

Lecture 15

CIS 341: COMPILERS

Announcements

- HW4: OAT v. 1.0
 - Parsing & basic code generation
 - **Due: Wednesday, March 23rd**



UNTYPED LAMBDA CALCULUS

(Untyped) Lambda Calculus

- The lambda calculus is a minimal programming language.
 - Note: we're writing `(fun x -> e)` lambda-calculus notation: $\lambda x. e$

Abstract syntax in OCaml:

```
type exp =  
  | Var of var          (* variables          *)  
  | Fun of var * exp    (* functions: fun x -> e *)  
  | App of exp * exp    (* function application *)
```

Concrete syntax:

```
exp ::=  
  | x          variables  
  | fun x -> exp functions  
  | exp1 exp2 function application  
  | ( exp )    parentheses
```

Operational Semantics

- Key operation: *capture-avoiding substitution*: $e_2\{e_1/x\}$
 - replaces all free occurrences of x in e_2 by e_1
 - must respect scope and alpha equivalence (renaming)
- *Reduction Strategies*
Various ways of *simplifying* (or “reducing”) lambda calculus terms.
 - *call-by-value evaluation*:
 - simplify the function argument *before* substitution
 - *does not* reduce under lambda (a.k.a. fun)
 - *call-by-name evaluation*:
 - *does not* simplify the argument before substitution
 - *does not* reduce under lambda
 - *weak-head normalization*:
 - does not simplify the argument before substitution
 - does not reduce under lambda
 - works on open terms, “suspending” reduction at variables
 - *normal order reduction*:
 - *does* reduce under lambda
 - first does weak-head normalization and then recursively continues to reduce
 - works on open terms – guaranteed to find a “normal form” if such a form exists

A “normal form” is one that has no substitution steps possible, i.e., there are no subterms of the form $(\text{fun } x \rightarrow e_1) e_2$ anywhere.

CBV Operational Semantics

- This is *call-by-value* semantics:
function arguments are evaluated before substitution

$$\frac{}{v \Downarrow v}$$

“Values evaluate to themselves”

$$\frac{\text{exp}_1 \Downarrow (\text{fun } x \rightarrow \text{exp}_3) \quad \text{exp}_2 \Downarrow v \quad \text{exp}_3\{v/x\} \Downarrow w}{\text{exp}_1 \text{ exp}_2 \Downarrow w}$$

“To evaluate function application: Evaluate the function to a value, evaluate the argument to a value, and then substitute the argument for the function. ”

CBN Operational Semantics

- This is *call-by-name* semantics:
function arguments are evaluated before substitution

$$\frac{}{v \Downarrow v}$$

“Values evaluate to themselves”

$$\frac{\text{exp}_1 \Downarrow (\text{fun } x \rightarrow \text{exp}_3) \quad \text{exp}_3\{\text{exp}_2/x\} \Downarrow w}{\text{exp}_1 \text{ exp}_2 \Downarrow w}$$

“To evaluate function application: Evaluate the function to a value, substitute the argument into the function body, and then keep evaluating.”



See fun.ml

Examples of encoding Booleans, integers, conditionals, loops, etc., in untyped lambda calculus.

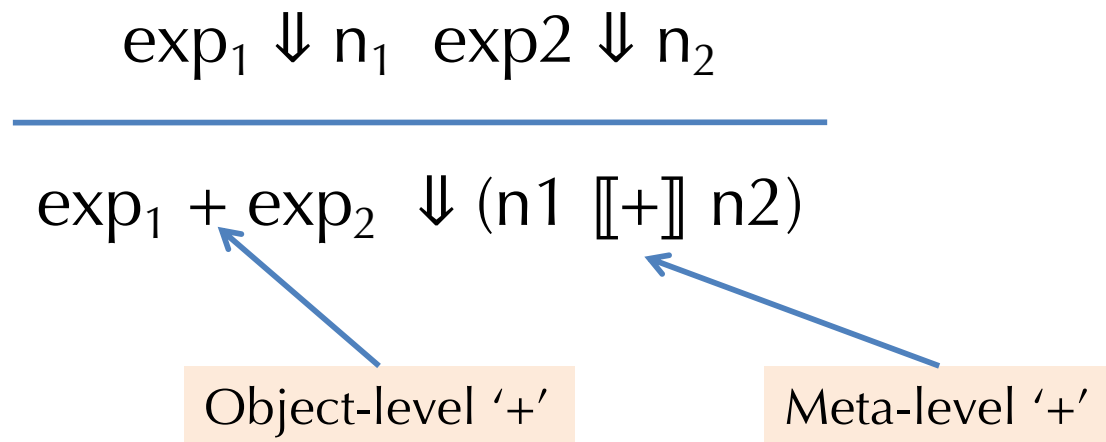
IMPLEMENTING THE INTERPRETER

Adding Integers to Lambda Calculus

$\text{exp} ::=$
| ...
| n *constant integers*
| $\text{exp}_1 + \text{exp}_2$ *binary arithmetic operation*

$\text{val} ::=$
| $\text{fun } x \rightarrow \text{exp}$ *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*





Scope, Types, and Context

STATIC ANALYSIS

Scope-Checking Lambda Calculus

- Consider how to identify “well-scoped” lambda calculus terms
 - Recall the free variable calculation
 - Given: G , a set of variable identifiers, e , a term of the lambda calculus
 - Judgment:* $G \vdash e$ means “the free variables of e are included in G ”
 $\text{fv}(e) \subseteq G$

$$\begin{aligned}\text{fv}(x) &= \{x\} \\ \text{fv}(\text{fun } x \rightarrow \text{exp}) &= \text{fv}(\text{exp}) \setminus \{x\} \quad ('x' \text{ is a bound in exp}) \\ \text{fv}(\text{exp}_1 \text{ exp}_2) &= \text{fv}(\text{exp}_1) \cup \text{fv}(\text{exp}_2)\end{aligned}$$

$$\frac{x \in G}{G \vdash x}$$

“the variable x is free”

$$\frac{G \vdash e_1 \quad G \vdash e_2}{G \vdash e_1 \text{ } e_2}$$

“ G contains the free variables of e_1 and e_2 ”

$$\frac{G \cup \{x\} \vdash e}{G \vdash \text{fun } x \rightarrow e}$$

“ x is available in the function body e ”

Scope-checking Code

- Compare the OCaml code to the inference rules:
 - structural recursion over syntax
 - the check either “succeeds” or “fails”

```
let rec scope_check (g:VarSet.t) (e:exp) : unit =  
  begin match e with  
    | Var x -> if VarSet.member x g then () else failwith (x ^ "not in scope")  
    | App(e1, e2) -> ignore (scope_check g e1); scope_check g e2  
    | Fun(x, e) -> scope_check (VarSet.union g (VarSet.singleton x)) e  
  end
```

$$\frac{x \in G}{G \vdash x}$$

$$\frac{G \vdash e_1 \quad G \vdash e_2}{G \vdash e_1 e_2}$$

$$\frac{G \cup \{x\} \vdash e}{G \vdash \text{fun } x \rightarrow e}$$

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 process?
 - In the compiler there's a mapping from variables to information we know about them.
 - This is "contextual information"
- How do we use that information to implement a scope checker?

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 ($G;L \vdash e : t$) by searching backwards through the rules.
- Compiling in a context is nothing more than a collection of inference rules specifying yet a different judgment ($G \vdash \text{src} \Rightarrow \text{target}$)
 - 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 $G;L \vdash e : \text{bool}$ 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	$G;L \vdash e : \text{bool}$	$G;L;rt \vdash s_1$	$G;L;rt \vdash s_2$
Conclusion	$G;L;rt \vdash \text{if } (e) s_1 \text{ else } s_2$		

- 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 oat-v1-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);
```

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

Example Derivation

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

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

Example Derivation

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

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

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*: $\llbracket C \vdash e : t \rrbracket$
 - 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?

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

- $\llbracket t \rrbracket$ is a target type
- $\llbracket e \rrbracket$ translates to a (potentially empty) sequence 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

- $C \vdash 341 + 5 : \text{int}$ what is $\llbracket C \vdash 341 + 5 : \text{int} \rrbracket$?

$\llbracket \vdash 341 : \text{int} \rrbracket = (\text{i64}, \text{Const } 341, [])$

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

 $\llbracket C \vdash 341 : \text{int} \rrbracket = (\text{i64}, \text{Const } 341, [])$

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

 $\llbracket C \vdash 341 + 5 : \text{int} \rrbracket = (\text{i64}, \%tmp, [\%tmp = \text{add i64 (Const } 341) (\text{Const } 5)])$

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 source 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?

$$\frac{x:t \in L}{G;L \vdash x:t} \quad \text{TYP_VAR}$$

as expressions
(which denote values)

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

as addresses
(which can be assigned)

Interpretation of Contexts

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

Interpretation of Variables

- Establish invariant for expressions:

$$\left[\frac{x:t \in L}{G;L \vdash x:t} \text{ TYP_VAR} \right] = (\%tmp, [\%tmp = \text{load } i64* \%id_x])$$

as expressions
(which denote values)

where $(i64, \%id_x) = \text{lookup } \llbracket L \rrbracket x$

- What about statements?

$$\left[\frac{x:t \in L \quad G;L \vdash exp:t}{G;L;rt \vdash x = exp; \Rightarrow L} \text{ TYP_ASSN} \right] = \text{stream @}$$

as addresses
(which can be assigned)

$[\text{store } \llbracket t \rrbracket \text{ opn}, \llbracket t \rrbracket * \%id_x]$

where $(t, \%id_x) = \text{lookup } \llbracket L \rrbracket x$
and $\llbracket G;L \vdash exp:t \rrbracket = (\llbracket t \rrbracket, \text{opn}, \text{stream})$

Other Judgments?

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

INVARIANT: stream is of the form:

```
stream' @  
[ %id_x = alloca  $\llbracket t \rrbracket$ ;  
  store  $\llbracket t \rrbracket$  opn,  $\llbracket t \rrbracket^* \%id\_x$  ]
```

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

- 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. $\llbracket C \vdash e_1 + e_2 : \text{int} \rrbracket = (\text{i64}, \%tmp, [\%tmp = \text{add } \llbracket e_1 \rrbracket \llbracket e_2 \rrbracket])$
- 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:

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

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

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

insns₁
right:
insn₂

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

insns₁
right:
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
```