

Lecture 18

# **CIS 4521/5521: COMPILERS**

# Announcements

- HW5: OAT v. 2.0
  - records, function pointers, type checking, array-bounds checks, etc.
  - Due: Wednesday, April 9<sup>th</sup>
  - Available soon (by Saturday morning)
  - **Start Early!**



# TYPECHECKING


# Recap

- A typechecking (static analysis) specification can be defined by collections of inference rules.
  - Each "judgment" form corresponds to a particular kind of analysis
  - The rule's premises and conclusion specify the intended checks

$$\frac{G \vdash e_1 : T \rightarrow S \quad G \vdash e_2 : T}{G \vdash e_1 e_2 : S} \quad \text{APPLICATION}$$

- Subtyping introduces a notion of subsumption (inclusion):

$$\frac{G \vdash e : T \quad T <: S}{G \vdash e : S} \quad \text{SUBSUMPTION}$$

*subtyping relation* 



# SUBTYPING OTHER TYPES

# Extending Subtyping to Other Types

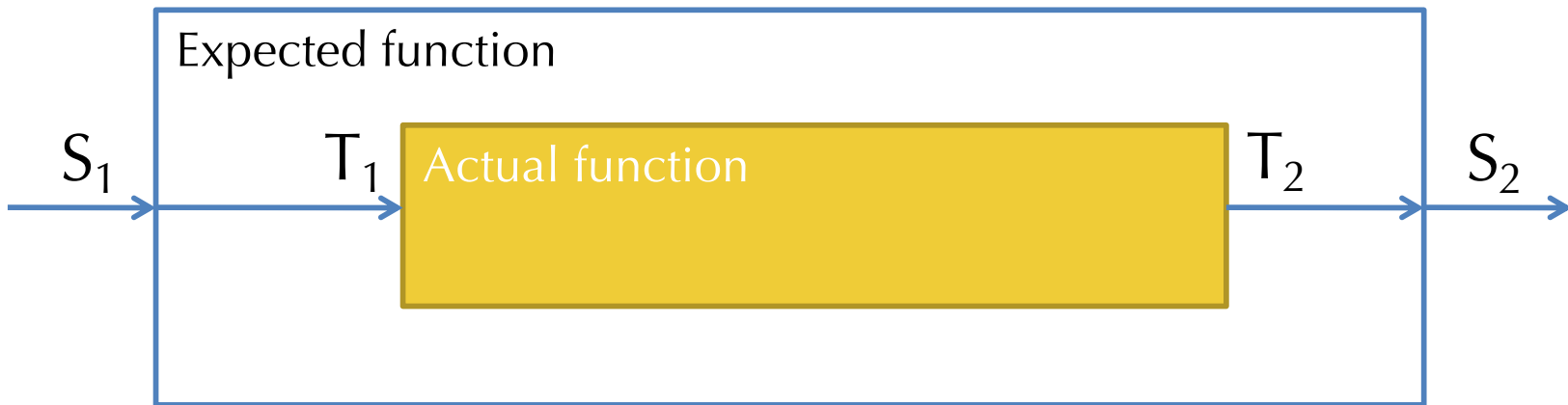
- What about subtyping for tuples?
  - Intuition: whenever a program expects something of type  $S_1 * S_2$ , it is sound to give it a  $T_1 * T_2$ .
  - Example:  $(\text{Pos} * \text{Neg}) <: (\text{Int} * \text{Int})$

$$\frac{T_1 <: S_1 \quad T_2 <: S_2}{(T_1 * T_2) <: (S_1 * S_2)}$$

- What about functions?
- When is  $T_1 \rightarrow T_2 <: S_1 \rightarrow S_2$  ?

# Subtyping for Function Types

- One way to see it:



- Need to convert an  $S_1$  to a  $T_1$  and  $T_2$  to  $S_2$ , so the argument type is *contravariant* and the output type is *covariant*.

$$\frac{S_1 <: T_1 \quad T_2 <: S_2}{(T_1 \rightarrow T_2) <: (S_1 \rightarrow S_2)}$$

# Immutable Records

- Record type:  $\{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_n:T_n\}$ 
  - Each  $\text{lab}_i$  is a label drawn from a set of identifiers.

RECORD

$$G \vdash e_1 : T_1 \quad G \vdash e_2 : T_2 \quad \dots \quad G \vdash e_n : T_n$$

---

$$G \vdash \{\text{lab}_1 = e_1; \text{lab}_2 = e_2; \dots ; \text{lab}_n = e_n\} : \{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_n:T_n\}$$

PROJECTION

$$G \vdash e : \{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_n:T_n\}$$

---

$$G \vdash e.\text{lab}_i : T_i$$

# Immutable Record Subtyping

- Depth subtyping:
  - Corresponding fields may be subtypes

DEPTH

$$T_1 <: U_1 \quad T_2 <: U_2 \quad \dots \quad T_n <: U_n$$

---

$$\{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_n:T_n\} <: \{\text{lab}_1:U_1; \text{lab}_2:U_2; \dots ; \text{lab}_n:U_n\}$$

- Width subtyping:
  - Subtype record may have *more* fields:

WIDTH

$$m \leq n$$

---

$$\{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_n:T_n\} <: \{\text{lab}_1:T_1; \text{lab}_2:T_2; \dots ; \text{lab}_m:T_m\}$$

# Depth & Width Subtyping vs. Layout

- Width subtyping (without depth) is compatible with "inlined" record representation as with C structs:

`{x:int; y:int; z:int} <: {x:int; y:int}`  
[Width Subtyping]



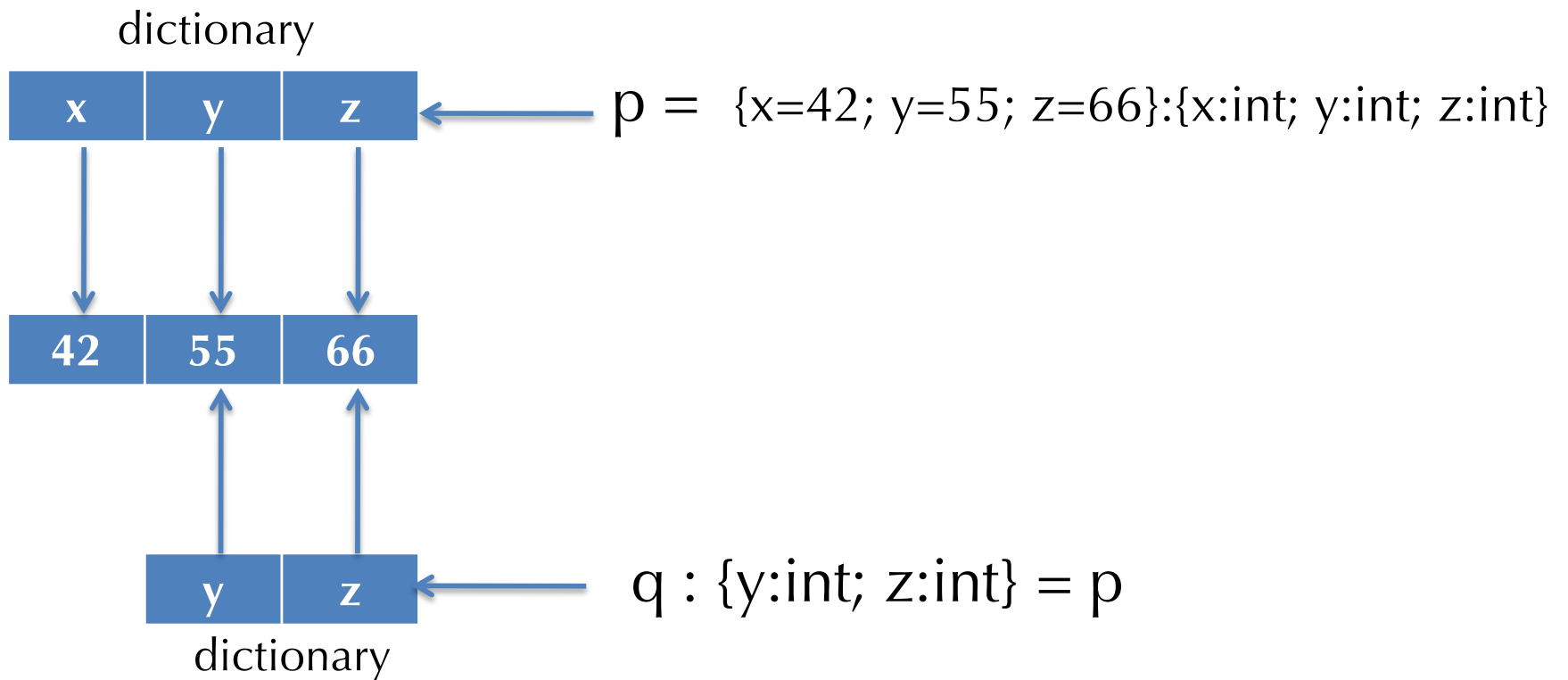
- The layout and underlying field indices for 'x' and 'y' are identical.
  - The 'z' field is just ignored
- Depth subtyping (without width) is similarly compatible, assuming that the space used by A is the same as the space used by B whenever  $A <: B$
- But... they don't mix without more work

# Immutable Record Subtyping (cont'd)

- Width subtyping assumes an implementation in which order of fields in a record matters:  
 $\{x:\text{int}; y:\text{int}\} \neq \{y:\text{int}; x:\text{int}\}$
- But:  $\{x:\text{int}; y:\text{int}; z:\text{int}\} <: \{x:\text{int}; y:\text{int}\}$ 
  - Implementation: a record is a struct, subtypes just add fields at the *end* of the struct.
- Alternative: allow permutation of record fields:  
 $\{x:\text{int}; y:\text{int}\} = \{y:\text{int}; x:\text{int}\}$ 
  - Implementation: compiler sorts the fields before code generation.
  - Need to know *all* of the fields to generate the code
- Permutation is not directly compatible with width subtyping:  
 $\{x:\text{int}; z:\text{int}; y:\text{int}\} = \{x:\text{int}; y:\text{int}; z:\text{int}\} <:/: \{y:\text{int}; z:\text{int}\}$

## If you want both:

- If you want permutability & dropping, you need to either copy (to rearrange the fields) or use a dictionary like this:





# MUTABILITY & SUBTYPING

# NULL

- What is the type of `null`?

- Consider:

```
int[] a = null;    // OK?  
int x   = null;    // not OK?  
string s = null;   // OK?
```

NULL

$G \vdash \text{null} : r$

- Null has any *reference type*
  - Null is generic
- What about type safety?
  - Requires defined behavior when dereferencing null  
e.g. Java's `NullPointerException`
  - Requires a safety check for every dereference operation  
(typically implemented using low-level hardware "trap" mechanisms.)

# Subtyping and References

- What is the proper subtyping relationship for references and arrays?
- Suppose we have NonZero as a type and the division operation has type:  $\text{Int} \rightarrow \text{NonZero} \rightarrow \text{Int}$ 
  - Recall that  $\text{NonZero} <: \text{Int}$
- Should  $(\text{NonZero ref}) <: (\text{Int ref})$  ?
- Consider this program:

```
Int bad(NonZero ref r) {  
  Int ref a = r;    (* OK because (NonZero ref <: Int ref*)  
  a := 0;           (* OK because 0 : Zero <: Int *)  
  return (42 / !r) (* OK because !r has type NonZero *)  
}
```

# Mutable Structures are Invariant

- Covariant reference types are unsound
  - As demonstrated in the previous example
- Contravariant reference types are also unsound
  - i.e. If  $T_1 <: T_2$  then  $\text{ref } T_2 <: \text{ref } T_1$  is also unsound
  - Exercise: construct a program that breaks contravariant references.
- Moral: Mutable structures are invariant:  
$$T_1 \text{ ref } <: T_2 \text{ ref} \quad \text{implies} \quad T_1 = T_2$$
- Same holds for arrays, OCaml-style mutable records, object fields, etc.
  - Note: Java and C# get this wrong. They allow covariant array subtyping, but then compensate by adding a dynamic check on every array update!

# Another Way to See It

- We can think of a reference cell as an immutable record (object) with two functions (methods) and some hidden state:
$$T \text{ ref} \simeq \{ \text{get}: \text{unit} \rightarrow T; \text{ set}: T \rightarrow \text{unit} \}$$
  - get returns the value hidden in the state.
  - set updates the value hidden in the state.
- When is  $T \text{ ref} <: S \text{ ref}$ ?
- Records are like tuples: subtyping extends pointwise over each component.
- $\{ \text{get}: \text{unit} \rightarrow T; \text{ set}: T \rightarrow \text{unit} \} <: \{ \text{get}: \text{unit} \rightarrow S; \text{ set}: S \rightarrow \text{unit} \}$ 
  - get components are subtypes:  $\text{unit} \rightarrow T <: \text{unit} \rightarrow S$
  - set components are subtypes:  $T \rightarrow \text{unit} <: S \rightarrow \text{unit}$
- From get, we must have  $T <: S$  (covariant return)
- From set, we must have  $S <: T$  (contravariant arg.)
- From  $T <: S$  and  $S <: T$  we conclude  $T = S$ .



# STRUCTURAL VS. NOMINAL TYPES

# Structural vs. Nominal Typing

- Is type equality / subsumption defined by the *structure* of the data or the *name* of the data?
- Example 1: type abbreviations (OCaml) vs. “newtypes” (a la Haskell)

```
(* OCaml: *)
type cents = int      (* cents = int in this scope *)
type age = int

let foo (x:cents) (y:age) = x + y
```

```
(* Haskell: *)
newtype Cents = Cents Integer  (* Integer and Cents are
                                isomorphic, not identical. *)
newtype Age = Age Integer

foo :: Cents -> Age -> Int
foo x y = x + y                (* Ill typed! *)
```

- Type abbreviations are treated “structurally”  
Newtypes are treated “by name”

# Nominal Subtyping in Java

- In Java, Classes and Interfaces must be named and their relationships *explicitly* declared:

```
(* Java: *)  
interface Foo {  
    int foo();  
}  
  
class C {          /* Does not implement the Foo interface */  
    int foo() {return 2;}  
}  
  
class D implements Foo {  
    int foo() {return 4521/5521;}  
}
```

- Similarly for inheritance: programmers must declare the subclass relation via the “**extends**” keyword.
  - Typechecker still checks that the classes are structurally compatible



See oat.pdf in HW5

# OAT'S TYPE SYSTEM

# OAT's Treatment of Types

- Primitive (non-reference) types:
  - `int`, `bool`
- Definitely non-null reference types: `R`
  - (named) mutable structs with (right-oriented) *width* subtyping
  - `string`
  - arrays (including length information, per HW4)
- Possibly-null reference types: `R?`
  - Subtyping: `R <: R?`
  - *Checked downcast* syntax `if?`:

```
int sum(int[]? arr) {  
    var z = 0;  
    if?(int[] a = arr) {  
        for(var i = 0; i < length(a); i = i + 1;) {  
            z = z + a[i];  
        }  
    }  
    return z;  
}
```

# OAT Features

- Named structure types with mutable fields
  - but using structural, width subtyping
- Typed function pointers
- Polymorphic operations: `length` and `== / !=`
  - need special case handling in the typechecker
- Type-annotated null values: `t null` always has type `t`?
- Definitely-not-null values means we need an "atomic" array initialization syntax
  - for example, `null` is not allowed as a value of type `int[]`, so to construct a record containing a field of type `int[]`, we need to initialize it
  - subtlety: `int[][]` cannot be initialized by default, but `int[]` can be

# 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;
```

# OAT "Returns" Analysis

- Typesafe, statement-oriented imperative languages like OAT (or Java) must ensure that a function (always) returns a value of the appropriate type.
  - Does the returned expression's type match the one declared by the function?
  - Do all paths through the code return appropriately?
- OAT's statement checking judgment
  - takes the expected return type as input: what type should the statement return (or `void` if none)
  - produces a boolean flag as output: does the statement definitely return?



# COMPILING WITH TYPES

# Compilation As Translating Judgments

- Consider the source typing judgment for source expressions:

$$C \vdash e : t$$

- How do we interpret this information in the target language?

$$\llbracket C \vdash e : t \rrbracket = ?$$

- $\llbracket C \rrbracket$  translates contexts
- $\llbracket t \rrbracket$  is a target type
- $\llbracket e \rrbracket$  translates to a (potentially empty) stream of instructions, that, when run, computes the result into some operand
- INVARIANT: if  $\llbracket C \vdash e : t \rrbracket = \text{ty}, \text{operand}, \text{stream}$   
then the type (at the target level) of the operand is  $\text{ty} = \llbracket t \rrbracket$

# Example

- $C \vdash 4521 + 5 : \text{int}$       what is  $\llbracket C \vdash 4521 + 5 : \text{int} \rrbracket$  ?

$\llbracket \vdash 4521 : \text{int} \rrbracket = (\text{i64}, \text{Const } 4521, [])$        $\llbracket \vdash 5 : \text{int} \rrbracket = (\text{i64}, \text{Const } 5, [])$

-----  
 $\llbracket C \vdash 4521 : \text{int} \rrbracket = (\text{i64}, \text{Const } 4521, [])$        $\llbracket C \vdash 5 : \text{int} \rrbracket = (\text{i64}, \text{Const } 5, [])$

-----  
 $\llbracket C \vdash 4521 + 5 : \text{int} \rrbracket = (\text{i64}, \%tmp, [\%tmp = \text{add i64 (Const } 4521) (\text{Const } 5)])$

# What about the Context?

- What is  $\llbracket C \rrbracket$ ?
- Source level  $C$  has bindings like:  $x:\text{int}, y:\text{bool}$ 
  - We think of it as a finite map from identifiers to types
- What is the interpretation of  $C$  at the target level?
- $\llbracket C \rrbracket$  maps source identifiers, “ $x$ ” to target types and  $\llbracket x \rrbracket$
- What is the interpretation of a variable  $\llbracket x \rrbracket$  at the target level?
  - How are the variables used in the type system?

# Interpretation of Contexts

- $\llbracket C \rrbracket$  = a map from source identifiers to types and target identifiers
- INVARIANT:  
     $x:t \in C$       means that
  - (1)     $\text{lookup } \llbracket C \rrbracket x = (\llbracket t \rrbracket^*, \%id\_x)$
  - (2)    the (target) type of  $\%id\_x$  is  $\llbracket t \rrbracket^*$     (a pointer to  $\llbracket t \rrbracket$ )

# Interpretation of Variables

$$\left[ \frac{x:t \in L}{H;G;L \vdash_{lhs} x : t; \overline{T}} \right] \begin{array}{l} \text{as addresses} \\ \text{(which can be assigned)} \end{array} = (T, \text{true}, \%id\_x, [])$$

where  $(T, \%id\_x) = \text{lookup } \llbracket L \rrbracket x$   
and by invariant:  $\overline{T} = \llbracket t \rrbracket^*$

$$\left[ \frac{H;G;L \vdash_{lhs} lhs : t; r}{H;G;L \vdash lhs : t} \right] \begin{array}{l} \text{lhs as expressions} \\ \text{(which are values)} \end{array} = (\%tmp, T, \text{stream @ } [\%tmp = \text{load } T * \%ptr])$$

where  
 $\llbracket H;G;L \vdash_{lhs} lhs : t; \overline{T} \rrbracket = (\llbracket t \rrbracket, \text{true}, ptr, \text{stream})$

# Interpretation of Assignment Stmts

$$\left[ \frac{\begin{array}{l} H;G;L \vdash_{lhs} lhs : t; \top \\ H;G;L \vdash exp : t' \\ H \vdash t' \leq t \end{array}}{H;G;L;rt \vdash lhs = exp; \Rightarrow L; \perp} \right] = (\llbracket H;G;L \rrbracket, \text{ptr\_code} @ \text{exp\_code} @ [\text{store } T \ \%e\_op, \%ptr])$$

assignment to a lhs

where

$$\llbracket H;G;L \vdash_{lhs} lhs : t; \top \rrbracket = (\llbracket t \rrbracket, \text{true}, \text{ptr}, \text{ptr\_code})$$

and

$$\llbracket H;G;L \vdash_{lhs} exp : t' \rrbracket = (\llbracket t' \rrbracket, \%e\_op, \text{exp\_code})$$

# Other Judgments?

- Statement:  
 $\llbracket H;G;L; \text{rt} \vdash \text{stmt} \Rightarrow C' \rrbracket = \llbracket C' \rrbracket, \text{stream}$
- Declaration:  
 $\llbracket H;G;L \vdash \text{var } x = \text{exp} \Rightarrow G;L,x:t \rrbracket = \llbracket G;L,x:t \rrbracket, \text{stream}$

INVARIANT: stream is of the form:

stream' @  
[E %id\_x = alloca  $\llbracket t \rrbracket$ ;  
I store  $\llbracket t \rrbracket$  opn,  $\llbracket t \rrbracket^* \text{\%id\_x}$  ]

and  $\llbracket H;G;L \vdash \text{exp} : t \rrbracket = (\llbracket t \rrbracket, \text{opn}, \text{stream}')$

- Rest follow similarly



# COMPILING CONTROL

# Translating while

- Consider translating “while(e) s”:
  - Test the conditional, if true jump to the body, else jump to the label after the body.

$\llbracket C; \text{rt} \vdash \text{while}(e) \ s \Rightarrow C' \rrbracket = \llbracket C' \rrbracket,$

```
lpre:
    opn =  $\llbracket C \vdash e : \text{bool} \rrbracket$ 
    %test = icmp eq i1 opn, 0
    br %test, label %lpost, label %lbody
lbody:
     $\llbracket C; \text{rt} \vdash s \Rightarrow C' \rrbracket$ 
    br %lpre
lpost:
```

- Note: writing  $\text{opn} = \llbracket C \vdash e : \text{bool} \rrbracket$  is pun
  - translating  $\llbracket C \vdash e : \text{bool} \rrbracket$  generates *code* that puts the result into **opn**
  - In this notation there is implicit collection of the code

# Translating if-then-else

- Similar to while except that code is slightly more complicated because if-then-else must reach a merge and the else branch is optional.

$$\llbracket C; \text{rt} \vdash \text{if } (e_1) s_1 \text{ else } s_2 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket$$

```
    opn =  $\llbracket C \vdash e : \text{bool} \rrbracket$ 
    %test = icmp eq i1 opn, 0
    br %test, label %else, label %then
then:
     $\llbracket C; \text{rt} \vdash s_1 \Rightarrow C' \rrbracket$ 
    br %merge
else:
     $\llbracket C; \text{rt} \vdash s_2 \Rightarrow C' \rrbracket$ 
    br %merge
merge:
```

# Connecting this to Code

- Instruction streams:
  - Must include labels, terminators, and “hoisted” global constants
- Must post-process the stream into a control-flow-graph
- See frontend.ml from HW4



# 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 %tmp4, label %else, label %then
```

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

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

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

# 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.  $\llbracket C \vdash e_1 + e_2 : \text{int} \rrbracket = (\text{i64}, \%tmp, [\%tmp = \text{add i64 } \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket])$
- But when the boolean expression we're compiling appears in a test, the program jumps to one label or another after the comparison but otherwise 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
- Make up new "judgement" that is similar to  $\llbracket C \vdash e : \text{bool} \rrbracket$  but has a different semantics. Call it  $\llbracket C \vdash e : \text{bool}@ \rrbracket$

# 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 the value:

$$\llbracket C \vdash e : \text{bool@} \rrbracket \text{ ltrue lfalse} = \text{stream}$$
$$\llbracket C, \text{rt} \vdash \text{if } (e) \text{ then } s_1 \text{ else } s_2 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket,$$

Notes:

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

```
    insns3
  then:
     $\llbracket s_1 \rrbracket$ 
    br %merge
  else:
     $\llbracket s_2 \rrbracket$ 
    br %merge
  merge:
```

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

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

---

$$\llbracket C \vdash \text{false} : \text{bool}@ \rrbracket \text{ ltrue lfalse} = [\text{br } \% \text{lfalse}] \quad \text{FALSE}$$

---

$$\llbracket C \vdash \text{true} : \text{bool}@ \rrbracket \text{ ltrue lfalse} = [\text{br } \% \text{ltrue}] \quad \text{TRUE}$$

---

$$\llbracket C \vdash e : \text{bool}@ \rrbracket \text{ lfalse ltrue} = \text{insns} \quad \text{NOT}$$

---

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

# Short Circuit Evaluation

Idea: build the logic into the translation

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

insns<sub>1</sub>  
right:  
insns<sub>2</sub>

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

insns<sub>1</sub>  
right:  
insns<sub>2</sub>

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
```