

$\vdash (\text{fun } (x:\text{int}) \rightarrow x \ 3) \ 5 : T$

# 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 cannot "go wrong" by trying to execute undefined code (such as `3 + (fun x -> 2)`)
  - Simply-typed lambda calculus is guaranteed to terminate! (i.e. it isn't Turing complete)



# Notes about this Typechecker

- The interpreter evaluates the body of a function only when it's applied.
- The typechecker always checks the body of the function
  - even if it's never applied
  - We *assume* the input has some type (say  $t_1$ ) and reflect this 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're going to get.
  - But we can calculate  $e_1$ 's type, check that  $e_2$  is an argument of the right type, and determine what type  $e_1$  will return.
- Question: Why is this an approximation?
- Question: What if `well_typed` always returns false?



oat.pdf

# TYPECHECKING OAT

# Checking Derivations

- A *derivation* or *proof tree* has (instances of) judgments as its nodes and edges that connect premises to a conclusion according to an inference rule.
- 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.
- Example1: Find a tree for the following program using the inference rules in oat.pdf:

```
var x1 = 0;  
var x2 = x1 + x1;  
x1 = x1 - x2;  
return(x1);
```

Example2: There is no tree for this ill-scoped program:

```
var x2 = x1 + x1;  
return(x2);
```

# Example Derivation

```
var x1 = 0;  
var x2 = x1 + x1;  
x1 = x1 - x2;  
return(x1);
```

$$\frac{\frac{\mathcal{D}_1 \quad \mathcal{D}_2 \quad \mathcal{D}_3 \quad \mathcal{D}_4}{G_0; \cdot; \text{int} \vdash \text{var } x_1 = 0; \text{var } x_2 = x_1 + x_1; x_1 = x_1 - x_2; \text{return } x_1; \Rightarrow \cdot, x_1:\text{int}, x_2:\text{int}}}{\vdash \text{var } x_1 = 0; \text{var } x_2 = x_1 + x_1; x_1 = x_1 - x_2; \text{return } x_1;} \begin{array}{l} [\text{STMTS}] \\ [\text{PROG}] \end{array}$$

# Example Derivation

$$\mathcal{D}_1 = \frac{\frac{\frac{}{G_0; \cdot \vdash 0 : \text{int}} [\text{INT}]}{G_0; \cdot \vdash 0 : \text{int}} [\text{CONST}]}{G_0; \cdot \vdash \text{var } x_1 = 0 \Rightarrow \cdot, x_1 : \text{int}} [\text{DECL}]}{G_0; \cdot; \text{int} \vdash \text{var } x_1 = 0; \Rightarrow \cdot, x_1 : \text{int}} [\text{SDECL}]$$

$$\mathcal{D}_2 = \frac{\frac{\frac{}{\vdash + : (\text{int}, \text{int}) \rightarrow \text{int}} [\text{ADD}]}{G_0; \cdot, x_1 : \text{int} \vdash x_1 : \text{int}} [\text{VAR}]}{G_0; \cdot, x_1 : \text{int} \vdash x_1 + x_1 : \text{int}} [\text{BOP}]}{\frac{\frac{}{G_0; \cdot, x_1 : \text{int}; \text{int} \vdash \text{var } x_2 = x_1 + x_1; \Rightarrow \cdot, x_1 : \text{int}, x_2 : \text{int}} [\text{DECL}]}{G_0; \cdot, x_1 : \text{int}; \text{int} \vdash \text{var } x_2 = x_1 + x_1; \Rightarrow \cdot, x_1 : \text{int}, x_2 : \text{int}} [\text{SDECL}]}$$

# Example Derivation

$$\mathcal{D}_3 = \frac{\frac{\frac{}{\vdash - : (\text{int}, \text{int}) \rightarrow \text{int}} \text{[ADD]} \quad \frac{x_1 : \text{int} \in \cdot, x_1 : \text{int}, x_2 : \text{int}}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int} \vdash x_1 : \text{int}} \text{[VAR]} \quad \frac{x_2 : \text{int} \in \cdot, x_1 : \text{int}, x_2 : \text{int}}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int} \vdash x_2 : \text{int}} \text{[VAR]}}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int} \vdash x_1 - x_2 : \text{int}} \text{[BOP]} \quad \frac{}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int}; \text{int} \vdash x_1 = x_1 - x_2; \Rightarrow \cdot, x_1 : \text{int}, x_2 : \text{int}} \text{[ASSN]}$$

$$\mathcal{D}_4 = \frac{\frac{x_1 : \text{int} \in \cdot, x_1 : \text{int}, x_2 : \text{int}}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int} \vdash x_1 : \text{int}} \text{[VAR]}}{G_0 ; \cdot, x_1 : \text{int}, x_2 : \text{int}; \text{int} \vdash \text{return } x_1; \Rightarrow \cdot, x_1 : \text{int}, x_2 : \text{int}} \text{[RET]}$$

# 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 behaviors:
  - 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...



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

Types as sets

Subsumption

## **TYPES, MORE GENERALLY**



# What are types, anyway?

- A *type* is can often be thought of as a predicate on the set of values in a system.
  - For example, 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 taken to be 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 *)
| PosT | 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:

P-INT

$i > 0$

$G \vdash i : \text{Pos}$

N-INT

$i < 0$

$G \vdash i : \text{Neg}$

ZERO

$G \vdash 0 : \text{Zero}$

- Same for booleans:

TRUE

$G \vdash \text{true} : \text{True}$

FALSE

$G \vdash \text{false} : \text{False}$

# What about “if”?

- Two cases are easy:

$$\frac{\boxed{\text{IF-T}} \quad G \vdash e_1 : \text{True} \quad G \vdash e_2 : T}{G \vdash \text{if } (e_1) \ e_2 \text{ else } e_3 : T} \qquad \frac{\boxed{\text{IF-F}} \quad G \vdash e_1 : \text{False} \quad E \vdash e_3 : T}{G \vdash \text{if } (e_1) \ e_2 \text{ else } e_3 : T}$$

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

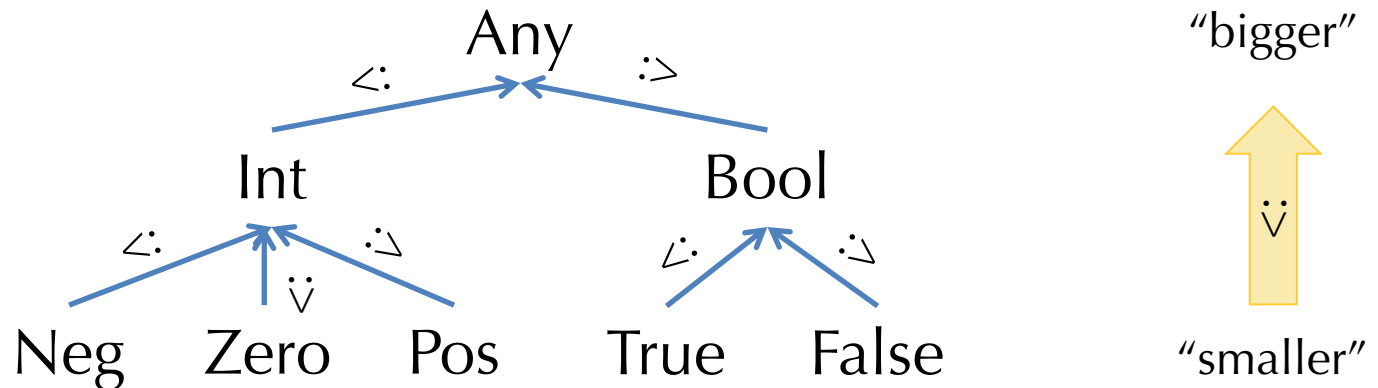
$x:\text{bool} \vdash \text{if } (x) \ 3 \text{ else } -1 : ?$

The true branch has type Pos and the false branch has type Neg.

- What should be the result type of the whole if?

# Subtyping and Upper Bounds

- If we think of types as sets of values, we have 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*:



- Given any two types  $T_1$  and  $T_2$ , we can calculate their *least upper bound* (LUB) according to the hierarchy.
  - Definition:  $\text{LUB}(T_1, T_2)$  is the smallest  $T$  such that  $T_1 <: T$  and  $T_2 <: T$
  - Example:  $\text{LUB}(\text{True}, \text{False}) = \text{Bool}$ ,  $\text{LUB}(\text{Int}, \text{Bool}) = \text{Any}$
- Note: might want to add types for “NonZero”, “NonNegative”, and “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

$$G \vdash e_1 : \text{bool} \quad E \vdash e_2 : T_1 \quad G \vdash e_3 : T_2$$

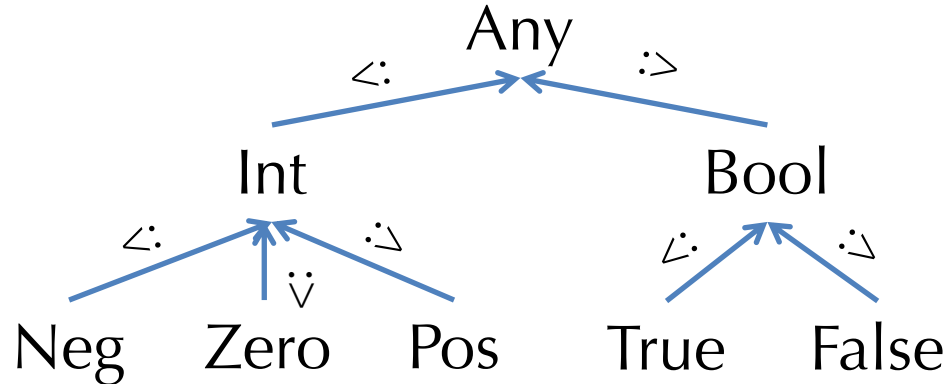
---

$$G \vdash \text{if } (e_1) e_2 \text{ else } e_3 : \text{LUB}(T_1, T_2)$$

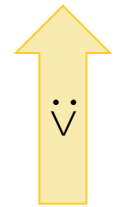
- Note:  $\text{LUB}(T_1, T_2)$  is the most precise type (according to the hierarchy) that can describe any value that has either type  $T_1$  or type  $T_2$ .
- In math notation,  $\text{LUB}(T_1, T_2)$  is sometimes written  $T_1 \vee T_2$
- LUB is also called the *join* operation.

# Subtyping Hierarchy

- A *subtyping hierarchy*:



“bigger”



“smaller”

- The subtyping relation is a *partial order*:
  - Reflexive:  $T <: T$  for any type  $T$
  - Transitive:  $T_1 <: T_2$  and  $T_2 <: T_3$  then  $T_1 <: T_3$
  - Antisymmetric: If  $T_1 <: T_2$  and  $T_2 <: T_1$  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.
  - e.g.,  $\text{Pos} <: \text{Int}$  is sound, since  $\{1, 2, 3, \dots\} \subseteq \{\dots, -3, -2, -1, 0, 1, 2, 3, \dots\}$
  - e.g.,  $\text{Int} <: \text{Pos}$  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 you 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 these specific subtypes of Int are *sound* for +

ADD

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

---

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