

# Crash Course on Notation in Programming Language Theory

Jeremy G. Siek

Indiana University, Bloomington

LambdaConf  
HOP Workshop  
June 2018

# Outline

- ▶ Sets, Tuples, Relations, and Functions
- ▶ Language Syntax and Grammars
- ▶ Operational Semantics
- ▶ Type Systems

# Sets

A **set** is a collection of objects.

Examples:

$$\emptyset \qquad \{0, 1, 2\}$$

Order and duplication doesn't matter:

$$\{0, 1, 2\} = \{2, 0, 1\} = \{1, 2, 1, 0, 1\}$$

Set operations:

$$\begin{aligned} 1 &\in \{0, 1, 2\} & 3 &\notin \{0, 1, 2\} \\ \{0, 1\} \cup \{1, 3\} &= \{0, 1, 3\} & \{0, 1\} \cap \{1, 3\} &= \{1\} \\ \{0, 1\} - \{1, 3\} &= \{0\} \end{aligned}$$

Sets can be infinite:

$$\mathbb{N} = \{0, 1, 2, \dots\} \qquad \mathbb{Z} = \{\dots, -2, -1, 0, 1, 2, \dots\}$$

Question: if  $X = \{0\} \cup Y$ , is it true that  $0 \notin Y$ ?

# Tuples

A **tuple** is a sequence of objects.

A **pair** is a tuple of two objects.

Example:

$$(0, 1, 2) \quad \text{or} \quad \langle 0, 1, 2 \rangle$$

Order and duplication matters:

$$(0, 1, 2) \neq (2, 0, 1)$$

$$(0, 1) \neq (0, 0, 1) \neq (0, 0, 1, 1)$$

Subscript with index to access  $n$ th object of the tuple:

$$(a, b, c)_0 = a$$

$$(a, b, c)_1 = b$$

$$(a, b, c)_2 = c$$

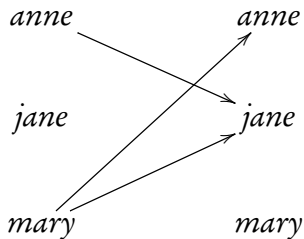
# Relations

A **relation** is a set of pairs.

Example:

$\{(anne, jane), (mary, anne), (mary, jane)\}$

Represents associations between entities:



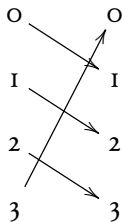
# Functions

A **function** is a relation that associates to each entity at most one other entity.

Example:

$$F = \{(0, 1), (1, 2), (2, 3), (3, 0)\}$$

Represents input/output:



$$F(0) = 1$$

$$F(1) = 2$$

$$F(2) = 3$$

$$F(3) = 0$$

Some functions are infinite:

$$Inc = \{(0, 1), (1, 2), (2, 3), (3, 4), \dots\}$$

# Definition by Rules

We can define an infinite set by a collection of rules.

Example: *Inc* is the set that contains only those elements specified by the following rules:

1.  $(0, 1) \in \textit{Inc}$ .
2. For any  $n$  and  $m$ , **if**  $(n, m) \in \textit{Inc}$ , **then**  $(n + 1, m + 1) \in \textit{Inc}$ .

It's OK for the rules to be recursive, for example,  $(n, m) \in \textit{Inc}$  in rule 2.

# Nonsensical Rules

Some collections of rules are nonsensical.

1.  $(0, 1) \in R$ .
2.  $(0, 1) \notin R$ .
- a. **If**  $(n, m) \in S$ , **then**  $(n, m + 1) \in S$ .
- b. **If**  $(n, m) \in S$ , **then**  $(n + 1, m) \in S$ .



# Nonsensical Rules

Some collections of rules are nonsensical.

1.  $(0, 1) \in R$ .
  2.  $(0, 1) \notin R$ .
- 
- a. **If**  $(n, m) \in S$ , **then**  $(n, m + 1) \in S$ .
  - b. **If**  $(n, m) \in S$ , **then**  $(n + 1, m) \in S$ .

Avoid using negation.

Include at least one non-recursive rule.

# Those Horizontal Lines

Just a notation for if-then rules.  
Premises go on top, conclusion on the bottom.

Recall *Inc*:

1.  $(0, 1) \in Inc$ .
2. For any  $n$  and  $m$ , **if**  $(n, m) \in Inc$ , **then**  
 $(n + 1, m + 1) \in Inc$ .

Definition of *Inc* via horizontal lines:

$$\frac{}{(0, 1) \in Inc} \qquad \frac{(n, m) \in Inc}{(n + 1, m + 1) \in Inc}$$

# Derivations Justify Membership

Recall *Inc*:

$$\text{Rule 1} \frac{}{(0, 1) \in \text{Inc}} \qquad \text{Rule 2} \frac{(n, m) \in \text{Inc}}{(n + 1, m + 1) \in \text{Inc}}$$

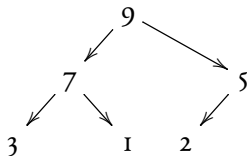
Is  $(2, 3) \in \text{Inc}$ ?

Yes, and here's why:

$$\begin{array}{c} \text{(Rule 1)} \frac{}{(0, 1) \in \text{Inc}} \\ \text{(Rule 2)} \frac{}{(1, 2) \in \text{Inc}} \\ \text{(Rule 2)} \frac{}{(2, 3) \in \text{Inc}} \end{array}$$

# Exercises

- ▶ Define the even natural numbers  $Even = \{0, 2, 4, 6, \dots\}$  using definition-by-rules.
- ▶ Give the derivation for  $4 \in Even$ .
- ▶ Define the set of all binary trees that satisfy the max-heap property (child's label is less-or-equal to parent's) using definition-by-rules. Examples:



# Outline

- ▶ Sets, Tuples, Relations, and Functions
- ▶ **Language Syntax and Grammars**
- ▶ Operational Semantics
- ▶ Type Systems

# Language Syntax and Grammars

Let's define a language of integer arithmetic, call it *Arith*.

A **language** is a set of programs, usually an infinite set.

Example:

$$Arith = \left\{ \begin{array}{l} 3, \\ 7, \\ 3 + 7, \\ -(3 + 7), \\ 4, \\ -(3 + 7) + 4, \\ \vdots \end{array} \right\}$$

# Syntax via Definition by Rules

Definition by rules to the rescue!

*Arith* is the set containing only those programs justified by the following rules:

- ▶ For any  $n \in \mathbb{Z}$ ,  $n \in \text{Arith}$ .
- ▶ For any  $e$ , **if**  $e \in \text{Arith}$ , **then**  $-e \in \text{Arith}$ .
- ▶ For any  $e_1$  and  $e_2$ , **if**  $e_1 \in \text{Arith}$  and  $e_2 \in \text{Arith}$ , **then**  $e_1 + e_2 \in \text{Arith}$ .
- ▶ For any  $e$ , **if**  $e \in \text{Arith}$ , **then**  $(e) \in \text{Arith}$ .

# Backus-Naur Form (BNF)

BNF is a notation for definition-by-rules that is specialized to programming languages.

A collection of BNF rules is called a **grammar**.

Arith ::= integer

Arith ::= "-" Arith

Arith ::= Arith "+" Arith

Arith ::= "(" Arith ")"

or equivalently

Arith ::= integer | "-" Arith | Arith "+" Arith  
          | "(" Arith ")"

derivation = parse tree



# Syntax Conventions in PL Theory

Instead of BNF:

$$\text{Arith} ::= \text{integer} \mid \text{"-"} \text{Arith} \mid \text{Arith} \text{"+"} \text{Arith} \\ \mid \text{"("} \text{Arith} \text{"}")}$$

We select some variables to range over elements of the sets,  
e.g.,  $n \in \mathbb{Z}$ ,  $m \in \mathbb{Z}$  and  $e \in \text{Arith}$ ,

then replace the set names with the variables:

$$e ::= n \mid -e \mid e + e$$

We omit the rule for parentheses; they are always allowed.

# Exercise

Give the PL-theory style syntax definition for a language *ArithPair* that includes integer arithmetic and pairs.

Example programs:

$$-(2 + 3, -4)_o \quad ((1 + 3, (5, 2 + 4))_1)_o$$

# Outline

- ▶ Sets, Tuples, Relations, and Functions
- ▶ Language Syntax and Grammars
- ▶ **Operational Semantics**
- ▶ Type Systems

# Operational Semantics

Define the meaning of programs by  
saying what happens when you run them.

Many styles of operational semantics:

- ▶ big-step semantics
- ▶ small-step semantics

# Big-step Semantics

A relation named *Eval* between programs and result values.

Notation:

$$e \Downarrow n \equiv (e, n) \in Eval$$

We define *Eval* as the set containing only those program-result pairs justified by the following rules.

$$\begin{array}{lll} \text{I} \frac{}{n \Downarrow n} & \text{N} \frac{e \Downarrow n}{-e \Downarrow -n} & \text{P} \frac{e_1 \Downarrow n \quad e_2 \Downarrow m}{e_1 + e_2 \Downarrow n + m} \end{array}$$

Example:

$$-(3 + 7) + 4 \Downarrow -6$$

# Derivation of a big-step

$$\begin{array}{c}
 \begin{array}{c}
 \text{I} \quad \frac{\quad}{3 \Downarrow 3} \quad \text{I} \quad \frac{\quad}{7 \Downarrow 7} \\
 \text{P} \quad \frac{\quad}{3 + 7 \Downarrow 10} \\
 \text{N} \quad \frac{\quad}{-(3 + 7) \Downarrow -10} \\
 \text{P} \quad \frac{\quad}{-(3 + 7) + 4 \Downarrow -6}
 \end{array}
 \quad
 \begin{array}{c}
 \text{I} \quad \frac{\quad}{4 \Downarrow 4}
 \end{array}
 \end{array}$$

# Small-step Semantics

A relation *Step* on programs that does just one computation.

Notation:

$$e \longrightarrow e' \equiv (e, e') \in \textit{Step}$$

We define *Step* as the set containing only those program-program pairs justified by the following rules.

$$\text{N} \frac{}{-n \longrightarrow -n} \quad \text{P} \frac{}{n+m \longrightarrow n+m}$$

$$\text{N}_1 \frac{e \longrightarrow e'}{-e \longrightarrow -e'} \quad \text{P}_1 \frac{e_1 \longrightarrow e'_1}{e_1 + e_2 \longrightarrow e'_1 + e_2} \quad \text{P}_2 \frac{e_2 \longrightarrow e'_2}{n + e_2 \longrightarrow n + e'_2}$$

Example:

$$-(3+7)+4 \longrightarrow -(10)+4 \longrightarrow -10+4 \longrightarrow -6$$

# Derivation of a small-step

$$\begin{array}{c}
 P \quad \text{-----} \\
 N_i \quad \frac{3 + 7 \longrightarrow 10}{-(3 + 7) \longrightarrow -(10)} \\
 P_i \quad \text{-----} \\
 -(3 + 7) + 4 \longrightarrow -(10) + 4
 \end{array}$$



# Exercises

- ▶ Define the set of values for the language *ArithPair*.
- ▶ Extend the big-step semantics to handle the language *ArithPair*.
- ▶ Extend the small-step semantics to handle the language *ArithPair*.

# Outline

- ▶ Sets, Tuples, Relations, and Functions
- ▶ Language Syntax and Grammars
- ▶ Operational Semantics
- ▶ **Type Systems**

# Type Systems

Consider a language *ArithBool* of integers, Booleans, and conditionals:

$$\begin{aligned} e \quad ::= \quad & n \mid -e \mid e + e \\ & \mid \text{true} \mid \text{false} \mid \neg e \mid e \vee e \\ & \mid \text{if } e \text{ then } e \text{ else } e \end{aligned}$$

A **type** classifies values, it is a set of values.

$$\begin{aligned} \text{Int} &= \{ \dots, -2, -1, 0, 1, 2, \dots \} \\ \text{Bool} &= \{ \text{true}, \text{false} \} \end{aligned}$$

- ▶ A type system **predicts** the type of the result value.  
( $-(3 + 7) + 4$ ) produces an Int,  
( $\text{true} \vee \neg \text{false}$ ) produces a Bool.
- ▶ A type system **enforces** that the arguments of an operation make sense.  
( $3 + \text{false}$ ) is ill-typed

# Programs can “go wrong”

Examples:

$$\nexists n. 3 + \text{false} \Downarrow n$$

$$\nexists n. \text{if } 0+0 \text{ then } 1 \text{ else } 3 \Downarrow n$$

In a small-step semantics, the reductions get “stuck”:

$$3 + \text{false} \not\rightarrow$$

$$\text{if } 0+0 \text{ then } 1 \text{ else } 3 \rightarrow \text{if } 0 \text{ then } 1 \text{ else } 3 \not\rightarrow$$

# Type System for *ArithBool*

*WellTyped* is a relation between programs and types.

Notation: let  $T$  range over types (Int and Bool).

$$\vdash e : T \equiv (e, T) \in \textit{WellTyped}$$

Type System: (definition-by-rules yet again!)

$$\begin{array}{c} \frac{}{\vdash n : \text{Int}} \quad \frac{\vdash e : \text{Int}}{\vdash -e : \text{Int}} \quad \frac{\vdash e_1 : \text{Int} \quad \vdash e_2 : \text{Int}}{\vdash e_1 + e_2 : \text{Int}} \\[10pt] \frac{}{\vdash \textit{true} : \text{Bool}} \quad \frac{}{\vdash \textit{false} : \text{Bool}} \quad \frac{\vdash e : \text{Bool}}{\vdash \neg e : \text{Bool}} \\[10pt] \frac{\vdash e_1 : \text{Bool} \quad \vdash e_2 : \text{Bool}}{\vdash e_1 \vee e_2 : \text{Bool}} \quad \frac{\vdash e_1 : \text{Bool} \quad \vdash e_2 : T \quad \vdash e_3 : T}{\vdash \textbf{if } e_1 \textbf{ then } e_2 \textbf{ else } e_3 : T} \end{array}$$

# Type Safety

Well-typed programs cannot “go wrong”.  
—Robin Milner (1978)

Let  $v$  range over values (integers and Booleans).

Theorem (Type Safety)

*If  $\vdash e : T$ , then  $e \Downarrow v$  and  $\vdash v : T$  for some  $v$ .*

# Rule Induction

Suppose set  $S$  is defined by a collection of rules such as

$$\frac{}{a_1 \in S} \quad \frac{a_2 \in S}{f(a_2) \in S} \quad \frac{a_3 \in S \quad a_4 \in S}{g(a_3, a_4) \in S}$$

You want to prove  $\forall x \in S. R(x)$ .

It is sufficient to prove:

- ▶  $R(a_1)$ ,
- ▶ if  $R(a_2)$ , then  $R(f(a_2))$ , and
- ▶ if  $R(a_3)$  and  $R(a_4)$ , then  $R(g(a_3, a_4))$ .

# Proof of Type Safety

By induction on the program  $e$ . The premise is that  $\vdash e : T$ .

- ▶ Case  $\boxed{e = n}$ :  $n \Downarrow n$  and  $\vdash n : \text{Int}$ .
- ▶ Case  $\boxed{e = -e'}$ : By the induction hypothesis,  $e' \Downarrow v$  and  $\vdash v : \text{Int}$ . So  $v$  is an integer. Thus,  $-e' \Downarrow -v$  and  $\vdash -v : \text{Int}$ .
- ▶ Case  $\boxed{e = \neg e'}$ : By the induction hypothesis,  $e' \Downarrow v$  and  $\vdash v : \text{Bool}$ . So  $v$  is *true* or *false*. Thus,  $\neg e' \Downarrow \neg v$  and  $\vdash \neg v : \text{Bool}$ .
- ▶  $\vdots$



# Exercises

- ▶ Devise a type system for the *ArithPair* language.
- ▶ Prove the Type Safety theorem for your type system.

# Suggested Reading

- ▶ Types and Programming Languages  
by Benjamin C. Pierce.
- ▶ Semantic Engineering with PLT Redex  
by Matthias Felleisen, Robert Bruce Findler, and  
Matthew Flatt.
- ▶ My blog: <http://siek.blogspot.com/>
- ▶ Papers in the annual ACM International Conference on  
Functional Programming (ICFP).
- ▶ Papers in the annual ACM Symposium on Principles of  
Programming Languages (POPL).

# Conclusion

- ▶ Infinite sets can be defined via definition-by-rules (aka. inductively defined sets).
- ▶ Use definition-by-rules to define:
  - ▶ language syntax,
  - ▶ operational semantics, and
  - ▶ type systems!
- ▶ Type Safety is the property that all well-typed programs cannot “go wrong” and the result value is of the expected type.



Recall *Inc*:

$$\text{Rule 1} \frac{}{(0, 1) \in \text{Inc}} \qquad \text{Rule 2} \frac{(n, m) \in \text{Inc}}{(n + 1, m + 1) \in \text{Inc}}$$

We capture the notion of applying the rules just once, starting from an arbitrary set  $X$ , with the following function.

$$F(X) = \{(0, 1)\} \cup \{(n + 1, m + 1) \mid (n, m) \in X\}$$

A set  $Y$  is **closed under**  $F$  if applying  $F$  to  $Y$  is already in  $Y$ :

$$F(Y) \subseteq Y$$

The intersection of all such  $Y$ 's is the set *Inc*:

$$\text{Inc} = \bigcap \{Y \mid F(Y) \subseteq Y\}$$

## Theorem (Knaster-Tarski)

*Suppose  $X = \bigcap \{Y \mid F(Y) \subseteq Y\}$  and  $F$  is a monotone function. Then  $X$  is a least fixed point of  $F$ .*

**Proof.** First we show that  $X$  is a fixed point of  $F$ .

$$F(X) = \bigcap \{F(Y) \mid F(Y) \subseteq Y\} \quad (1)$$

$$F(X) \subseteq X \quad (2)$$

$$F(F(X)) \subseteq F(X) \quad F \text{ monotone} \quad (3)$$

$$X \subseteq F(X) \quad \text{def. of } X \quad (4)$$

$$X = F(X) \quad \text{from (2), (4)} \quad (5)$$

Next we show that  $X$  is the least fixed point.

Suppose  $X'$  is another fixed point.

$$X' = F(X')$$

$$F(X') \subseteq X'$$

$$X \subseteq X' \quad \text{def. of } X$$

So by the Knaster-Tarski Theorem,  $Inc$  is the least fixed point of  $F$ , that is,

$$Inc = F(Inc)$$

for any  $X'$ , if  $X' = F(X')$ , then  $Inc \subseteq X'$

# Definitional Interpreter

A recursive procedure that performs the computation.

Example interpreter for *Arith*:

$$eval : Arith \rightarrow \mathbb{Z}$$

$$eval(n) = n$$

$$eval(-e) = -eval(e)$$

$$eval(e_1 + e_2) = eval(e_1) + eval(e_2)$$

Example run:

$$eval(-(3 + 7) + 4) = -6$$



# Definitional interpreter for *ArithBool*

$$eval : ArithBool \rightarrow \mathbb{Z}$$

$$\vdots$$

$$eval(true) = true$$

$$eval(false) = false$$

$$eval(\neg e) = \neg eval(e)$$

$$eval(e_1 \vee e_2) = eval(e_1) \vee eval(e_2)$$

$$eval(\mathbf{if} \ e_1 \ \mathbf{then} \ e_2 \ \mathbf{else} \ e_3) = \begin{cases} eval(e_2) & \text{if } eval(e_1) = true \\ eval(e_3) & \text{if } eval(e_1) = false \end{cases}$$

*eval* is a partial function.