

# Homework 1: Abstract Syntax and Rule Definitions

15-814: Types and Programming Languages

Fall 2015

TA: Evan Cavallo (ecavallo@cs.cmu.edu)

Out: 9/10/15

Due: 9/24/15

Notes:

- **Correction 9/15/15:** Fixed the treatment of contexts and added  $\Delta \vdash e \text{val}_e$  requirement to canonical forms lemma.
- Welcome to 15-814's first homework assignment!
- Please submit your work as a PDF file to [ecavallo@cs.cmu.edu](mailto:ecavallo@cs.cmu.edu) and include the phrase "15-814 Homework 1" in the subject line of your email.
- We will be using Piazza to answer questions as well as make clarifications, corrections and announcements. Please sign up for the class on Piazza to make sure you don't miss anything: [piazza.com/cmu/fall2015/15814](http://piazza.com/cmu/fall2015/15814)

In this homework, we will examine a simple C-like language from the perspective of abstract syntax and type theory.

## 1 Arithmetic

We'll start with a simple expression language similar to that covered in class. There are two types, integers `int` and booleans `bool`. The language consists of literals, a few arithmetic and comparison operators, and a conditional expression.

$$\begin{aligned} \tau &::= \text{int} \mid \text{bool} \\ e &::= x \mid n \mid \text{true} \mid \text{false} \mid e + e \mid e * e \mid e \leq e \mid \text{if } e \text{ then } e \text{ else } e \end{aligned}$$

Since this language has no variable-binding constructs, translating the concrete presentation into abstract syntax is quite simple.

concrete	abstract
$n$	<code>num</code> [ $n$ ]
<code>true</code>	<code>true</code>
<code>false</code>	<code>false</code>
$e_1 + e_2$	<code>plus</code> ( $e_1$ ; $e_2$ )
$e_1 * e_2$	<code>times</code> ( $e_1$ ; $e_2$ )
$e_1 \leq e_2$	<code>leq</code> ( $e_1$ ; $e_2$ )
<code>if</code> $e_1$ <code>then</code> $e_2$ <code>else</code> $e_3$	<code>if</code> ( $e_1$ ; $e_2$ ; $e_3$ )

We define the typing judgment  $\Gamma \vdash e :_e \tau$  and operational semantics judgments  $\Delta \vdash e \text{val}_e e'$  and  $\Delta \vdash e \mapsto_e e'$  in Appendix A. Type contexts  $\Gamma$  and evaluation contexts  $\Delta$  are defined with the following grammar:

$$\begin{aligned} \Gamma &::= \cdot \mid \Gamma, x : \tau \\ \Delta &::= \cdot \end{aligned}$$

For the moment, the only evaluation context is the empty context  $\cdot$ , but we will change this as we build on the language.

Before we add anything new, however, let's check some properties to make sure you're comfortable with rule induction. For the proofs in this assignment, you may condense similar cases in arguments by induction, but be clear and rigorous in your reasoning.

**Task 1** Prove the following inversion lemma:

(If Inversion) If  $\Gamma \vdash \text{if}(e_1; e_2; e_3) :_e \tau$ , then  $\Gamma \vdash e_1 :_e \text{bool}$ ,  $\Gamma \vdash e_2 :_e \tau$ , and  $\Gamma \vdash e_3 :_e \tau$ .

This seems immediate, but really follows from the induction principle for the typing judgment. (Tip: Prove that for all  $\Gamma, e, \tau$  such that  $\Gamma \vdash e :_e \tau$ , if  $e = \text{if}(e_1; e_2; e_3)$  for some  $e_1, e_2, e_3$ , then  $\Gamma \vdash e_1 :_e \text{bool}$ ,  $\Gamma \vdash e_2 :_e \tau$ , and  $\Gamma \vdash e_3 :_e \tau$ .)

In general, an inversion lemma is one which recovers the premises of a rule from its conclusion. In the rest of the assignment, you may use these without proof, but be sure to note explicitly when you apply them and check for yourself that they hold.

**Task 2** Prove unicity of typing for this language.

(Unicity of Typing) For any  $\Gamma, e, \tau, \tau'$  such that  $\Gamma \vdash e :_e \tau$  and  $\Gamma \vdash e :_e \tau'$ , we have  $\tau = \tau'$ . You may assume that any variable appears at most once in a given context.

## 2 Functions

### 2.1 Function Blocks

On top of the expression language, we now add a syntax for programs  $p$ , which consist of a series of one-argument function definitions and a final expression computation. (In such a small language, it would be useful to have multi-argument functions, but we will restrict ourselves for the sake of simplicity.) In order to make use of these definitions, we add syntax for function calls to the expression language. The sort of functions consists only of function variables, which we denote here with the letter  $u$ .

$$\begin{aligned} f &::= u \\ e &::= \dots \mid f(e) \\ p &::= \text{fun } u (x : \tau) \{e\}; p \mid \text{result } e \end{aligned}$$

We allow calls to previously defined functions to appear in the definitions of later functions as well as in the result expression. As an example, the program

```
fun not (b : bool) {if b then false else true};
fun abs (x : int) {if not(0 ≤ x) then -1 * x else x};
fun square (z : int) {z * z};
fun cube (z : int) {square(z) * z};
result cube(abs(3)) + cube(abs(-2))
```

should compute to 35.

**Task 3** Give a definition in abstract syntax for the language constructs described above. Include the arities (with sorts) of the operators, per Section 1.2 of PFPL.

Before giving the precise operational semantics of this language, we have an ambiguity to settle. Consider the following program.

```
fun u (x : int) {x};
fun v (y : int) {u(y)};
fun u (z : int) {z + 1};
result v(1)
```

The result expression  $v(1)$  will be computed in an environment of function definitions. Executing the call  $v(1)$  looks up the definition of  $v$ , which then leads to a call to  $u$ . But which of the two functions named  $u$  does this call refer to? Depending on how the environment behaves, it may be the first or the second, and so this program might compute either 1 or 2. The first evaluation strategy, in which the  $u$  that  $v$  calls remains the  $u$  available when  $v$  was defined, is called *static scoping* (sometimes *lexical scoping*). The second, wherein the call to  $u$  in  $v$  refers to whichever  $u$  was most recently added to the environment, is called *dynamic scoping*.

**Task 4** We want to respect the identification convention:

*Abstract binding trees are always identified up to  $\alpha$ -equivalence.*

*With this in mind, which of the two strategies is correct? Explain.*

Now, let's move on to defining the statics and dynamics of this language precisely. We will need two typing judgments to determine if programs are well-typed. Besides the judgment  $\Gamma \vdash e :_e \tau$  for typing expressions we have already mentioned, there should be a judgment  $\Gamma \vdash p :_p \tau$  which states that a program has result of type  $\tau$ . We also need to add a new kind of hypothesis to our type contexts:

$$\Gamma ::= \dots \mid u :_f \tau_1 \Rightarrow \tau_2$$

The assumption  $u :_f \tau_1 \Rightarrow \tau_2$  is used to denote that  $u$  is a function variable in scope which takes an argument of type  $\tau_1$  and returns a result of type  $\tau_2$ .

**Task 5** Starting with the rules in Appendix A, complete the definition of the typing judgment  $\Gamma \vdash e :_e \tau$ , and define the judgement  $\Gamma \vdash p :_p \tau$ .

**Task 6** Give a structural operational semantics for this language by completing the definition of  $\Delta \vdash e \mapsto_e e'$  and defining a judgement  $\Delta \vdash p \mapsto_p p'$ . Assume that the value judgement  $\Delta \vdash p \text{ val}_p$  is defined by the single rule

$$\frac{\Delta \vdash e \text{ val}_e}{\Delta \vdash (\text{result } e) \text{ val}_p} \text{ (RES-V)}$$

*In order to track the environment of function definitions, we will add a new kind of assumption to our evaluation context:*

$$\Delta ::= \dots \mid u \stackrel{\text{def}}{=} x.e$$

*The assumption  $u \stackrel{\text{def}}{=} x.e$  expresses that the function variable  $u$  is currently in scope and that it takes an argument  $x$  and computes the expression  $e$ .*

*Your definition of the semantics should satisfy the properties of progress and preservation.*

(Progress) *If  $\vdash p :_p \tau$ , then either  $\vdash p \text{ val}_p$  or there exists  $p'$  such that  $\vdash p \mapsto_p p'$ .*

(Preservation) *If  $\vdash p :_p \tau$  and  $\vdash p \mapsto_p p'$ , then  $\vdash p' :_p \tau$ .*

**Task 7** Prove progress for the rules you have specified. You will want to state and prove an analogous theorem for expressions first. For programs, it may be helpful to prove the following more general theorem:

*If  $\Delta ::_f \Gamma$  and  $\Gamma \vdash p :_p \tau$ , then either  $\Delta \vdash p \text{ val}_p$  or there exists  $p'$  such that  $\Delta \vdash p \mapsto_p p'$ .*

Here  $\Delta ::_f \Gamma$  is a judgment on contexts defined by the following rules:

$$\frac{}{\cdot ::_f \cdot} \text{ (ENV-NIL)} \quad \frac{\Delta ::_f \Gamma}{(\Delta, f \stackrel{\text{def}}{=} x.e) ::_f (\Gamma, f :_f \tau_1 \Rightarrow \tau_2)} \text{ (ENV-CONS)}$$

This judgment expresses that the environment  $\Delta$  and type context  $\Gamma$  reference the same function variables, i.e. that they agree on what is in scope. This condition is enough to guarantee a canonical forms lemma, which you will need for your proof:

(Canonical Forms Lemma) *Let contexts  $\Delta ::_f \Gamma$  and an expression  $e$  be given. Assume that  $\Delta \vdash e \text{ val}_e$ . If  $\Gamma \vdash e :_e \text{ int}$ , then  $e = \text{num}[n]$  for some  $n \in \mathbb{Z}$ . If  $\Gamma \vdash e :_e \text{ bool}$ , then either  $e = \text{true}$  or  $e = \text{false}$ .*

This lemma enumerates the forms that well-typed values can take (the eponymous canonical forms). You are not required to prove this lemma.

We can make this language a bit more useful by allowing function definitions to refer recursively to the functions they define. We can then write, for example, the factorial function:

```
fun fact (x : int) {if x ≤ 0 then 1 else x * fact(x + (-1))};
result fact(5)
```

**Task 8** *How does this extension change your answers to Tasks 3, 5, and 6? Give the syntax and rules where different. (You should also make a change to the form of the assumption  $u \stackrel{\text{def}}{=} x.e$ .)*

## 2.2 Functions as Expressions

While we can write programs with a predefined list of functions in this language, there is no way of manipulating functions within expressions. For example, we cannot write a function which takes a function as an argument or one which returns a function. One way of solving this is by adding an operator which makes an expression of a function name, in the style of C's function pointers.

$$e := \dots \mid \text{fn}(f) \mid *e(e)$$

We add a corresponding type of functions which has expressions  $\text{fn}(f)$  as its canonical forms.

$$\tau := \dots \mid \tau \Rightarrow \tau$$

Now it is possible to write functions which act on functions. For example, the following program should return 9.

```
fun fact (x : int) {if x ≤ 0 then 1 else x * fact(x + (-1))};
fun onetothree (y : int ⇒ int) {*y(1) + *y(2) + *y(3)};
result onetothree(fn(fact))
```

Naively, we might extend our language to accommodate these new constructs by adding new statics and dynamics rules for the  $\text{fn}(f)$  and  $*e(e)$ . However, simply doing this will break either progress or preservation.

**Task 9** *Which property (progress or preservation) fails? Why? Describe two solutions to the problem, one by changing the rules for the judgments  $\Delta \vdash p \text{ val}_p$  and  $\Delta \vdash p \mapsto_p p'$  and one by changing the rules for the judgment  $\Delta \vdash p :_p \tau$ .*

These issues, as well as the ambiguity discussed in Task 4, arise essentially because we have separated functions and expressions into distinct sorts. Because of this, we cannot substitute function definitions for function names until we evaluate a particular function call. As such, we have to maintain an environment of function definitions to substitute at time-of-call. If we instead have first-class functions, constructible within the expression language, the situation is much simpler.

**Task 10** *Return to the core expression language. We add a new, distinct function type which contains recursive function expressions.*

$$\tau ::= \dots \mid \tau \rightarrow \tau$$

*Define abstract syntax for operators for constructing and applying recursive function terms. (In other words, define a version of  $\lambda$ -abstraction and application which allows for self-reference.) Give statics and dynamics rules for these operators, assuming that we have only the expression language and no outer layer of function definitions.*

### 3 Structs

In this section, we will add C-style structs to our language. As an example, take the following C definition of a type of binary trees, which are intended to be labeled with integers at the leaves. (If you are not familiar with C syntax, you can find simple tutorials online or ask Evan to clarify.)

```
typedef struct tree_data tree;

struct tree_data {
    int leaf_value;
    tree *left;
    tree *right;
};
```

By convention, either both `left` and `right` are null, in which case the tree is a leaf tagged with the value in `leaf_value`, or neither is null, in which case the tree is a node and `leaf_value` is ignored.

To add these sort of data structures to our language, we'll decompose the problem into three parts. First, we need a way to build types with multiple fields. Second, we need types that have a "null" value. Finally, we need types that can refer to themselves recursively.

The first two are fairly easily solved by adding product and sum types to the language.

$$\begin{aligned} \tau & ::= \dots \mid \tau \times \tau \mid \mathbf{unit} \mid \tau + \tau \\ e & ::= \dots \mid \langle e, e \rangle \mid \langle \rangle \mid \mathbf{fst} \ e \mid \mathbf{snd} \ e \mid \mathbf{in}[l](e) \mid \mathbf{in}[r](e) \mid \mathbf{case} \ e \ \{x.e; x.e\} \end{aligned}$$

A canonical element of the product type  $\tau_1 \times \tau_2$  is a pair  $\langle e_1, e_2 \rangle$  where  $e_1 :_e \tau_1$  and  $e_2 :_e \tau_2$ ; we can represent a struct type with multiple fields by chaining products. The unit type `unit` has a single canonical element  $\langle \rangle$ , and can be used to represent an empty struct. We think of `unit` as a nullary product, a tuple with zero fields. The canonical elements of the sum type  $\tau_1 + \tau_2$  are either of the form  $\mathbf{in}[l](e_1)$  with  $e_1 :_e \tau_1$  or  $\mathbf{in}[r](e_2)$  with  $e_2 :_e \tau_2$ . We may therefore represent a possibly-null field of type  $\tau$  with the type  $\tau + \mathbf{unit}$ .

Adding recursively defined types is more complicated, but we can approach it in much the same way that we approached recursive function blocks in the previous section. We add type variables to the type grammar, and type declarations to the program grammar:

$$\begin{aligned}\tau &::= \dots | t \\ p &::= \dots | \text{type } t = \tau; p\end{aligned}$$

In the type declaration above,  $t$  may appear in  $\tau$ . Much as we track an environment of function definitions for evaluation, we can use a context of type definitions for typechecking in the presence of these declarations. We won't go into the details in this homework, as we'll be covering recursively-defined types later in the course.

**Task 11** Give a type declaration in this language which corresponds to the C type `tree` defined above. This definition should match the C struct's behavior exactly.

**Task 12** The C struct definition of `tree` assumes that the programmer follows the conventions described above, and it is possible to construct a `tree` which is well-typed but corresponds to no legal tree. Give an alternate definition of `tree` in our language which statically enforces the conventions.

## A Base Expression Language

### A.1 Statics

$$\begin{aligned}\frac{}{\Gamma, x :_e \tau \vdash x :_e \tau} \text{ (HYP)} \quad & \frac{}{\Gamma \vdash \text{num}[n] :_e \text{int}} \text{ (NUM)} \quad & \frac{}{\Gamma \vdash \text{true} :_e \text{bool}} \text{ (TRUE)} \\ \\ \frac{}{\Gamma \vdash \text{false} :_e \text{bool}} \text{ (FALSE)} \quad & \frac{\Gamma \vdash e_1 :_e \text{int} \quad \Gamma \vdash e_2 :_e \text{int}}{\Gamma \vdash \text{plus}(e_1; e_2) :_e \text{int}} \text{ (PLUS)} \\ \\ \frac{\Gamma \vdash e_1 :_e \text{int} \quad \Gamma \vdash e_2 :_e \text{int}}{\Gamma \vdash \text{times}(e_1; e_2) :_e \text{int}} \text{ (TIMES)} \quad & \frac{\Gamma \vdash e_1 :_e \text{int} \quad \Gamma \vdash e_2 :_e \text{int}}{\Gamma \vdash \text{leq}(e_1; e_2) :_e \text{bool}} \text{ (LEQ)} \\ \\ \frac{\Gamma \vdash e_1 :_e \text{bool} \quad \Gamma \vdash e_2 :_e \tau \quad \Gamma \vdash e_3 :_e \tau}{\Gamma \vdash \text{if}(e_1; e_2; e_3) :_e \tau} \text{ (IF)}\end{aligned}$$

### A.2 Dynamics

$$\begin{aligned}\frac{}{\Delta \vdash \text{num}[n] \text{val}_e} \text{ (NUM-V)} \quad & \frac{}{\Delta \vdash \text{true} \text{val}_e} \text{ (TRUE-V)} \quad & \frac{}{\Delta \vdash \text{false} \text{val}_e} \text{ (FALSE-V)} \\ \\ \frac{\Delta \vdash e_1 \mapsto_e e'_1}{\Delta \vdash \text{plus}(e_1; e_2) \mapsto_e \text{plus}(e'_1; e_2)} \text{ (PLUS-S1)} \quad & \frac{\Delta \vdash e_1 \text{val}_e \quad \Delta \vdash e_2 \mapsto_e e'_2}{\Delta \vdash \text{plus}(e_1; e_2) \mapsto_e \text{plus}(e_1; e'_2)} \text{ (PLUS-S2)} \\ \\ \frac{}{\Delta \vdash \text{plus}(\text{num}[n]; \text{num}[m]) \mapsto_e \text{num}[n + m]} \text{ (PLUS-I)} \quad & \frac{\Delta \vdash e_1 \mapsto_e e'_1}{\Delta \vdash \text{times}(e_1; e_2) \mapsto_e \text{times}(e'_1; e_2)} \text{ (TIMES-S1)} \\ \\ \frac{\Delta \vdash e_1 \text{val}_e \quad \Delta \vdash e_2 \mapsto_e e'_2}{\Delta \vdash \text{times}(e_1; e_2) \mapsto_e \text{times}(e_1; e'_2)} \text{ (TIMES-S2)} \quad & \frac{}{\Delta \vdash \text{times}(\text{num}[n]; \text{num}[m]) \mapsto_e \text{num}[nm]} \text{ (TIMES-I)} \\ \\ \frac{\Delta \vdash e_1 \mapsto_e e'_1}{\Delta \vdash \text{leq}(e_1; e_2) \mapsto_e \text{leq}(e'_1; e_2)} \text{ (LEQ-S1)} \quad & \frac{\Delta \vdash e_1 \text{val}_e \quad \Delta \vdash e_2 \mapsto_e e'_2}{\Delta \vdash \text{leq}(e_1; e_2) \mapsto_e \text{leq}(e_1; e'_2)} \text{ (LEQ-S2)}\end{aligned}$$

$$\frac{n \leq m}{\Delta \vdash \text{leq}(\text{num}[n]; \text{num}[m]) \mapsto_e \text{true}} \text{ (LEQ-I1)} \quad \frac{n > m}{\Delta \vdash \text{leq}(\text{num}[n]; \text{num}[m]) \mapsto_e \text{false}} \text{ (LEQ-I2)}$$

$$\frac{\Delta \vdash e_1 \mapsto_e e'_1}{\Delta \vdash \text{if}(e_1; e_2; e_3) \mapsto_e \text{if}(e'_1; e_2; e_3)} \text{ (IF-S)} \quad \frac{}{\Delta \vdash \text{if}(\text{true}; e_2; e_3) \mapsto_e e_2} \text{ (IF-I1)}$$

$$\frac{}{\Delta \vdash \text{if}(\text{false}; e_2; e_3) \mapsto_e e_3} \text{ (IF-I2)}$$