

Lecture 16

EECS 483: COMPILER CONSTRUCTION

Announcements

- Midterm
 - Grades will be released after we review the results on Wednesday's class (3/20).
- HW4: OAT v.1.0
 - Parsing & translation to LLVM IR
 - Helps to start early!
 - **Due: Tuesday, March 26th**

Personal Announcement

- My wife and I are having a baby
 - Due date: March 30, but could be any day now
 - We will have guest lectures (by our GSI Eric or a guest professor) for at least 4 lectures after the baby comes.
 - I will have my office hours remotely, will not be available for scheduled office hours.

Implementing First-Class Functions

- First attempt: Functions as Code
 - Represent a function value as its code
 - What about *local function definitions*?
 - `let f = fun x -> fun y -> x + y`
 - Every time we call `f 0`, `f 5`, `f 256`, we get a *different* function
 - In a substitution-based interpreter, we substitute a value in and get a different term each time we call the function.
 - Compilation
 - Infeasible to implement all of these possible functions statically in memory, impossible if the domain of the function is infinite!
 - Requires *runtime code generation*, which itself has high runtime overhead.
 - Usually not used for arbitrary first class functions, only in specialized situations

Implementing First-Class Functions

- Closures
 - Consider
 - `let f = fun x -> fun y -> x + y`
 - Each function `f 0`, `f 5`, `f 256`,... is implemented by substitution:
 - `(fun y -> x + y){ 0 / x }`
 - `(fun y -> x + y){ 5 / x }`
 - `(fun y -> 256 + y){ 256 / x }`
 - Idea: represent a first-class function as a ***pair*** of
 - A piece of code with **free** variables
 - An **environment** that provides all of the values of the free variables
 - In compilation
 - The code with free variables can be a code pointer to code that takes the environment as an argument
 - `fun (env, y) -> env.x + y`
 - The environment can be implemented in multiple ways
 - Array: fast access in the function
 - Linked list: more sharing between different closures

See fun.ml, cc.ml from lec15.zip

CLOSURES AND CLOSURE CONVERSION

Closure Conversion Summary

- A **closure** is a pair of an environment and a code pointer
 - the environment is a map data structure binding variables to values
 - environment could just be a list of the values (with known indices)
- Building a closure value:
 - code pointer is a function that takes an extra argument for the environment: $A \rightarrow B$ becomes $(\text{Env} * A \rightarrow B)$
 - body of the closure “projects out” then variables from the environment
 - creates the environment map by bundling the free variables
- Applying a closure:
 - project out the environment, invoke the function (pointer) with the environment and its “real” argument
- Hoisting:
 - Once closure converted, all functions can be lifted to the top level



Scope, Types, and Context

SEMANTIC ANALYSIS

Compilation in a Nutshell

Source Code

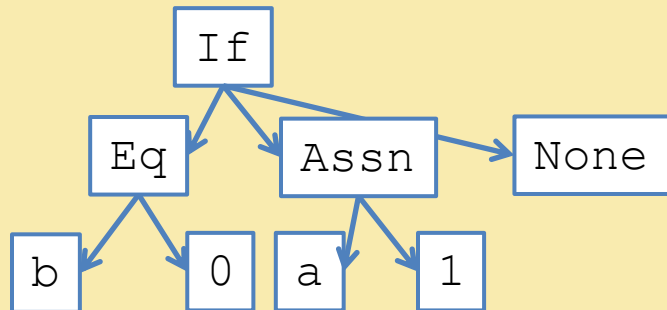
(Character stream)

```
if (b == 0) { a = 1; }
```

Token stream:

if	(b	==	0)	{	a	=	0	;	}
----	---	---	----	---	---	---	---	---	---	---	---

Abstract Syntax Tree:



Assembly Code

```
11: cmpq %eax, $0
   jeq 12
   jmp 13
12:
...
```

Lexical Analysis

Parsing

Analysis & Transformation

Backend

Intermediate code:

```
11: %cnd = icmp eq i64 %b, 0
   br i1 %cnd, label %12, label %13
12: store i64* %a, 1
   br label %13
13:
```

Most of the Remainder of the Course

Source Code

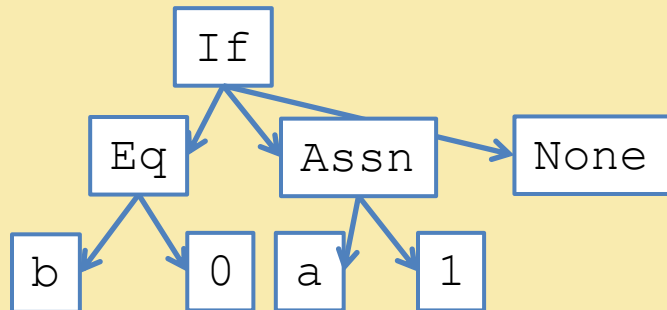
(Character stream)

```
if (b == 0) { a = 1; }
```

Token stream:

if	(b	==	0)	{	a	=	0	;	}
----	---	---	----	---	---	---	---	---	---	---	---

Abstract Syntax Tree:



Assembly Code

```
11: cmpq %eax, $0
    jeq 12
    jmp 13
12:
...
```

Lexical Analysis

Parsing

Analysis & Transformation

Backend

Intermediate code:

```
11: %cnd = icmp eq i64 %b, 0
    br i1 %cnd, label %12, label %13
12: store i64* %a, 1
    br label %13
13:
```

Static Program Analysis

- Static program analysis is analysis of a program at compile-time
- Used for two main purposes in the compiler:
 - Last stage of the frontend: “Type checking” or “Semantic Analysis”
 - Not every program that passes parsing is valid
 - `int main() { return x; }`
 - `int main() { return “hello world”; }`
 - If the type checker fails, the program is rejected, like a parse error
 - After the program passes the frontend, we consider it well-formed and will compile it.
 - During optimization: “static analysis”
 - We can do more optimizations if we know more about the program
 - Are these equivalent programs?
 - `int main() { int y = f(); return 0; }`
 - `int main() { return 0; }`
 - We can optimize the first to the second if we establish that `f` is side-effect free.
 - Since they take place after the frontend, the analysis never rejects the program
- Next few weeks: type checking, after that optimization and analysis

Type Checking as Grammar

	Specification	Implementation
Lexing	Regular Expressions	DFA
Parsing	CFG LL(1) grammars LR(1) grammars	Pushdown automata Recursive descent Shift/reduce parser
Type checking	Inference rules	Manual recursive descent

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?

Inference Rules

- We can read a judgment $G \vdash e$ as
“the expression e is well scoped and has free variables in G ”
- For any environment G , expression e , and statements s_1, s_2 .

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

holds if $G \vdash e$ and $G \vdash s_1$ and $G \vdash s_2$ all hold.

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

Premises	$G \vdash e$	$G \vdash s_1$	$G \vdash s_2$
Conclusion	$G \vdash \text{if } (e) s_1 \text{ else } s_2$		

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

Judgments

- A *judgment* is a (meta-syntactic) notation that *names* a relation among one or more sets.
 - The sets are usually built from object-language syntax elements and other “math” sets (e.g., integers, natural numbers, etc.)
 - We usually describe them using metavariables that range over the sets.
 - Often use domain-specific notation to ease reading.
 - The meaning of judgments, *i.e.*, which sets they represent, is defined by (collections of) inference rules
- Example: When we say “ $G \vdash e$ is a judgment where G is a context of variables and e is a term, defined by these [...] inference rules” that is shorthand for this “math speak”:
 - Let Var be the set of all (syntactic) variables
 - Let Exp be the set $\{e \mid e \text{ is a term of the untyped lambda calculus}\}$
 - Let $\mathcal{P}(\text{Var})$ be the (finite) powerset of variables (set of all finite sets)
 - Define $\text{well-scoped} \subseteq (\mathcal{P}(\text{Var}), \text{Exp})$ to be a relation satisfying the properties defined by the associated inference rules [...]
 - Then “ $G \vdash e$ ” is notation that means that $(G, e) \in \text{well-scoped}$

Scope-Checking Lambda Calculus

- Consider how to identify “well-scoped” lambda calculus terms
 - Given: G , a set of variable identifiers, e , a term of the lambda calculus
 - Judgment*: $G \vdash e$ “the free variables of e are included in G ”

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

“the variable x is free, but in scope”

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

“ 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
```

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

VAR

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

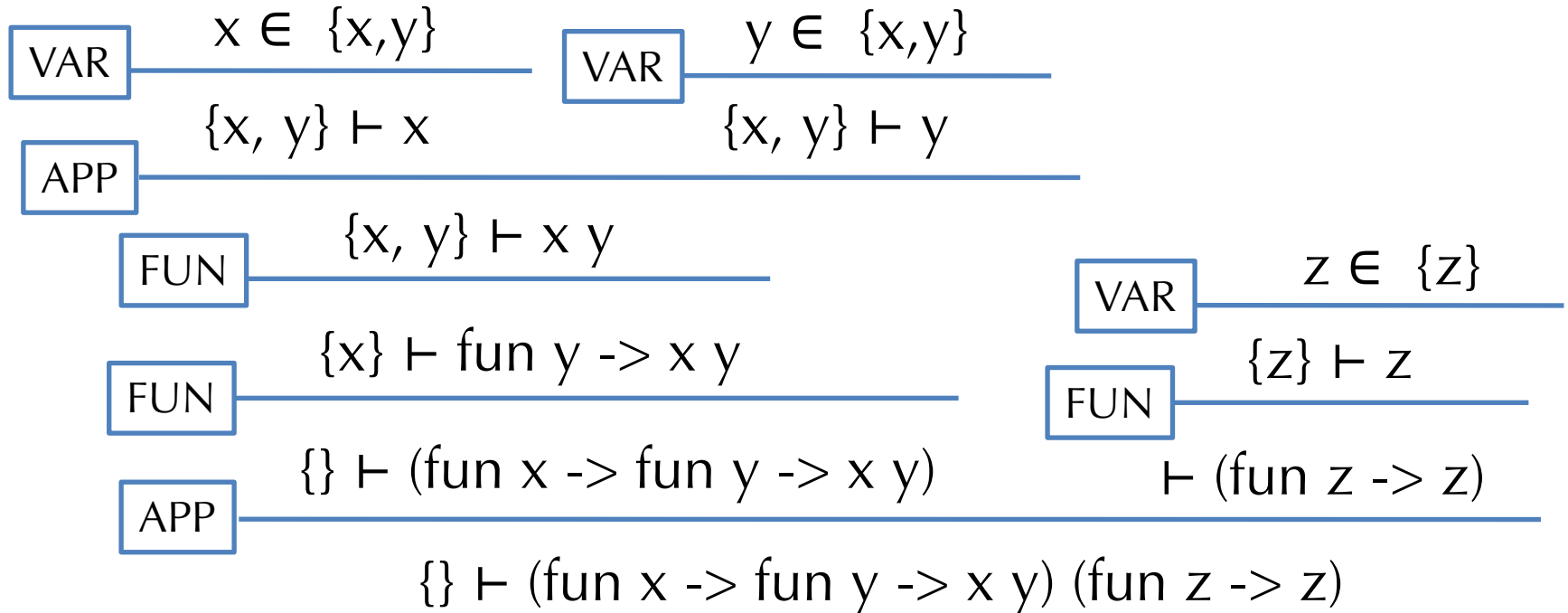
APP

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

FUN

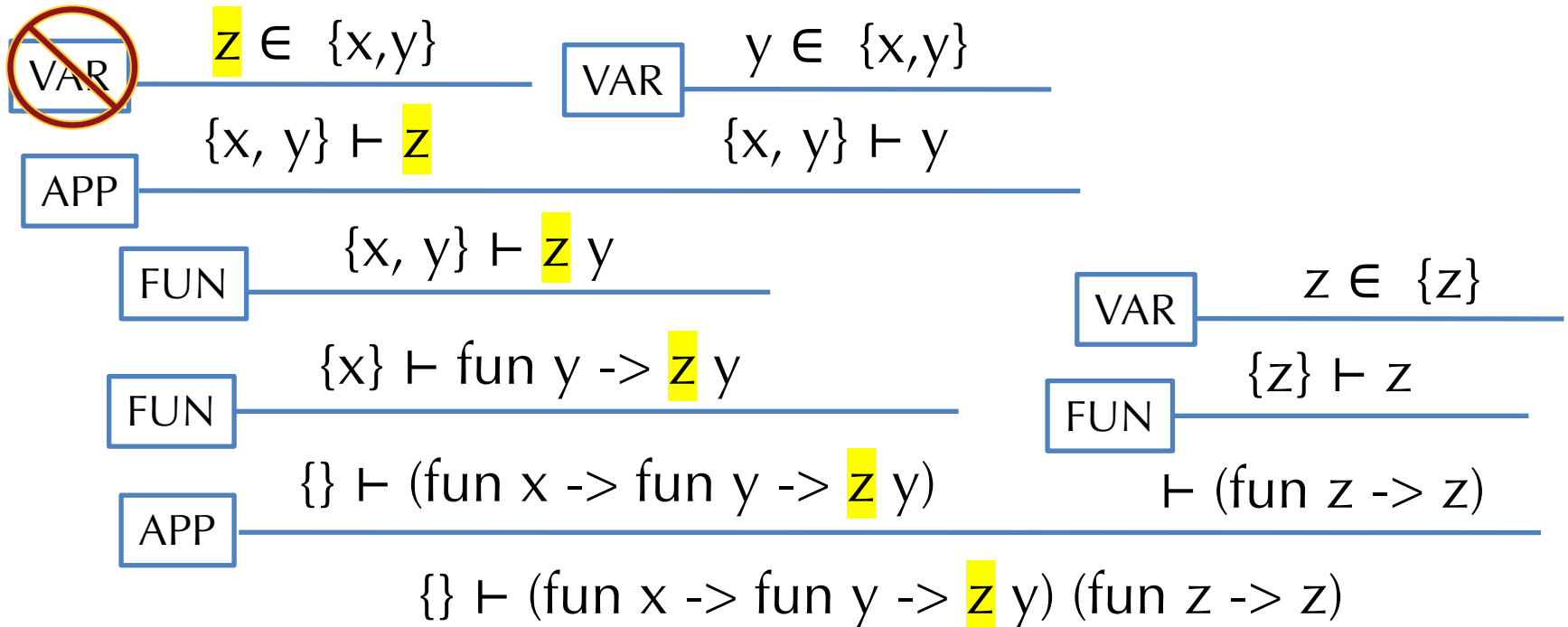
- The inference rules are a *specification* of the intended behavior of this scope checking code.
 - they don't specify the order in which the premises are checked

Example Derivation Tree



- Note: the OCaml function `scope_check` verifies the existence of this tree. The structure of the recursive calls when running `scope_check` is the same shape as this tree!
- Note that $x \in E$ is implemented by the function `VarSet.mem`

Example Failed Derivation



- This program is *not* well scoped
 - The variable z is not bound in the body of the left function.
 - The typing derivation fails because the VAR rule cannot succeed
 - (The other parts of the derivation are OK, though!)

Uses of the inference rules

- We can do proofs by induction on the structure of the derivation.
- For example:

Lemma: If $G \vdash e$ then $\text{fv}(e) \subseteq G$.

Proof.

By induction on the derivation that $G \vdash e$.

- case: VAR then we have $e = x$ (for some variable x) and $x \in G$. But $\text{fv}(e) = \text{fv}(x) = \{x\}$, but then $\{x\} \subseteq G$.

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

- case: APP then we have $e = e_1 e_2$ (for some $e_1 e_2$) and, by induction, we have $\text{fv}(e_1) \subseteq G$ and $\text{fv}(e_2) \subseteq G$, so $\text{fv}(e_1 e_2) = \text{fv}(e_1) \cup \text{fv}(e_2) \subseteq G$

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

- case: FUN then we have $e = (\text{fun } x \rightarrow e_1)$ for some x, e_1 and, by induction, we have $\text{fv}(e_1) \subseteq G \cup \{x\}$, but then we also have $\text{fv}(\text{fun } x \rightarrow e_1) = \text{fv}(e_1) \setminus \{x\} \subseteq ((G \cup \{x\}) \setminus \{x\}) \subseteq G$

$$\frac{G \cup \{x\} \vdash e_1}{G \vdash \text{fun } x \rightarrow e_1}$$

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

See tc.ml

STATICALLY RULING OUT PARTIALITY: TYPE CHECKING

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*

$\text{exp}_1 \Downarrow n_1 \quad \text{exp}_2 \Downarrow n_2$

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

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

NOTE: there are no rules for the case where exp_1 or exp_2 evaluate to functions! The semantics is *undefined* in those cases.

Type Checking / Static Analysis

- Recall the interpreter from the Eval3 module:

```
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., $3/0$, $3 + \text{true}$, ...)
- A compiler can't generate sensible code for this case.
 - A naïve implementation might “add” an integer and a function pointer

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

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

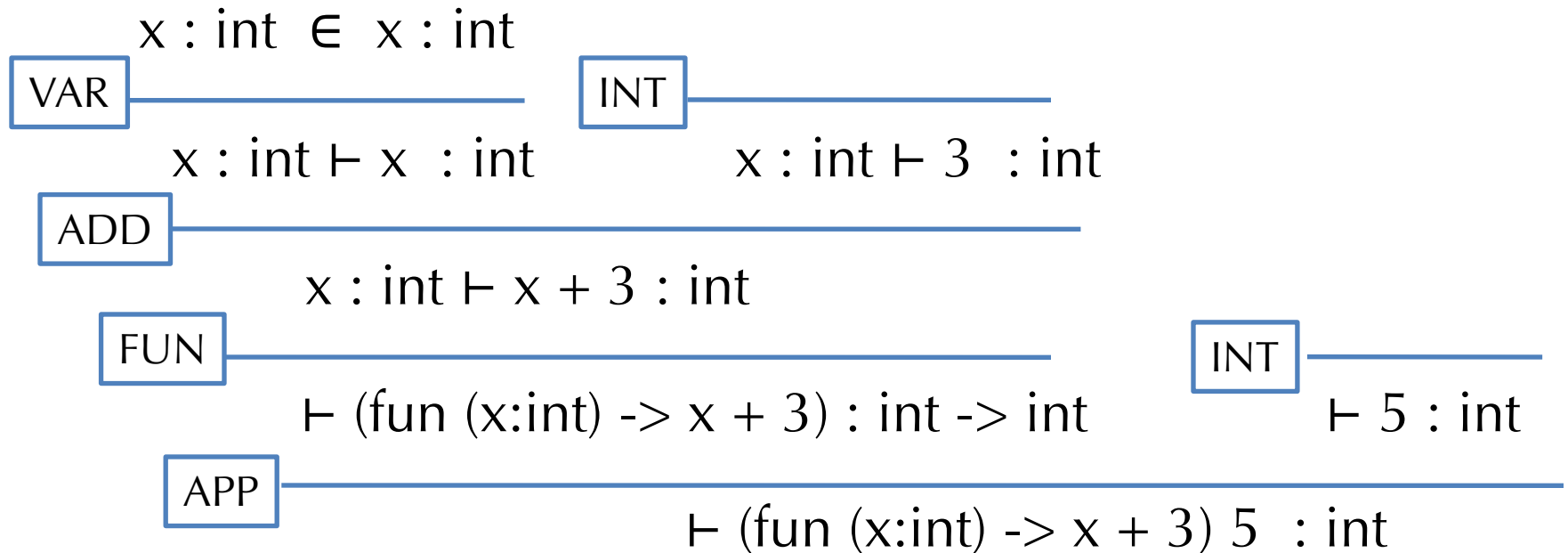
- Note how these rules correspond to the code.

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`

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

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

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

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)

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...

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 ($E \vdash e : t$) 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 rules are very similar in structure to the typechecking rules
- Strong mathematical foundations
 - The “Curry-Howard-Lambek correspondence”:
 - Programming Language : Logic : Category theory : Order Theory
 - Programs : Proof : Morphism : Inequality
 - Type : Proposition : Object : Element
 - See EECS 490, 590, my 598 if you're interested in type systems!



COMPILING

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

- $C \vdash 341 + 5 : \text{int}$ 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) 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, $\llbracket t \rrbracket^* \text{\%id_x}$]

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

$$\llbracket C \vdash e : \text{bool@} \rrbracket \text{ ltrue lfalse} = \text{stream}$$
$$\llbracket C, \text{rt} \vdash \text{if } (e) \text{ then } s1 \text{ else } s2 \Rightarrow C' \rrbracket = \llbracket C' \rrbracket,$$

Notes:

- takes two extra arguments: a “true” branch label and a “false” branch label.
- Doesn’t “return a value”

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

- 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

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

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

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

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

Short Circuit Evaluation

Idea: build the logic into the translation

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

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
```



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

Types as sets

Subsumption

TYPES, MORE GENERALLY

Arrays

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

NEW

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

$$E \vdash \text{new } T[e_1](e_2) : T[]$$

e_1 is the size of the newly allocated array. e_2 initializes the elements of the array.

INDEX

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

$$E \vdash e_1[e_2] : T$$

UPDATE

$$E \vdash e_1 : T[] \quad E \vdash e_2 : \text{int} \quad E \vdash e_3 : T$$

$$E \vdash e_1[e_2] = e_3 \text{ ok}$$

Note: These rules don't ensure that the array index is in bounds – that 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: $T \text{ ref}$

REF

$$E \vdash e : T$$

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

DEREF

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

$$E \vdash !e : T$$

ASSIGN

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

$$E \vdash e_1 := e_2 : \text{unit}$$

Note the similarity with the rules for arrays...