Lecture 15

# CIS 341: COMPILERS

# Announcements

- HW4:  OAT v. 1.0
  - Parsing & basic code generation
  - **Due: March 28th**


- No lecture on Thursday, March 22
  - Dr. Z will be away

# Adding Integers to Lambda Calculus

exp ::=
    | ...
    | n                                  *constant integers*
    | $exp_1 + exp_2$              *binary arithmetic operation*

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

$n\{v/x\}$        = n              *constants have no free vars.*
$(e_1 + e_2)\{v/x\}$ $= (e_1\{v/x\} + e_2\{v/x\})$    *substitute everywhere*

$$\frac{exp_1 \Downarrow n_1 \quad exp2 \Downarrow n_2}{exp_1 + exp_2 \;\; \Downarrow (n1 \; [\![+]\!] \; n2)}$$

Object-level '+'             Meta-level '+'

Compiling lambda calculus to straight-line code.

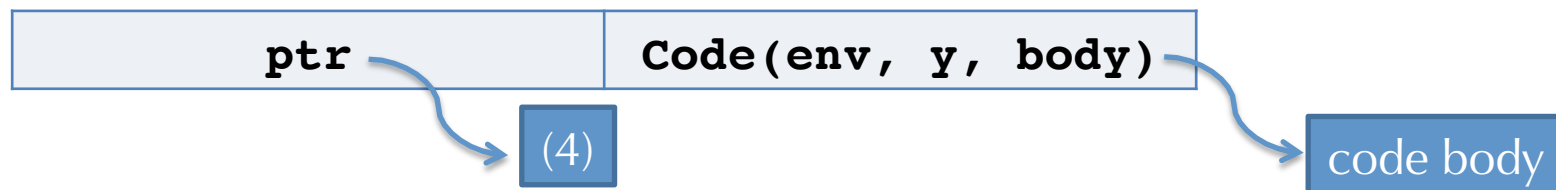Representing evaluation environments at runtime.

# CLOSURE CONVERSION

# Compiling First-class Functions

- To implement first-class functions on a processor, there are two problems:
  - First: we must implement substitution of free variables
  - Second: we 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.
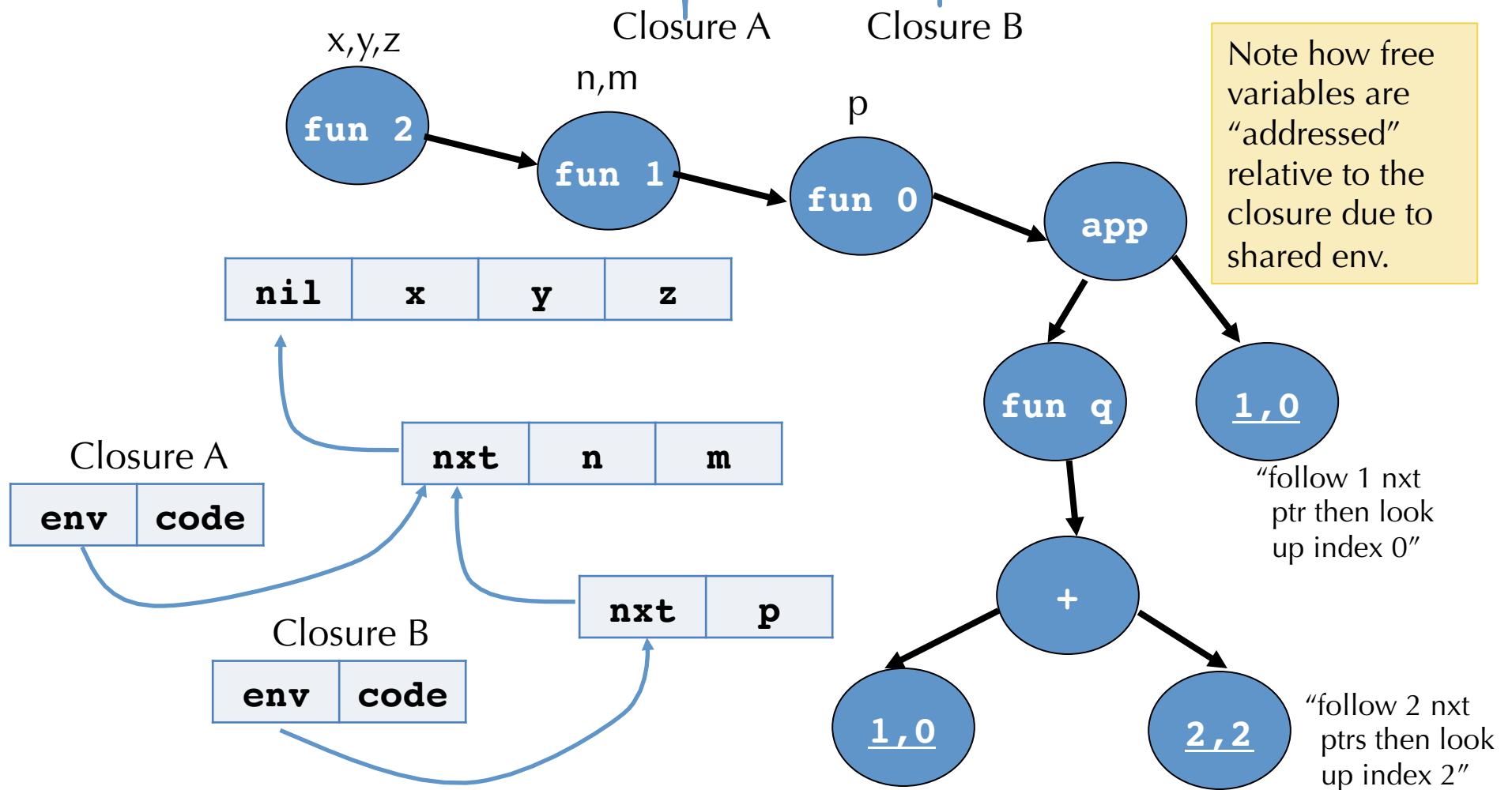
| ptr | Code(env, y, body) |
|-----|--------------------|

(4)

code body

- 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

- As we saw, the simple closure conversion algorithm 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 datastructure 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

x,y,z

n,m

p

Note how free variables are "addressed" relative to the closure due to shared env.

**fun 2** → **fun 1** → **fun 0** → **app**

| nil | x | y | z |
|-----|---|---|---|

Closure A

| env | code |
|-----|------|

| nxt | n | m |
|-----|---|---|

Closure B

| env | code |
|-----|------|

| nxt | p |
|-----|---|

**fun q**

**1,0**

"follow 1 nxt ptr then look up index 0"

**+**

**1,0**        **2,2**

"follow 2 nxt ptrs then look up index 2"

Scope, Types, and Context

# STATIC ANALYSIS

# Variable Scoping

- Consider the problem of determining whether a programmer-declared variable is in scope.

- Issues:
  - Which variables are available at a given point in the program?
  - Shadowing – is it permissible to re-use the same identifier, or is it an error?

- Example: The following program 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?

# Contexts and Inference Rules

- Need to keep track of contextual information.
  - What variables are in scope?
  - What are their types?

- How do we describe this?
  - In the compiler there's a mapping from variables to information we know about them.

# Why Inference Rules?

- They are a compact, precise way of specifying language properties.
    - E.g. ~20 pages for full Java vs. 100's of pages of prose Java Language Spec.

- Inference rules correspond closely to the recursive AST traversal that implements them

- Type checking (and type inference) is nothing more than attempting to prove a different judgment ( G;L ⊢ e : t ) by searching backwards through the rules.
- Compiling in a context is nothing more than a collection of inference rules specifying yet a different judgment ( G ⊢ src ⇒ target )
    - Moreover, the compilation judgment is similar to the typechecking judgment

- Strong mathematical foundations
    - The "Curry-Howard correspondence":  Programming Language ~ Logic, Program ~ Proof, Type ~ Proposition
    - See CIS 500 next Fall if you're interested in type systems!

# Inference Rules

- We can read a judgment $G;L \vdash e : t$ as
  "the expression e is well typed and has type t"
- For any environment G, expression e, and statements $s_1$, $s_2$.

$$G;L;rt \vdash \text{if } (e)\ s_1 \text{ else } s_2$$

holds if $G;L \vdash e : bool$ and $G;L;rt \vdash s_1$ and $G;L;rt \vdash s_2$
all hold.

- More succinctly: we summarize these constraints as an *inference rule*:

Premises $\left\{ \rule{0pt}{20pt}\right.$

$$G;L \vdash e : bool \qquad G;L;rt \vdash s_1 \qquad G;L;rt \vdash s_2$$

Conclusion $\left\{ \rule{0pt}{20pt}\right.$

$$G;L;rt \vdash \text{if } (e)\ s_1 \text{ else } s_2$$

- This rule can be used for *any* substitution of the syntactic
  metavariables G, e, $s_1$ and $s_2$.

# 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 oat0-defn.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);
```

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

# Example Derivation

$$\mathcal{D}_1 \quad = \quad \cfrac{\cfrac{\cfrac{\cfrac{}{G_0\,;\cdot \vdash 0 : \mathtt{int}}\;[\mathbf{INT}]}{G_0\,;\cdot \vdash 0 : \mathtt{int}}\;[\mathbf{CONST}]}{G_0\,;\cdot \vdash \mathtt{var}\ x_1 = 0 \Rightarrow \cdot, x_1\!:\!\mathtt{int}}\;[\mathbf{DECL}]}{G_0\,;\cdot\,;\mathtt{int} \vdash \mathtt{var}\ x_1 = 0\,; \Rightarrow \cdot, x_1\!:\!\mathtt{int}}\;[\mathbf{SDECL}]$$

$$\mathcal{D}_2 \quad = \quad \cfrac{\cfrac{\cfrac{\cfrac{}{\vdash + : (\mathtt{int}, \mathtt{int}) \rightarrow \mathtt{int}}[\mathbf{ADD}] \quad \cfrac{x_1\!:\!\mathtt{int} \in \cdot, x_1\!:\!\mathtt{int}}{G_0\,;\cdot, x_1\!:\!\mathtt{int} \vdash x_1 : \mathtt{int}}[\mathbf{VAR}] \quad \cfrac{x_1\!:\!\mathtt{int} \in \cdot, x_1\!:\!\mathtt{int}}{G_0\,;\cdot, x_1\!:\!\mathtt{int} \vdash x_1 : \mathtt{int}}\genfrac{}{}{0pt}{}{[\mathbf{VAR}]}{[\mathbf{BOP}]}}{G_0\,;\cdot, x_1\!:\!\mathtt{int} \vdash x_1 + x_1 : \mathtt{int}}}{G_0\,;\cdot, x_1\!:\!\mathtt{int}\,;\mathtt{int} \vdash \mathtt{var}\ x_2 = x_1 + x_1\,; \Rightarrow \cdot, x_1\!:\!\mathtt{int}, x_2\!:\!\mathtt{int}}[\mathbf{DECL}]}{G_0\,;\cdot, x_1\!:\!\mathtt{int}\,;\mathtt{int} \vdash \mathtt{var}\ x_2 = x_1 + x_1\,; \Rightarrow \cdot, x_1\!:\!\mathtt{int}, x_2\!:\!\mathtt{int}}[\mathbf{SDECL}]$$

# Example Derivation

$$x_1\texttt{:int} \in \cdot, x_1\texttt{:int}, x_2\texttt{:int} \ ;$$

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

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

# Why Inference Rules?

- They are a compact, precise way of specifying language properties.
    - E.g. ~20 pages for full Java vs. 100's of pages of prose Java Language Spec.

- Inference rules correspond closely to the recursive AST traversal that implements them

- Compiling in a context is nothing more an "interpretation" of the inference rules that specify typechecking*:  $[\![C \vdash e : t]\!]$
    - Compilation follows the typechecking judgment

- Strong mathematical foundations
    - The "Curry-Howard correspondence":  Programming Language ~ Logic, Program ~ Proof, Type ~ Proposition
    - See CIS 500 next Fall if you're interested in type systems!

*Here (and later) we'll write context C for `G;L`, the combination of the global and local contexts.

# 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?
$$[\![C \vdash e : t]\!] = \quad ?$$

- $[\![t]\!]$ is a target type

- $[\![e]\!]$ translates to a (potentially empty) sequence of instructions, that, when run, computes the result into some operand

- INVARIANT:  if   $[\![C \vdash e : t]\!]$ = ty, operand , stream
               then the type (at the target level) of the operand is ty=$[\![t]\!]$

# Example

- C ⊢ 341 + 5 : int          what is ⟦ C ⊢ 341 + 5 : int⟧  ?

⟦ ⊢ 341 : int ⟧ = (i64, Const 341, [])          ⟦⊢ 5 : int⟧ = (i64, Const 5, [])

--------------------------------------          ---------------------------------------

⟦C ⊢ 341 : int⟧ = (i64, Const 341, [])          ⟦C ⊢ 5 : int⟧ = (i64, Const  5, [])

-----------------------------------------------------------------------------------

⟦C ⊢ 341 + 5 : int⟧ =   (i64, %tmp, [%tmp = add i64 (Const 341) (Const 5)])

# What about the Context?

- What is $[\![C]\!]$?
- Source level C has bindings like:   x:int, y:bool
  - We think of it as a finite map from identifiers to types


- What is the interpretation of C at the target level?


- $[\![C]\!]$ maps source identifiers, "x" to source types and $[\![x]\!]$


- What is the interpretation of a variable $[\![x]\!]$ at the target level?
  - How are the variables used in the type system?

$$\frac{x:t \in L}{G\,;L \vdash x : t} \quad \text{TYP\_VAR}$$

as expressions
(which denote values)

$$\frac{x:t \in L \quad G\,;L \vdash exp : t}{G\,;L\,;rt \vdash x = exp\,; \Rightarrow L} \quad \text{TYP\_ASSN}$$

as addresses
(which can be assigned)

# Interpretation of Contexts

- ⟦C⟧ = a map from source identifiers to types and target identifiers

- INVARIANT:
    x:t ∈ C        means that

        (1)     lookup ⟦C⟧ x = (t, `%id_x`)
        (2)     the (target) type of `%id_x`  is ⟦t⟧*     (a pointer to ⟦t⟧)

# Interpretation of Variables

- Establish invariant for expressions:

$$\left[\!\!\left[\dfrac{x:t \in L}{G\,;L \vdash x : t} \;\; \text{TYP\_VAR}\right]\!\!\right]$$

as expressions
(which denote values)

$=$ `(%tmp,  [%tmp = load i64* %id_x])`

where (`i64`, `%id_x`) = lookup ⟦L⟧ x

- What about statements?

$$\left[\!\!\left[\dfrac{x:t \in L \quad G\,;L \vdash exp : t}{G\,;L\,;rt \vdash x = exp;\; \Rightarrow L} \;\; \text{TYP\_ASSN}\right]\!\!\right]$$

as addresses
(which can be assigned)

$=$ stream @

`[store ⟦t⟧ opn, ⟦t⟧* %id_x]`

where (t, `%id_x`) = lookup ⟦L⟧ x
and ⟦G;L ⊢ exp : t⟧ = (⟦t⟧, `opn`, stream)

# Other Judgments?

- Statement:
  $\llbracket C; rt \vdash stmt \Rightarrow C' \rrbracket = \llbracket C' \rrbracket$ , stream


- Declaration:
  $\llbracket G;L \vdash t\ x = exp \Rightarrow G;L,x:t \rrbracket = \llbracket G;L,x:t \rrbracket$, stream

  INVARIANT:   stream is of the form:
  ```
  stream' @
  [ %id_x = alloca 〚t〛;
     store 〚t〛 opn, 〚t〛* %id_x ]
  ```

  and    $\llbracket G;L \vdash exp : t \rrbracket = (\llbracket t \rrbracket,$ `opn`, 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.

$⟦$`C;rt` $⊢$ `while(e) s` $⇒$ `C'`$⟧$ $=$ $⟦$`C'`$⟧$,

```
lpre:
    opn = ⟦C ⊢ e : bool⟧
    %test = icmp eq i1 opn, 0
    br %test, label %lpost, label %lbody
lbody:
    ⟦C;rt ⊢ s ⇒ C'⟧
    br %lpre
lpost:
```

- Note: writing  `opn` $=$ $⟦$`C` $⊢$ `e : bool`$⟧$   is pun
  - translating  $⟦$C $⊢$ e : bool$⟧$ 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.

$$[\![C;rt \vdash \text{if } (e_1) \text{ } s_1 \text{ else } s_2 \Rightarrow C'}]\!] = [\![C']\!]$$

```
        opn = ⟦C ⊢ e : bool⟧
        %test = icmp eq i1 opn, 0
        br %test, label %else, label %then
then:
         ⟦C;rt ⊢ s₁ ⇒ C'⟧
         br %merge
else:
        ⟦C; rt s₂ ⇒ C'⟧
         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 $[\![e]\!]$ to produce a value
  - $[\![C \vdash e : t]\!]$ = (ty, operand, stream)
  - e.g.  $[\![C \vdash e_1 + e_2 : \text{int}]\!]$   = (i64, `%tmp`,   [`%tmp = add` $[\![e_1]\!]$ $[\![e_2]\!]$])

- But when the 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

# Idea: Use a different translation for tests

Usual Expression translation:

$\llbracket C \vdash e : t \rrbracket$ = (ty, operand, stream)

Conditional branch translation of booleans,
   without materializing the value:

$\llbracket C \vdash e : bool@ \rrbracket$ ltrue lfalse = stream

$\llbracket C, rt \vdash$ if (e) then s1 else s2 $\Rightarrow$ 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:
    ⟦s1⟧
    br %merge
else:
    ⟦s₂⟧
    br %merge
merge:
```

where

$\llbracket C, rt \vdash s_1 \Rightarrow C' \rrbracket$ = $\llbracket C' \rrbracket$, $insns_1$
$\llbracket C, rt \vdash s_2 \Rightarrow C'' \rrbracket$ = $\llbracket C'' \rrbracket$, $insns_2$
$\llbracket C \vdash e : bool@ \rrbracket$ then else = $insns_3$

# Short Circuit Compilation: Expressions

- ⟦C ⊢ e : bool@⟧ ltrue lfalse = insns

$$\frac{\rule{0pt}{0pt}}{\text{⟦}C \vdash \text{false} : \text{bool@⟧ ltrue lfalse} = \texttt{[br \%lfalse]}} \quad \text{FALSE}$$

$$\frac{\rule{0pt}{0pt}}{\text{⟦}C \vdash \text{true} : \text{bool@⟧ ltrue lfalse} = \texttt{[br \%ltrue]}} \quad \text{TRUE}$$

$$\frac{\text{⟦}C \vdash e : \text{bool@⟧ lfalse ltrue} = \text{insns}}{\text{⟦}C \vdash \text{!}e : \text{bool@⟧ ltrue lfalse} = \text{insns}} \quad \text{NOT}$$

# Short Circuit Evaluation

Idea:  build the logic into the translation

$$\frac{[\![C \vdash e1 : bool@]\!]\ ltrue\ right = insns_1 \quad [\![C \vdash e2 : bool@]\!]\ ltrue\ lfalse = insns_2}{[\![C \vdash e1 \,|\, e2 : bool@]\!]\ ltrue\ lfalse =}$$

```
    insns₁
right:
    insn₂
```

$$\frac{[\![C \vdash e1 : bool@]\!]\ right\ lfalse = insns_1 \quad [\![C \vdash e2 : bool@]\!]\ ltrue\ lfalse = insns_2}{[\![C \vdash e1 \,\&\, e2 : bool@]\!]\ ltrue\ lfalse =}$$

```
    insns₁
right:
    insn₂
```

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
```