

Lecture 16

CIS 341: COMPILERS

Announcements

- HW4: OAT v. 1.0
 - Parsing & basic code generation
 - **Due: Wednesday, March 28th**
- HW5: OAT v. 2.0
 - records, function pointers, type checking, array-bounds checks, etc.
 - Due: Wednesday, April 11th



Scope, Types, and Context

STATIC ANALYSIS

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.

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?

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



See tc.ml

STATICALLY RULING OUT PARTIALITY: TYPE CHECKING

Notes about this Typechecker

- In the interpreter, we only evaluate the body of a function when it's applied.
- In the typechecker, we always check 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 also determine what type e_1 will return.
- Question: Why is this an approximation?
- Question: What if `well_typed` always returns `false`?

Contexts and Inference Rules

- Need to keep track of contextual information.
 - What variables are in scope?
 - What are their types?
 - What information do we have about each syntactic construct?
- What relationships are there among the syntactic objects?
 - e.g. is one type a subtype of another?
- How do we describe this information?
 - In the compiler there's a mapping from variables to information we know about them – the "context".
 - The compiler has a collection of (mutually recursive) functions that follow the structure of the syntax.

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. $x : \text{int}, b : \text{bool} \vdash x : \text{int}$

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 correspondence”: Programming Language ~ Logic, Program ~ Proof, Type ~ Proposition
 - See CIS 500 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 ”
- We can read a judgment $G;L \vdash s$ as “the statement s is well formed”
- For any environment G , expression e , and statements s_1, s_2 .

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

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

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

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

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

Simply-typed Lambda Calculus

- For the language in “tc.ml” we have five inference rules:

INT

$$\frac{}{E \vdash i : \text{int}}$$

VAR

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

ADD

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

FUN

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

APP

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

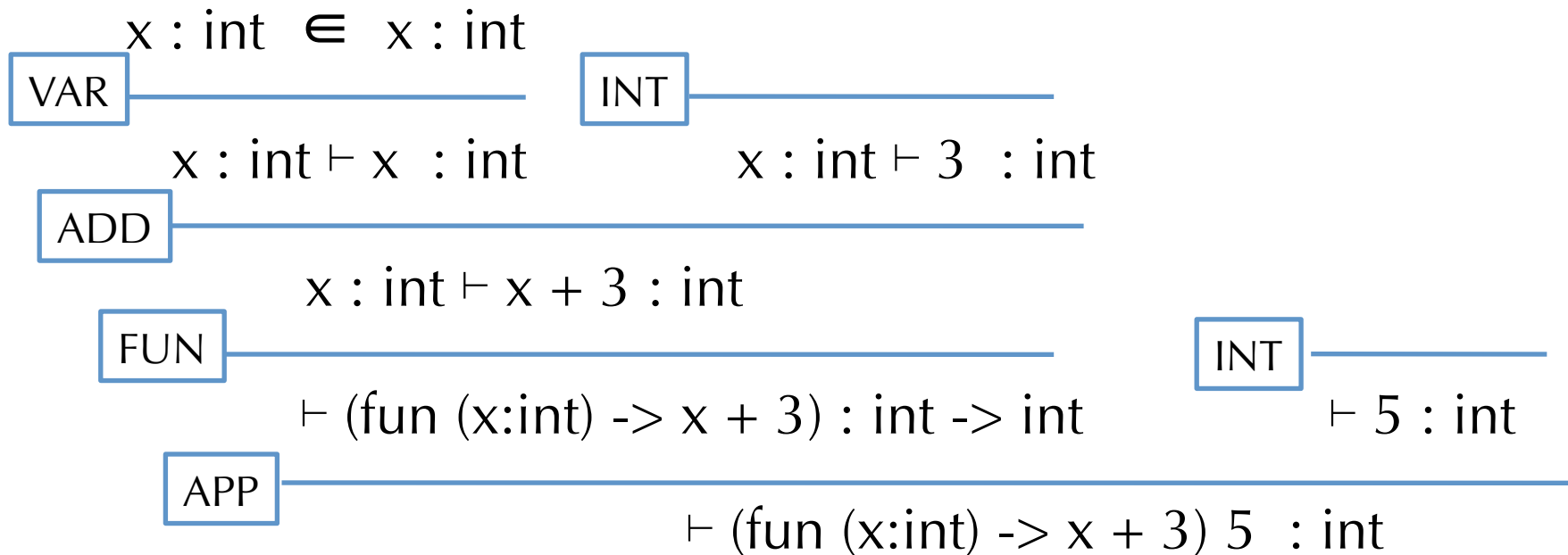
- Note how these rules correspond to the code.

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`

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);
```

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

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...



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

Types as sets

Subsumption

TYPES, MORE GENERALLY

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

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

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

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

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

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

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,  $\llbracket t \rrbracket^* \text{\%id\_x}$  ]
```

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

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...