

Concepts of Programming Languages

Syntax

Lecturer: Gabriele Keller

Tutor: Liam O'Connor

University of New South Wales

School of Computer Sciences & Engineering

Sydney, Australia

COMP 3161/9161 Week 3

Overview

- So far
 - judgements and inference rules
 - rule induction
 - grammars specified using inference rules
- This week
 - relations and inference rules
 - first-order abstract syntax
 - higher-order abstract syntax
 - substitution

Judgements revisited

- A judgement states that a certain property holds for a specific object (which corresponds to a set membership)
- **More generally**, judgements express **a relationship between a number of objects** (n -ary relations)
- **Examples:**
 - *4 divides 16* (binary relationship)
 - *ail is a substring of mail* (binary)
 - *3 plus 5 equals 8* (tertiary)
- A n -ary relation implicitly defines sets of n -tuples
 - divides: $\{(2, 0), (2, 2), (2, 4), \dots (3, 0), (3, 3), (3, 6), \dots, (4, 0), (4, 4), (4, 8), \dots\}$
 - substring: $\{("", "mail"), ("m", "mail"), ("ma", "mail"), ("ai", "mail"), \dots\}$
 - plus_equal: $\{(0, 0, 0), (0, 1, 1), (0, 2, 2), \dots, (1, 2, 3), (2, 2, 4), (3, 2, 5), \dots\}$

Relations

Definition: A binary relation R is

symmetric, iff for all a, b , aRb implies bRa

reflexive, iff for all a , aRa holds

transitive, iff for all a, b, c , aRb and bRc implies aRc

Definition:

A relation which is symmetric, reflexive, and transitive is called equivalence relation.

Concrete Syntax

$$\begin{array}{c} \frac{e_1 \text{ SExpr} \quad e_2 \text{ PExpr}}{e_1 + e_2 \text{ SExpr}} \qquad \frac{e \text{ PExpr}}{e \text{ SExpr}} \\[10pt] \frac{e_1 \text{ PExpr} \quad e_2 \text{ FExpr}}{e_1 * e_2 \text{ PExpr}} \qquad \frac{e \text{ FExpr}}{e \text{ PExpr}} \\[10pt] \frac{}{i \text{ FExpr}} \quad i \in \text{Int} \qquad \frac{e \text{ SExpr}}{(e) \text{ FExpr}} \end{array}$$

- the inference rules for *SExpr* defined the **concrete syntax** of a simple language, including precedence and associativity
- the concrete syntax of a language is designed with the human user in mind
- not adequate for internal representation during compilation

Concrete vs abstract syntax

- Example:

- ▶ $1 + 2 * 3$

- ▶ $1 + (2 * 3)$

- ▶ $(1) + ((2) * (3))$

- ▶ what is the problem?

- Concrete syntax contains too much information

- ▶ these expressions all have different derivations, but semantically, they represent the same arithmetic expression

- After parsing, we're just interested in three cases: an expression is either

- ▶ an addition

- ▶ a multiplication or

- ▶ a number

Concrete vs abstract syntax

- we use terms of the form

(Operator arg₁ arg₂)

to represent parsed programs unambiguously; e.g.,

`Plus (Num 1) (Times (Num 2) (Num 3))`

- we define the **abstract syntax** of arithmetic expressions as follows:

$$\frac{t_1 \text{ expr} \quad t_2 \text{ expr}}{(Times\ t_1\ t_2)\ \text{expr}}$$
$$\frac{t_1 \text{ expr} \quad t_2 \text{ expr}}{(Plus\ t_1\ t_2)\ \text{expr}}$$
$$\frac{\text{---}}{(Num\ i)\ \text{expr}} \quad i \in Int$$

Concrete vs abstract syntax

- Parsers

- ▶ check if the program (sequence of tokens) is derivable from the rules of the *concrete syntax*
- ▶ turn the derivation into an *abstract syntax tree*

- Transformation rules

- ▶ we formalise this with inference rules as a binary relation \leftrightarrow :

We write

$$e \text{ } SExpr \leftrightarrow t \text{ } expr$$

iff the (concrete syntax) expression e corresponds to the (abstract syntax) expression t .

Usually, many different concrete expressions correspond to a single abstract expression

Concrete vs abstract syntax

- Example:

★ $1 + 2 * 3$ *SExpr* \leftrightarrow Plus (Num 1) (Times (Num 2) (Num 3)) *expr*

★ $1 + (2 * 3)$ *SExpr* \leftrightarrow Plus (Num 1) (Times (Num 2) (Num 3)) *expr*

★ $(1) + ((2)*(3))$ *SExpr* \leftrightarrow Plus (Num 1) (Times (Num 2) (Num 3)) *expr*

Concrete vs abstract syntax

- **Formal definition:** we define a **parsing relation** \leftrightarrow formally as an extension of the structural rules of the concrete syntax.

$$\frac{e_1 \text{ SExpr} \leftrightarrow e_1' \text{ expr} \quad e_2 \text{ PExpr} \leftrightarrow e_2' \text{ expr}}{e_1 + e_2 \text{ SExpr} \leftrightarrow (\text{Plus } e_1' \ e_2') \text{ expr}}$$

$$e_1 + e_2 \text{ SExpr} \leftrightarrow (\text{Plus } e_1' \ e_2') \text{ expr}$$

$$\frac{e_1 \text{ PExpr} \leftrightarrow e_1' \text{ expr} \quad e_2 \text{ FExpr} \leftrightarrow e_2' \text{ expr}}{e_1 * e_2 \text{ PExpr} \leftrightarrow (\text{Times } e_1' \ e_2') \text{ expr}}$$

$$e_1 * e_2 \text{ PExpr} \leftrightarrow (\text{Times } e_1' \ e_2') \text{ expr}$$

$$\frac{}{i \text{ FExpr} \leftrightarrow (\text{Num } i) \text{ expr}} \quad i \in \text{Int}$$

$$\frac{e \text{ PExpr} \leftrightarrow e' \text{ expr}}{e \text{ SExpr} \leftrightarrow e' \text{ expr}}$$

$$e \text{ SExpr} \leftrightarrow e' \text{ expr}$$

$$\frac{e \text{ FExpr}}{e \text{ PExpr}}$$

$$e \text{ PExpr}$$

$$\frac{e \text{ SExpr} \leftrightarrow e' \text{ expr}}{(e) \text{ FExpr} \leftrightarrow e' \text{ expr}}$$

$$(e) \text{ FExpr} \leftrightarrow e' \text{ expr}$$

Concrete vs abstract syntax

- Formal definition: we define a parsing relation \leftrightarrow formally as an extension of the structural rules of the concrete syntax.

$$\frac{e_1 \text{ SExpr} \leftrightarrow e_1' \text{ expr} \quad e_2 \text{ PExpr} \leftrightarrow e_2' \text{ expr}}{e_1 + e_2 \text{ SExpr} \leftrightarrow (\text{Plus } e_1' \ e_2') \text{ expr}}$$

$$e_1 + e_2 \text{ SExpr} \leftrightarrow (\text{Plus } e_1' \ e_2') \text{ expr}$$

$$\frac{e_1 \text{ PExpr} \leftrightarrow e_1' \text{ expr} \quad e_2 \text{ FExpr} \leftrightarrow e_2' \text{ expr}}{e_1 * e_2 \text{ PExpr} \leftrightarrow (\text{Times } e_1' \ e_2') \text{ expr}}$$

$$e_1 * e_2 \text{ PExpr} \leftrightarrow (\text{Times } e_1' \ e_2') \text{ expr}$$

$$\frac{}{i \text{ FExpr} \leftrightarrow (\text{Num } i) \text{ expr}} \quad i \in \text{Int}$$

$$\frac{e \text{ PExpr} \leftrightarrow e' \text{ expr}}{e \text{ SExpr} \leftrightarrow e' \text{ expr}}$$

$$e \text{ SExpr} \leftrightarrow e' \text{ expr}$$

$$\frac{e \text{ FExpr} \leftrightarrow e' \text{ expr}}{e \text{ PExpr} \leftrightarrow e' \text{ expr}}$$

$$e \text{ PExpr} \leftrightarrow e' \text{ expr}$$

$$\frac{e \text{ SExpr} \leftrightarrow e' \text{ expr}}{(e) \text{ FExpr} \leftrightarrow e' \text{ expr}}$$

$$(e) \text{ FExpr} \leftrightarrow e' \text{ expr}$$

The translation relation \leftrightarrow

- The binary **syntax translation relation**

- ▶ $e \leftrightarrow e'$

can be viewed as **translation function**

- ▶ input is e
 - ▶ output is e'
 - ▶ derivations are unambiguously determined by e
 - *since the grammar of the concrete syntax was unambiguous*
 - ▶ e' is unambiguously determined by the derivation
 - *for each concrete syntax term, there is only one rule we can apply at each step*

The translation relation \leftrightarrow

- Derive the abstract syntax as follows:
 - bottom up**, decompose the concrete expression e according to the left hand side of \leftrightarrow
 - top down**, synthesise the abstract expression e' according to the right hand side of each \leftrightarrow from the rules used in the derivation.
- Example:** derivation for $1 + 2 * 3$ (we abbreviate $SExpr$, $PExpr$, $FExpr$ with S , P , F respectively, and $expr$ with e)

$$\begin{array}{c}
 \begin{array}{ccc}
 \text{1 Int} & & \text{2 Int} \\
 \hline
 \text{1 } F \leftrightarrow (\text{Num } 1) e & \text{2 } F \leftrightarrow (\text{Num } 2) e & \text{3 Int} \\
 \hline
 \text{1 } P \leftrightarrow (\text{Num } 1) e & \text{2 } P \leftrightarrow (\text{Num } 2) e & \text{3 } F \leftrightarrow (\text{Num } 3) e \\
 \hline
 \text{1 } S \leftrightarrow (\text{Num } 1) e & \text{2 } * \text{3 } P \leftrightarrow (\text{Times } (\text{Num } 2) (\text{Num } 3)) e & \\
 \hline
 \text{1 + 2 * 3 } S \leftrightarrow \text{Plus } (\text{Num } 1) (\text{Times } (\text{Num } 2) (\text{Num } 3)) e
 \end{array}
 \end{array}$$

Parsing and inference rules

- The parsing problem

Given a sequence of tokens s $SExpr$, find t such that

$$s \text{ } SExpr \leftrightarrow t \text{ } expr$$

- Requirements

A parser should be

- ▶ **total** for all expressions that are correct according to the concrete syntax, that is
 - there must be a $t \text{ } expr$ for every $s \text{ } SExpr$
- ▶ **unambiguous**, that is for every t_1 and t_2 with
 - $s \text{ } SExpr \leftrightarrow t_1 \text{ } expr$ and $s \text{ } SExpr \leftrightarrow t_2 \text{ } expr$
we have $t_1 = t_2$

Parsing and pretty printing

- The parsing problem

Given a sequence of tokens s $SExpr$, find t such that

$$s \text{ } SExpr \leftrightarrow t \text{ } expr$$

- What about the inverse?

- given $t \text{ } expr$, find $s \text{ } SExpr$

- The inverse of parsing is **unparsing**

- ▶ **unparsing** is often ambiguous

- ▶ **unparsing** is often partial (not total)

- Pretty printing

- unparsing together with appropriate formatting is called **pretty printing**

- due to the ambiguity of unparsing, this will usually not reproduce the original program (but a semantically equivalent one)

Parsing and pretty printing

Example

Given the abstract syntax term

```
Times (Num 3) (Times (Num 4) (Num 5))
```

pretty printing may produce the string

*“3 * 4 * 5” or “(3 * 4) * 5”*

- ▶ it's best to choose the most simple, readable representation
- ▶ but usually, this requires extra effort

Bindings

- Local variable bindings (let)

Let's extend our simple expression language with one feature

- ▶ variables and variable bindings

- ▶ `let v = e1 in e2`

- Example:

```
let
  x = 3
in x + 1
```

```
let x = 3
in let y = x + 1
   in x + y
```

- Concrete syntax (adding two new rules):

$$\frac{id \text{ Ident}}{id \text{ FExpr}}$$

$$\frac{e_1 \text{ SExpr} \quad e_2 \text{ SExpr}}{\text{let } id = e_1 \text{ in } e_2 \quad \text{FExpr}}$$

Bindings

- First order abstract syntax:

$$\frac{}{(Num\ i)\ expr} \quad i \in Int$$

$$\frac{t_1\ expr \quad t_2\ expr}{(Times\ t_1\ t_2)\ expr}$$

$$\frac{t_1\ expr \quad t_2\ expr}{(Plus\ t_1\ t_2)\ expr}$$

$$\frac{id\ Ident}{(Var\ id)\ expr}$$

$$\frac{(Var\ id)\ expr \quad t_1\ expr \quad t_2\ expr}{(Let\ id\ t_1\ t_2)\ expr}$$

Bindings

- Scope

- ▶ $\text{let } x = e_1 \text{ in } e_2$ introduces -or binds- the variable x for use within its **scope** e_2
- ▶ we call the occurrence of x in the left-hand side of the binding its **binding occurrence** (or defining occurrence)
- ▶ occurrences of x in e_2 are **usage occurrences**
- ▶ finding the binding occurrence of a variable is called **scope resolution**

- Two types of scope resolution

- ▶ **static scoping**: scoping resolution happens at compile time
- ▶ **dynamic scoping**: resolution happens at run time (discussed later in the course)

Bindings

Example:

```
let
  x = y
in let y = 2
   in x  scope of y  scope of x
```

Out of scope variable: the first occurrence of *y* is ***out of scope***

Bindings

Example:

```
let
  x = 5
in let x = 3
   in x + x
```

Shadowing: the inner binding of `x` is shadowing the outer binding

► *static scoping*

```
const int b = 5;
int foo()
{
    int a = b + 5;
    return a;
}

int bar()
{
    int b = 2;
    return foo();
}

int main()
{
    foo(); // returns 10
    bar(); // returns 10
    return 0;
}
```

► *dynamic scoping:*

```
const int b = 5;
int foo()
{
    int a = b + 5;
    return a;
}

int bar()
{
    int b = 2;
    return foo();
}

int main()
{
    foo(); // returns 10
    bar(); // returns 7
    return 0;
}
```

Bindings

Example:

what is the difference between these two expressions?

```
let
  x = 3
in x + 1
```

```
let
  y = 3
in y + 1
```

α -equivalence:

- ▶ they only differ in the choice of the bound variable names
- ▶ we call them α -equivalent
- ▶ we call the process of consistently changing variable names α -renaming
- ▶ the terminology is due to a conversion rule of the λ -calculus
- ▶ we write $e_1 \equiv_\alpha e_2$ if two expressions are α -equivalent
- ▶ the relation \equiv_α is an equivalence relation

Substitution

- Free variables

- ▶ a free variable is one without a binding occurrence

- ▶ `let x = 1 in x + y` *y is free in this expression*

- **Substitution:** replacing all occurrences of a free variable x in an expression e by another expression e' is called substitution

- **Example:** substituting x with $2 * y$ in

- ▶ `5 * x + 7` yields

- ▶ `5 * (2 * y) + 7`

Substitution

- We have to be careful when applying substitution:

– let $y = 5$ in $y * x + 7$

– let $z = 5$ in $z * x + 7$

α -equivalent

- ▶ substitute x by $2 * y$ in both

– let $y = 5$ in $y * (2 * y) + 7$

– let $z = 5$ in $z * (2 * y) + 7$

not α -equivalent anymore!

- ▶ the free variable y of $2 * y$ is **captured** in the first expression

Substitution

- **Capture-free substitution:** to substitute e' for x in e we require the free variables in e' to be different from the variables in e
- We can always arrange for a substitution to be capture free
 - ▶ use α -renaming of e' (the expression replacing the variable)
 - ▶ change all variable names that occur in e and e'
 - ▶ or use fresh variable names

Higher-order abstract syntax

- A problem with (first-order) abstract syntax

$$\begin{array}{c} \frac{}{(Num\ i)\ expr} \quad i \in Int \\[1em] \frac{t_1\ expr \quad t_2\ expr}{(Times\ t_1\ t_2)\ expr} \qquad \frac{t_1\ expr \quad t_2\ expr}{(Plus\ t_1\ t_2)\ expr} \\[1em] \frac{id\ Ident}{(Var\ id)\ expr} \qquad \frac{(Var\ id)\ expr \quad t_1\ expr \quad t_2\ expr}{(Let\ id\ t_1\ t_2)\ expr} \end{array}$$

- Defining and usage occurrence of variables are treated the same
 - ▶ abstract syntax doesn't differentiate between binding and using occurrence of a variable
 - ▶ it's difficult to identify α -equivalent expressions
 - ▶ variables are just terms, like numbers

Higher-order abstract syntax

- **Higher-order abstract syntax** has variables and abstraction as special constructs
- A term of the form **$x.t$** is called **an abstraction**
- Structure of a higher-order term: a higher-order term can have one of four forms:
 - (1) a constant (e.g., int, string)
 - (2) a variable x
 - (3) (Operator $t_1 \dots t_n$)
 - ▶ Num 4
 - ▶ Plus x (Num 4)
 - (4) $x.t$ (i.e., the variable x is bound in term t)
 - ▶ $x.$ Plus x (Num 1)
 - ▶ $x.y.$ Plus x y

Higher-order abstract syntax

- Higher-order abstract syntax for let-expressions

first-order

$$\frac{id \text{ Ident}}{(Var \ id) \ expr} \quad \frac{var(id) \ expr \ t_1 \ expr \ t_2 \ expr}{(Let \ id \ t_1 \ t_2) \ expr}$$

higher-order

$$\frac{id \text{ Ident}}{id \ expr} \quad \frac{t_1 \ expr \quad t_2 \ expr}{(Let \ t_1 \ id.t_2) \ expr}$$

- Mapping of concrete to higher-order syntax

$$\frac{e_1 \ SExpr \leftrightarrow t_1 \ expr \quad e_2 \ SExpr \leftrightarrow t_2 \ expr}{let \ id = e_1 \ in \ e_2 \ end \ SExpr \leftrightarrow (Let \ t_1 \ id.t_2) \ expr}$$

$$\frac{id \text{ Ident}}{id \ FExpr \leftrightarrow id \ expr}$$

- Example:

$$let \ x = 5 \ in \ x+y \ SExpr \leftrightarrow (Let \ (Num \ 5) \ (x.Plus \ x \ y)) \ expr$$

Substitution

Definition: *A notation for substitution*

We write

$$t[x:=t']$$

to denote a term t where all the occurrences of x have been replaced by the term t' .

- **Example:**

$$(\text{Plus } x \ y) \ [x := (\text{Num } 1)] \ = \ (\text{Plus } (\text{Num } 1) \ y)$$

Definition: *Renaming*

If we replace a variable in the binding and the body of an abstraction, it is called **renaming**, and the resulting term is α -equivalent to the original term:

$$x.t \equiv_{\alpha} y.t \ [x:=y]$$

if y doesn't occur free in t (or $y \notin FV(t)$)

Substitution

- A inductive definition of $FV(t)$:

$$FV(x) = \{x\}$$

$$FV(Op\ t_1 \dots t_n) = FV(t_1) \cup \dots \cup FV(t_n)$$

$$FV(x.t) = FV(t) \setminus \{x\}$$

- Substituting one variable by another:

$$x[x:=y] = y$$

$$z[x:=y] = z, \text{ if } z \neq x$$

$$(Op\ t_1 \dots t_n) [x:=y] = Op\ (t_1 [x:=y]) \dots (t_n [x:=y])$$

$$x.t [x:=y] = x.t$$

$$z.t [x:=y] = z. (t [x:=y]), \text{ if } x \neq z, y \neq z$$

$$y.t [x:=y] = \text{undefined, if } x \neq y$$

Substitution

- Substituting a variable by a term u :

$$x [x:=u] = u$$

$$z [x:=u] = z, \text{ if } z \neq x$$

$$(Op \ t_1 \ \dots \ t_n) [x:=u] = Op \ (t_1 [x:=u]) \ \dots \ (t_n [x:=u])$$

$$x.t [x:=u] = x.t$$

$$z.t [x:=u] = z. (t [x:=u]), \text{ if } x \neq z, z \notin FV(u)$$

$$y.t [x:=u] = \text{undefined}, \text{ if } y \in FV(u)$$

The untyped λ -Calculus

- Introduced in the 1936 by mathematician Alonzo Church
- Very simple, turing complete formalism for computations
- λ -terms:

$$\frac{}{id \text{ } \lambda\text{-term}}$$

$$\frac{t \text{ } \lambda\text{-term}}{\lambda \text{ } id.t \text{ } \lambda\text{-term}}$$

$$\frac{t \text{ } \lambda\text{-term} \quad s \text{ } \lambda\text{-term}}{(t \text{ } s) \text{ } \lambda\text{-term}}$$

where id is an identifier

The untyped λ -Calculus

- The calculus has three rules:

- ▶ α -conversion

- if $t \equiv_{\alpha} s$, then the two terms are equivalent in the calculus

- ▶ β -reduction

- $(\lambda x. t)s$ can be reduced to $t[x := s]$

- ▶ η -conversion

- $\lambda x. (f\ x)$ is equivalent to f if x is not free in f

The untyped λ -Calculus

- The simple untyped λ -calculus has no constants, conditional, built-in operations
- Natural numbers, arithmetic, logic operations can all be encoded in terms of the simple λ -calculus
 - ▶ $0 := \lambda f . \lambda x . x$
 - ▶ $1 := \lambda f . \lambda x . (f \ x)$
 - ▶ $2 := \lambda f . \lambda x . (f \ (f \ x))$
 - ▶ $3 := \lambda f . \lambda x . (f \ (f \ (f \ x)))$
 - ▶ successor function (+1): $\lambda n . \lambda f . \lambda x . (f \ ((n \ f) \ x))$

The untyped λ -Calculus

- Boolean values

- True: $\lambda x.\lambda y.x$

- False: $\lambda x.\lambda y.y$

- If: $\lambda c.\lambda t.\lambda e.c\ t\ e$

Example: Boolean Expressions

- Consider a language of boolean expressions:

$$bexpr ::= True \mid False \mid \neg bexpr \mid (bexpr) \mid$$
$$bexpr \wedge bexpr \mid bexpr \vee bexpr$$

- Assume the usual precedence and associativity rules apply
- What would the rules for the concrete syntax of this language look like?
- What are the rules for abstract syntax of the language?
- What happens if we introduce variables?