510: Programming Languages Product and Sum Types

David Walker Fall, 2002

Overview

The ML datatype mechanism combines

- sum, or disjoint union, types;
- recursive types;
- abstract types

into a single mechanism.

Overview

Datatype values are built using **constructors**.

- *e.g.*, 3::nil.
- *e.g.*, node(empty, 1, empty)

Datatype values are decomposed using **pattern matching**.

fun depth (node (t1, $_-$, t2)) = 1 + max(depth t1, depth t2)

Overview

To analyze these features of ML, we'll start with these types:

- **Product**, or **tuple**, types.
- Sum, or disjoint union, types.

Then we'll add recursive and, later, abstract types.

Product Types

Product, or **tuple**, types give you structured data.

- Nullary products: unit. Sole value is ().
- Binary products: $\tau_1 * \tau_2$. Values are ordered pairs.
- *n*-ary products: $\tau_1 * \cdots * \tau_n$. Values are ordered *n*-tuples.
- Labelled products, or records: {name:string, salary:float}. Elements are labelled tuples.

We'll formalize binary and nullary products.

Product Types: Abstract Syntax

Adding product types to MinML is easy.

Types τ ::= unit | $\tau_1 * \tau_2$ Exp's e ::= () | check e_1 is () in e_2 end | (e_1, e_2) | split e_1 as (x, y) in e_2 end Values v ::= () | (v_1, v_2)

The variables x and y are bound within e_2 in the expression split e_1 as (x, y) in e_2 end.

Product Types: Static Semantics

$\overline{\Gamma \vdash () : unit}$

 $\frac{\Gamma \vdash e_1 : \texttt{unit} \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \texttt{check} \, e_1 \, \texttt{is} \, \texttt{()} \, \texttt{in} \, e_2 \, \texttt{end} : \tau_2}$

 $\frac{\Gamma \vdash e_1 : \tau_1 \quad \Gamma \vdash e_2 : \tau_2}{\Gamma \vdash (e_1, e_2) : \tau_1 * \tau_2}$

 $\frac{\Gamma \vdash e_1 : \tau_1 * \tau_2 \quad \Gamma, x : \tau_1, y : \tau_2 \vdash e_2 : \tau}{\Gamma \vdash \texttt{split} \, e_1 \texttt{ as } (x, y) \texttt{ in } e_2 \texttt{ end } : \tau}$

Product Types: Dynamic Semantics

 $\overline{\texttt{check}\,(\texttt{)}\,\texttt{is}\,(\texttt{)}\,\texttt{in}\,e\,\texttt{end}\mapsto e}$

$$e_1\mapsto e_1'$$

check e_1 is () in e_2 end \mapsto check e_1' is () in e_2 end

Product Types: Dynamic Semantics

$$\frac{e_1 \mapsto e_1'}{(e_1, e_2) \mapsto (e_1', e_2)}$$

$$\frac{e_2 \mapsto e_2'}{(v_1, e_2) \mapsto (v_1, e_2')}$$

 $\overline{\texttt{split}(v_1,v_2)\texttt{ as }(x,y)\texttt{ in }e\texttt{ end }\mapsto \{v_1,v_2/x,y\}e}$

$$\frac{e_1 \mapsto e'_1}{\operatorname{split} e_1 \operatorname{as} (x, y) \operatorname{in} e_2 \operatorname{end} \mapsto \operatorname{split} e'_1 \operatorname{as} (x, y) \operatorname{in} e_2 \operatorname{end}}$$

Product Types: Example

ML code:

fun ifact (0, a) = a
 | ifact (n, a) = ifact (n-1, n*a)

MinML code:

fun ifact (p:int*int) is
split p as (n, a) in
if n=0 then a else ifact (-(n,1), *(n, a)) f

Product Types: Example

The split construct provides a **single layer** of pattern matching.

- No nested tuples.
- No possibility of failure.

Product Types: Safety

Preservation:

- By induction on evaluation.
- Using substitution lemma for **split**.

Progress:

- Canonical forms of product type are pairs.
- Can always split a pair of the right type.

Sum Types

Sum, or disjoint union, types give you choices.

- Nullary: void, with **no** elements.
- Binary: τ₁+τ₂. Values are either a value of type τ₁ tagged in1, or a value of type τ₂ tagged inr.
- *n*-ary: $\tau_1 + \cdots + \tau_n$.
- Labelled: [present:string, absent:unit].

We'll consider nullary and binary sums.

Sum Types: Abstract Syntax

Types
$$\tau$$
 := void | $\tau_1 + \tau_2$

$$\begin{aligned} Exp's \quad e \quad &::= \quad \operatorname{inl}_{\tau_1 + \tau_2}(e_1) \mid \operatorname{inr}_{\tau_1 + \tau_2}(e_2) \mid \\ & \operatorname{case}_{\tau} e_0 \text{ of inl}(x : \tau_1) \Rightarrow e_1 \mid \operatorname{inr}(y : \tau_2) \Rightarrow e_2 \text{ end} \end{aligned}$$

$$Val's v ::= inl_{\tau_1+\tau_2}(v_1) | inr_{\tau_1+\tau_2}(v_2)$$

In the expression

 $\operatorname{case}_{\tau} e_0 \operatorname{of} \operatorname{inl}(x : \tau_1) \Longrightarrow e_1 | \operatorname{inr}(y : \tau_2) \Longrightarrow e_2 \operatorname{end},$

the variable x is bound in e_1 and the variable y is bound in e_2 .

Sums: Informal Description

The type $\tau_1 + \tau_2$ is the **disjoint union** of τ_1 and τ_2 .

- Values of each type τ_1 and τ_2 are included within it.
- Elements are **tagged** with inl or inr to indicate where they came from.

Thus int+int is quite different from int!

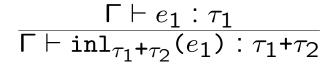
- Elements are inl(n) and inr(n).
- Disjoint union is different from ordinary set union!

Sums: Informal Description

The case construct provides non-nested, exhaustive pattern matching over a sum type:

```
case e:int+int
  of inl(x:int) => +(x,1)
   | inr(y:int) => -(y,1)
```

Sums: Static Semantics



$$\frac{\Gamma \vdash e_2 : \tau_2}{\Gamma \vdash \operatorname{inr}_{\tau_1 + \tau_2}(e_2) : \tau_1 + \tau_2}$$

 $\frac{\Gamma \vdash e_0 : \tau_1 + \tau_2 \quad \Gamma, x_1 : \tau_1 \vdash e_1 : \tau \quad \Gamma, x_2 : \tau_2 \vdash e_2 : \tau}{\Gamma \vdash \mathsf{case}_\tau e_0 \text{ of } \mathsf{inl}(x_1 : \tau_1) \mathrel{=>} e_1 \mid \mathsf{inr}(x_2 : \tau_2) \mathrel{=>} e_2 \text{ end}}$ $: \tau$

Sums: Dynamic Semantics

$$\frac{e \mapsto e'}{\operatorname{inl}_{\tau_1 + \tau_2}(e) \mapsto \operatorname{inl}_{\tau_1 + \tau_2}(e')}$$

$$\frac{e \mapsto e'}{\operatorname{inr}_{\tau_1 + \tau_2}(e) \mapsto \operatorname{inr}_{\tau_1 + \tau_2}(e')}$$

 $\operatorname{case}_{\tau} \operatorname{inl}_{\tau_1 + \tau_2}(v) \text{ of } \operatorname{inl}(x_1 : \tau_1) \Longrightarrow e_1 \mid \operatorname{inr}(x_2 : \tau_2) \Longrightarrow e_2 \text{ end}$ $\mapsto \{v/x_1\}e_1$

 $\operatorname{case}_{\tau} \operatorname{inr}_{\tau_1 + \tau_2}(v) \text{ of } \operatorname{inl}(x_1 : \tau_1) \Longrightarrow e_1 \mid \operatorname{inr}(x_2 : \tau_2) \Longrightarrow e_2 \text{ end} \\ \mapsto \{v/x_2\}e_2$

Booleans are **definable** from sums!

- bool = unit+unit.
- true = inl(()), false = inr(()).
- if e then e_1 else e_2 fi = case e of inl(x_1 :unit) => e_1 | inr(x_2 :unit) => e_2 end.

In fact any **non-recursive** data type is similarly definable.

datatype T = A | B | C of int

- T = unit+(unit+int).
- A = inl(()).
- B = inr(inl(())).
- C(n) = inr(inr(n)).

Pattern matching corresponds to case analysis:

case e
of A => a
| B => b
| C(z) => c

Corresponding MinML code:

```
case e
of inl(w:unit) => a
| inr(x:unit+int) =>
case x
of inl(y:unit) => b
| inr(z:int) => c
```

Sums: Safety

Preservation: by induction on evaluation.

Progress: by induction on typing.

- Canonical forms of type $\tau_1 + \tau_2$: $\operatorname{inl}_{\tau_1 + \tau_2}(v_1)$ or $\operatorname{inr}_{\tau_1 + \tau_2}(v_2)$.
- Proof by induction on typing.

The exhaustiveness of case is crucial for progress!

Unit and Void

The type unit has **one** element, (). The type void has **no** elements! Consequently,

- If a function has type int→void, it must not terminate for any argument.
- If a function has type int→unit, it might return, but the result has to be ().

(Some languages use void when they mean unit)

Many languages have a so-called **null pointer** or **null object**.

- The value null in Java.
- The cast (T *)0 in C.

The "null pointer" is used to model the **ab**-**sence** of a value.

- Often as a default initial value for variables.
- As a "base case" for complex data structures.

The null pointer is a standard source of bugs.

- Null pointer exception in Java.
- Bus error in C.

Standard languages have no ability to track whether a pointer is null.

- Must check for null on each access.
- Explicit null checks do not change the type.

But these problems never arise in ML! Why?

- Absence of "pointer mentality" valueoriented programming.
- Without pointers there are no null pointers!

Why are there no null pointers in ML?

- Sum types obviate the need for them!
- SML: datatype 'a option = NONE | SOME of 'a

In ML there is a **type distinction** between

- A genuine value of type τ , and
- An **optional** value of type τ option.

The key to this is the presence of **sum types**.

- Case analysis **changes the type** from τ option to τ .
- The type system tracks whether a value is present or not! There is no need for a NONE check!

Skeletal ML code for working with options:

fun dispatch (x :
$$\tau$$
 option) =
case x
of NONE => e_0
| SOME (x' : τ) => e_1

Within e_1 the variable x' is **known** not to be "null"!

Skeletal Java code for working with null pointers:

if (x == null)

$$s_1$$

else
 s_2

Within s_2 the type of x is still Object and might still be null!

A harder case:

if (MyMethod(x)) s_1 else s_2

The compiler cannot (in general) track that MyMethod returning false implies that x is non-null!

Summary

Products support structured data.

• Similar to **struct**'s in C, but with automatic allocation and no "pointers".

Sums support alternative data.

- Choice of two distinguishable alternatives.
- Case analysis propagates type change.