Lecture 14
CIS 341: COMPILERS

Announcements

- Midterm Exam:
 - Not quite done grading yet!
 - Will be available Thursday

- Project 4 is available from the course web pages
 - Due on Thursday, March 21st.
 - As usual, start early and ask questions if you get stuck
 - Note: revised version of LL intermediate representation to be more compliant with "real" LLVM IR

Untyped lambda calculus Substitution Evaluation

FIRST-CLASS FUNCTIONS

"Functional" languages

- Languages like ML, Haskell, Scheme, Python, C#, (maybe eventually Java?) include *first-class* functions.
- Functions can be passed as arguments (e.g. map or fold)
- Functions can be returned as values (e.g. compose)
- Functions nest: inner function can refer to variables bound in the outer function

```
let add = fun x -> fun y -> x + y
let inc = add 1
let dec = add -1
let compose = fun f -> fun g -> fun x -> f (g x)
let id = compose inc dec
```

• How do we implement such functions?

Free Variables and Scoping

```
let add = fun x \rightarrow fun y \rightarrow x + y
let inc = add 1
```

- The result of add 1 is a function
- After calling **add**, we can't throw away its argument (or its local variables) because those are needed in the function returned by add.
- We say that the variable x is *free* in fun $y \rightarrow x + y$
 - Free variables are defined in an outer scope
- We say that the variable y is *bound* by "fun y" and its scope is the body "x + y" in the expression fun y -> x + y
- A term with no free variables is called *closed*.
- A term with one or more free variables is called *open*.

(Untyped) Lambda Calculus

- The lambda calculus is a minimal programming language.
 - Note: we're writing (fun x –> e) lambda-calculus notation: $\lambda x. e$
- It has variables, functions, and function application.
 - That's it! (Though for examples, I'll add int and +)
 - It's Turing Complete.
 - It's the foundation for a *lot* of research in programming languages.
 - Basis for "functional" languages like Scheme, ML, Haskell, etc.

Abstract syntax in OCaml:



Concrete syntax.	exp ::=	
	X	variables
	fun x -> exp	functions
	$\exp_1 \exp_2$	function application
CIS 341: Compilers	(exp)	parentheses

Values and Substitution

• The only values of the lambda calculus are (closed) functions:

- To *substitute* a (closed) value \mathbf{v} for some variable \mathbf{x} in an expression \mathbf{e}
 - Replace all *free occurrences* of x in e by v.
 - In OCaml: written subst v x e
 - In Math: written $e\{v/x\}$
- Function application is interpreted by substitution:

(fun x -> fun y -> x + y) 1

- = subst 1 x (fun y \rightarrow x + y)
- = (fun y -> 1 + y)

Lambda Calculus Operational Semantics

• Substitution function (in Math):

$x\{v/x\}$	= V
y{v/x}	= y
$(fun x \rightarrow exp)\{v/x\}$	= (fun x -> exp)
$(fun y \rightarrow exp)\{v/x\}$	$= (\operatorname{fun} y \rightarrow \exp\{v/x\})$
$(e_1 e_2)\{v/x\}$	$= (e_1\{v/x\} e_2\{v/x\})$

(replace the free x by v) (assuming $y \neq x$) (x is bound in exp) (assuming $y \neq x$) (substitute everywhere)

• Examples:

 $x y \{(fun z \rightarrow z)/y\} \Rightarrow x (fun z \rightarrow z)$

 $(fun x \rightarrow x y)\{(fun z \rightarrow z) / y\} \Rightarrow (fun x \rightarrow x (fun z \rightarrow z))$

 $(\operatorname{fun} x \rightarrow x) \{ (\operatorname{fun} z \rightarrow z) / x \} \Rightarrow (\operatorname{fun} x \rightarrow x) / x \text{ is not free!}$

Free Variable Calculation

• An OCaml function to calculate the set of free variables in a lambda expression:



- A lambda expression e is *closed* if free_vars e returns
 VarSet.empty
- In mathematical notation:

 $\begin{aligned} & \text{fv}(x) &= \{x\} \\ & \text{fv}(\text{fun } x \text{ -> } exp) &= \text{fv}(exp) \setminus \{x\} \quad ('x' \text{ is a bound in exp}) \\ & \text{fv}(exp_1 exp_2) &= \text{fv}(exp_1) \cup \text{fv}(exp_2) \end{aligned}$

Operational Semantics

- Specified using just two inference rules with judgments of the form exp ↓ val
 - Read this notation a as "program exp evaluates to value val"
 - This is *call-by-value* semantics: function arguments are evaluated before substitution

$v \Downarrow v$

"Values evaluate to themselves"

$$\exp_1 \Downarrow (\text{fun } x \rightarrow \exp_3) \qquad \exp_2 \Downarrow v \qquad \exp_3\{v/x\} \Downarrow w$$

 $\exp_1 \exp_2 \Downarrow w$

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

Adding Integers to Lambda Calculus

$$exp ::= | ... | n | exp_1 + exp_2 val ::= | fun x -> exp | n | n | n | n | n n x -> exp | n x -> exp x -= exp$$

 $exp_{1} \Downarrow n_{1} exp{2} \Downarrow n_{2}$ $exp_{1} + exp_{2} \Downarrow (n1 [+] n2)$ Object-level '+' Meta-level '+'

Zdancewic CIS 341: Compilers

How to Implement?

- Code in fun.ml shows:
 - A substitution-based interpreter
 - Two environment-based interpreters (one broken)

• We'll come back to compilation after discussing typechecking....

Ruling out ill-defined programs at compile time.

TYPE CHECKING

Zdancewic CIS 341: Compilers

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 + true, ...)
- A compiler can't generate sensible code for this case.
 - A naïve implementation might "add" an integer and a pointer

What to do?

- Don't worry about it... e.g. C, C++
 - Result: segmentation faults, bus errors, etc.
- Make all operations total (i.e. defined everywhere)... e.g. Scheme / Perl
 - 3 + true \rightarrow 42, ... (language specifies behavior)
 - Result: unpredictable answers
- Raise a "runtime type error"... e.g. Python, Ruby, and other dynamically typed languages
 - Result: failure at deployment time
- Try to rule out ill-formed programs... e.g. Java, C#, ML, Haskell
 - 3 + true \rightarrow compiler error:
 - "This expression has type bool but is here used with type int"
 - Result: predictable programs, but it's harder to "get programs running"
- How do you know you've ruled out all ill-formed programs?

Type Soundness

- Build a model of the programming language
 - One model: an interpreter
 - Another model: constructed in mathematics
 - Usually defined via the abstract syntax
- Model defines where the language operations are partial
 - Partiality is different for different languages: e.g. "foo" + "bar" is meaningful in Java but not OCaml
- Construct a function: well_typed : Ast -> unit
 - When well_typed e succeeds, running e will definitely not trigger one of the undefined operation cases (i.e. e is type safe)
 - When well_typed e aborts with an exception, running e might trigger an undefined operation (i.e. e is not type safe)
- *Prove* that the well_typed function is correct.
 - Such proofs are sometimes difficult, but doable for real languages (e.g. SML, Java)

Typechecking

- How do we implement the function well_typed?
- Big idea: "approximate" the interpreter:
 - Problem is partiality in the language semantics as defined by the interpreter.
 - Instead of interpreting the program, write a function called typecheck that computes a type for the program (rather than the answer obtained by running the program).
 - Behavior of typecheck is guided by what the interpreter would do.
- See "tc.ml"

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.)
 - Because of this, we must *assume* the input has some type (say t₁) and reflect this in the type of the function (t₁ -> t₂).
- 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?

Defining Type Systems Mathematically

- In the OCaml implementation we have:
 typecheck (env:environment) (e:exp):ty
 - Where exp is the type of abstract syntax and environment is a list of var * ty pairs.
 - The result of typecheck is a type
- We can abstract this function in math as a relation: The notation: $E \vdash e : t$ means typecheck C = t
 - "In the environment E, program e is well-typed and has type t"
 - "e : t" is a type judgment
- Simple examples: ⊢ 3 : int ⊢ true : bool ⊢ "hello" : string
- Bigger examples: $\vdash (2 * 3) + 5$: int \vdash if (true) 3 else 4: int

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: $x_1 : t_1, x_2 : t_2, ..., x_n : t_n$
- For example: $x : int, b : bool \vdash if (b) 3 else x : int$
- What do we need to know to decide whether "if (b) 3 else x" has type int in the environment x : int, b : bool?
 - b must be a bool i.e. $x : int, b : bool \vdash b : bool$
 - 3 must be an int i.e. $x : int, b : bool \vdash 3 : int$
 - x must be an int i.e. $x : int, b : bool \vdash x : int$

Generalizing 'if' & Inference Rules

• For any environment E, expressions e_1 , e_2 , e_3 , and type T the judgment

 $E \vdash if(e_1) e_2 else e_3 : T$

is true if $E \vdash e_1$: bool, and $E \vdash e_2$: T, and $E \vdash e_3$: T are all true.

• More succinctly: we summarize this as an *inference rule*:

Premises
$$E \vdash e_1 : bool \quad E \vdash e_2 : T \quad E \vdash e_3 : T$$

Conclusion $E \vdash if (e_1) e_2 else e_3 : T$

- This rule holds for *any* substitution of the syntactic metavariables E, e_1 , e_2 , e_3 , and T

Simply-typed Lambda Calculus

• For the language in "tc.ml" we have five inference rules:



• 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 (fun (x:int) -> x + 3) 5 : 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 $x : int \in E$ is implemented by the function **lookup**