

Liam O'Connor

CSE, UNSW (and data61)

Semester 2 2018

Concrete Syntax

Arithmetic Expressions

$$\begin{array}{c}
 \frac{i \in \mathbb{Z}}{i \text{ Atom}} \\
 \frac{a \text{ SExp}}{(a) \text{ Atom}} \\
 \frac{e \text{ Atom}}{e \text{ PExp}} \\
 \frac{e \text{ PExp}}{e \text{ SExp}} \\
 \frac{a \text{ Atom} \quad b \text{ PExp}}{a \times b \text{ PExp}} \\
 \frac{a \text{ PExp} \quad b \text{ SExp}}{a + b \text{ SExp}}
 \end{array}$$

All the syntax we have seen so far is *concrete syntax*. Concrete syntax is described by judgements on *strings*, which describe the actual text input by the programmer.

Abstract Syntax

Working with concrete syntax directly is *unsuitable* for both compiler implementation and proofs. Consider:

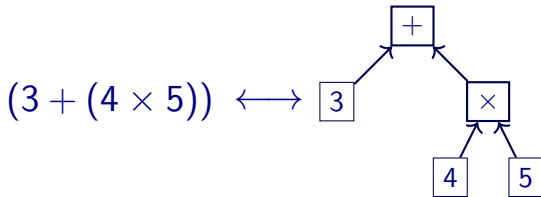
- $3 + (4 \times 5)$
- $3 + 4 \times 5$
- $(3 + (4 \times 5))$

TIMTOWTDI¹ makes life harder for us. Different derivations represent the same semantic program. We would like a representation of programs that is as simple as possible, removing any extraneous information. Such a representation is called *abstract syntax*.

¹“There is more than one way to do it”.

Abstract Syntax

Typically, the *abstract syntax* of a program is represented as a **tree** rather than as a string.



Writing trees in our inference rules would rapidly become unwieldy, however. We shall define a **term** language in which to express trees.

Terms

Definition

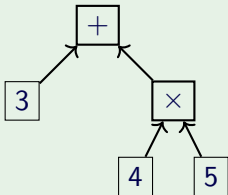
In this course, a *term* is a structure that can either be a *symbol*, like `Plus` or `Times` or `3`; or a *compound*, which consists of an *symbol* followed by one or more *argument subterms*, all in parentheses.

$$t ::= \text{Symbol} \mid (\text{Symbol } t_1 t_2 \dots)$$

These particular terms are also known as *s-expressions*. Terms can equivalently be thought of a subset of `Haskell` where the only kinds of expressions allowed are literals and data constructors.

Term Examples

Example



(Plus (Num 3) (Times (Num 4) (Num 5)))

Armed with an appropriate Haskell data declaration, this can be implemented straightforwardly:

```
data Exp = Plus Exp Exp  
         | Times Exp Exp  
         | Num Int
```

Concrete to Abstract

Concrete Syntax

$$\begin{array}{c}
 \frac{i \in \mathbb{Z}}{i \text{ Atom}} \quad \frac{a \text{ SExp}}{(a) \text{ Atom}} \quad \frac{e \text{ Atom}}{e \text{ PExp}} \quad \frac{e \text{ PExp}}{e \text{ SExp}} \\
 \frac{a \text{ Atom} \quad b \text{ PExp}}{a \times b \text{ PExp}} \quad \frac{a \text{ PExp} \quad b \text{ SExp}}{a + b \text{ SExp}}
 \end{array}$$

Abstract Syntax

$$\frac{i \in \mathbb{Z}}{(\text{Num } i) \text{ AST}} \quad \frac{a \text{ AST} \quad b \text{ AST}}{(\text{Plus } a \ b) \text{ AST}} \quad \frac{a \text{ AST} \quad b \text{ AST}}{(\text{Times } a \ b) \text{ AST}}$$

Now we have to specify a *relation* to connect the two!

Relations

Up until now, most judgements we have used have been *unary* — corresponding to a set of satisfying objects.

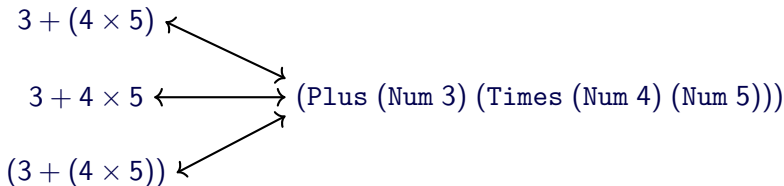
It's also possible for a judgement to express a relationship between *two objects* (a *binary* judgement) or a *number of objects* (an *n-ary* judgement).

Example (Relations)

- 4 *divides* 16 (binary)
- mail *is an anagram of* liam (binary)
- 3 *plus* 5 *equals* 8 (ternary)

n-ary judgements where $n \geq 2$ are sometimes called *relations*, and correspond to an *n*-tuple of satisfying objects.

Parsing Relation



$$i \in \mathbb{Z}$$

$$\frac{i \in \mathbb{Z}}{i \text{ Atom} \longleftrightarrow (\text{Num } i) \text{ AST}}$$

$$\frac{a \text{ Atom} \longleftrightarrow a' \text{ AST} \quad b \text{ PExp} \longleftrightarrow b' \text{ AST}}{a \times b \text{ PExp} \longleftrightarrow (\text{Times } a' b') \text{ AST}}$$

$$a \times b \text{ PExp} \longleftrightarrow (\text{Times } a' b') \text{ AST}$$

$$\frac{a \text{ PExp} \longleftrightarrow a' \text{ AST} \quad b \text{ SExp} \longleftrightarrow b' \text{ AST}}{a + b \text{ SExp} \longleftrightarrow (\text{Plus } a' b') \text{ AST}}$$

$$a + b \text{ SExp} \longleftrightarrow (\text{Plus } a' b') \text{ AST}$$

$$\frac{a \text{ SExp} \longleftrightarrow a \text{ AST}}{(a) \text{ Atom} \longleftrightarrow a \text{ AST}} \quad \frac{e \text{ Atom} \longleftrightarrow a \text{ AST}}{e \text{ PExp} \longleftrightarrow a \text{ AST}} \quad \frac{e \text{ PExp} \longleftrightarrow a \text{ AST}}{e \text{ SExp} \longleftrightarrow a \text{ AST}}$$

$$(a) \text{ Atom} \longleftrightarrow a \text{ AST} \quad e \text{ PExp} \longleftrightarrow a \text{ AST} \quad e \text{ SExp} \longleftrightarrow a \text{ AST}$$

Relations as Algorithms

The *parsing relation* \longleftrightarrow is an extension of our existing concrete syntax rules. Therefore it is **unambiguous**, just as those rules are. Furthermore, the abstract syntax for a particular concrete syntax can be **unambiguously** determined **solely** by looking at the left hand side of \longleftrightarrow .

An Algorithm

To determine the abstract syntax corresponding to a particular concrete syntax:

- 1 Derive the left hand side of the \longleftrightarrow (the concrete syntax) **bottom-up** until reaching axioms.
- 2 Fill in the right hand side of the \longleftrightarrow (the abstract syntax) **top-down**, starting at the axioms.

This process of converting concrete to abstract syntax is called *parsing*.

Example

Rules

$$\begin{array}{c}
 \frac{i \in \mathbb{Z}}{i \mathbf{A} \longleftrightarrow (\text{Num } i)} \quad \frac{a \mathbf{S} \longleftrightarrow a'}{(a) \mathbf{A} \longleftrightarrow a'} \quad \frac{e \mathbf{A} \longleftrightarrow a}{e \mathbf{P} \longleftrightarrow a} \quad \frac{e \mathbf{P} \longleftrightarrow a}{e \mathbf{S} \longleftrightarrow a} \\
 \\
 \frac{a \mathbf{A} \longleftrightarrow a' \quad b \mathbf{P} \longleftrightarrow b'}{a \times b \mathbf{P} \longleftrightarrow (\text{Times } a' b')} \quad \frac{a \mathbf{P} \longleftrightarrow a' \quad b \mathbf{S} \longleftrightarrow b'}{a + b \mathbf{S} \longleftrightarrow (\text{Plus } a' b')}
 \end{array}$$

$$\begin{array}{c}
 \frac{1 \mathbf{A} \longleftrightarrow (\text{Num } 1) \mathbf{AST}}{1 \mathbf{P} \longleftrightarrow (\text{Num } 1) \mathbf{AST}} \quad \frac{\frac{2 \mathbf{A} \longleftrightarrow (\text{Num } 2) \mathbf{AST}}{2 \times 3 \mathbf{P} \longleftrightarrow (\text{Times } (\text{Num } 2) (\text{Num } 3)) \mathbf{AST}}}{2 \times 3 \mathbf{S} \longleftrightarrow (\text{Times } (\text{Num } 2) (\text{Num } 3)) \mathbf{AST}} \quad \frac{\frac{3 \mathbf{A} \longleftrightarrow (\text{Num } 3) \mathbf{AST}}{3 \mathbf{P} \longleftrightarrow (\text{Num } 3) \mathbf{AST}}}{1 + 2 \times 3 \mathbf{S} \longleftrightarrow (\text{Plus } (\text{Num } 1) (\text{Times } (\text{Num } 2) (\text{Num } 3))) \mathbf{AST}}
 \end{array}$$

The Inverse

What about the inverse operation to **parsing**?

Unparsing

Unparsing, also called *pretty-printing*, is the process of starting with the **abstract syntax** on the right hand side of the parsing relation \longleftrightarrow and attempting to synthesise a concrete syntax on the left.

Problem

There are **many** concrete syntaxes for a given abstract syntax. The algorithm is *non-deterministic*.

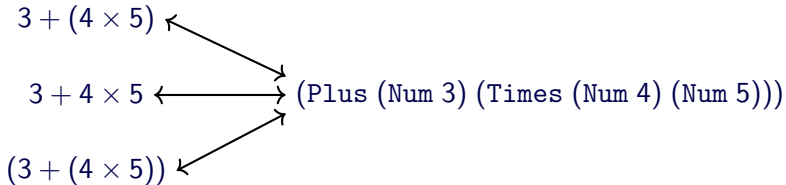
While it is desirable to have:

$$\textit{parse} \circ \textit{unparse} = \textit{id}$$

It is not usually true that:

$$\textit{unparse} \circ \textit{parse} = \textit{id}$$

Example



Going from **right to left** requires some formatting guesswork to produce readable code.

Algorithms to do this can get quite involved!

Let's implement a parser for arithmetic. to coding

Adding Let

Let us extend our arithmetic expression language with **variables**, including a `let` construct to give them values.

Concrete Syntax

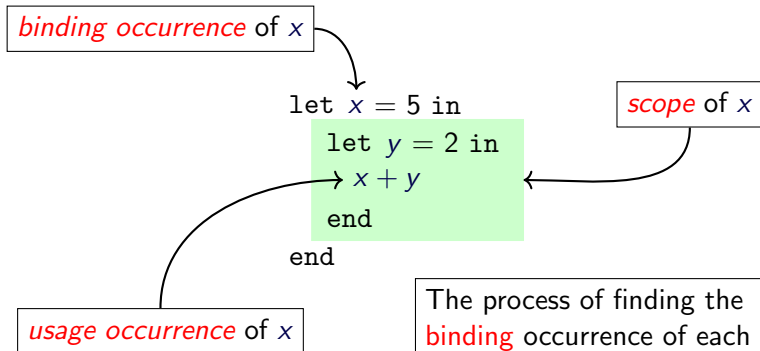
$$\frac{x \text{ Ident}}{x \text{ Atom}} \quad \frac{x \text{ Ident} \quad e_1 \text{ SExp} \quad e_2 \text{ SExp}}{\text{let } x = e_1 \text{ in } e_2 \text{ end Atom}}$$

Example

```
let x = 3 in
  x + 4
end
```

```
let x = 3 in
  let y = 4 in x + y end
end
```

Scope



The process of finding the **binding** occurrence of each used variable is called **scope resolution**. Usually this is done **statically**. If no binding can be found, an **out of scope** error is raised.

Shadowing

What does this program evaluate to?

```
let x = 5 in
```

```
  let x = 2 in
```

```
    x + x
```

```
  end
```

```
end
```



x is *shadowed* here

This program results in 4.

α -equivalence

What is the difference between these two programs?

let $x = 5$ in	let $a = 5$ in
let $x = 2$ in	let $y = 2$ in
$x + x$	$y + y$
end	end
end	end

They are **semantically** identical, but differ in the choice of **bound variable names**. Such expressions are called **α -equivalent**.

We write $e_1 \equiv_{\alpha} e_2$ if e_1 is α -equivalent to e_2 . The relation \equiv_{α} is an **equivalence relation**. That is, it is **reflexive**, **transitive** and **symmetric**.

The process of consistently renaming variables that preserves α -equivalence is called **α -renaming**.

Substitution

A variable x is *free* in an expression e if x occurs in e but is not bound in e .

Example (Free Variables)

The variable x is free in $x + 1$, but not in `let $x = 3$ in $x + 1$ end.`

A *substitution*, written $e[x := t]$ (or $e[t/x]$ in some other courses), is the replacement of all *free* occurrences of x in e with the term t .

Example (Simple Substitution)

$(5 \times x + 7)[x := y \times 4]$ is the same as $(5 \times (y \times 4) + 7)$.

Problems with substitution

Consider these two α -equivalent expressions.

`let y = 5 in y × x + 7 end`

and

`let z = 5 in z × x + 7 end`

What happens if you apply the substitution $[x := y \times 3]$ to both expressions? You get two **non- α -equivalent** expressions!

`let y = 5 in y × (y × 3) + 7 end`

and

`let z = 5 in z × (y × 3) + 7 end`

This problem is called **capture**.

Variable Capture

Capture can occur for a substitution $e[x := t]$ whenever there is a bound variable in the expression e with the same name as a free variable occurring in t .

Fortunately

It is **always possible** to avoid capture.

- **α -rename** the offending bound variable to an unused name, or
- If you have access to the free variable's definition, renaming the free variable, or
- Use a **different abstract syntax representation** that makes capture impossible (More on this next lecture).

To whiteboard Let's define *capture-avoiding substitution* for our **AST** abstract syntax.