Programming Languages and Techniques (CIS1200)

Lecture 5

**Datatypes and Trees** 

#### CIS 1200 Announcements

- HW01 is due **TOMORROW** at midnight
  - Mandatory in-person check-in (15-20 minutes) with your recitation TAs after you submit
  - look for them to coordinate
- HW02 will be released soon
- If you would prefer a less-crowded recitation you are welcome to switch
  - Section 208, 217, 218, and 222 have fewer students
  - Please change directly on Path@Penn before the add deadline tomorrow
  - If you have issues with registration, please mail <u>cis1200@seas.upenn.edu</u>

#### How to get help at Office Hours

1. Go to office hours location (schedule available on <u>course website</u>)

#### TA Office Hours Schedule

Monday	Tuesday	Wednesday	Thursday	Friday	Saturday	Sunday
5-10pm	2-4pm (Virtual OHQ), 5-10pm	6-8pm	7-10pm	None	4-6pm	2-6pm
Towne 217	Towne 303 (5-7pm), Towne 217 (7-10pm)	Towne 217	Towne 217		Towne 217	Towne 217

- 2. Enter the office hours queue via OHQ
  - a. Can be accessed via the course website, you may need to click "add course" and search "CIS 1200"
- 3. Fill out the template

Ask a Question	
Question *	e (examples include retracing logic with examples, drawing a
picture, using the de under etc.)	
Describe Yourself	Make sure to mark which
Beside the window, wearing a red hoodie	question you're struggling with and clearly explain how you have tried to solve the
	problem on your own before you will be helped!

4. Wait to be seen

## **Reminder: No Laptops during Lecture**

- Laptops *closed*... minds *open*
  - Although this is a computer science class, the use of electronic devices – laptops, phones, etc., during lecture (*except for participating in quizzes*) is *prohibited*
- Why?
  - Device users tend to surf/chat/email/ game/text/tweet/etc.
  - They also distract those around them
  - More effective to take notes by hand
  - You'll get plenty of time in front of your computer while working on the homework :-)



#### Recap: Lists, Recursion, & Tuples

## A General Pattern: Structural Recursion Over Lists

Structural recursion builds an answer from smaller components:

The branch for [] calculates the value (f []) directly.

- this is the *base case* of the recursion

The branch for hd::rest calculates (f (hd::rest)) given hd and (f rest).

- this is the *inductive case* of the recursion

#### 5: Given the definition below, which of the following is correct?



**(()** 0

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app



```
f [1; 2] [3;4]
⇒ 1 :: (f [2] [3;4])
⇒ 1 :: 2 :: (f [] [3;4])
⇒ 1 :: 2 :: [3;4]
= [1;2;3;4]
```

#### 5: What is the type of this expression?

int	
	0%
int list	
	0%
int list list	
	0%
(int * int list) list	
	0%
int * (int list) * (int list list)	
	0%
(int * int * int) list	
	0%
none (expression is ill typed)	
	0%

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

#### What is the type of this expression?

#### (1, [1], [[1]])

- 1. int list
- 2. int list list
- 3. (int \* int list) list
- 4. int \* (int list) \* (int list list)
- 5. (int \* int \* int) list
- 6. none (expression is ill typed)

#### 5: What is the type of this expression?





Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

What is the type of this expression?

[ (1,true); (0, false) ]

- 1. int \* bool
- 2. int list \* bool list
- 3. (int \* bool) list
- 4. (int \* bool) list list
- 5. none (expression is ill typed)

Answer: 3

## **Topics for Today**

## Types for Structured Data

- Like most programming languages, OCaml offers a variety of ways of creating and manipulating structured data
- We have already seen:
  - *primitive datatypes* (int, string, bool, ... )
  - lists (int list, string list, string list list, ...)
  - tuples (int \* int, int \* string, …)
- Today:
  - type abbreviations
  - user-defined datatypes

## **Type Abbreviations**

## A Handy Feature: Type Abbreviations

OCaml lets us *name* (i.e., give a synonym for) an existing type



- A type abbreviation is **interchangeable** with its definition
- Abbreviations are useful for naming important concepts

let profit (attendees:int) : money = ...

#### **Datatypes and Trees**

## HW 2 Case Study: Evolutionary Trees

- Problem: reconstruct evolutionary trees\* from DNA data.
  - What are the relevant abstractions?
  - How can we use the language features to define them?
  - How do the abstractions help shape the program?



\*Interested? Check this out: Dawkins: The Ancestor's Tale: A Pilgrimage to the Dawn of Evolution

#### **DNA Computing Abstractions**

- Nucleotide
  - Adenine (A), Guanine (G), Thymine (T), or Cytosine (C)
- Helix
  - a sequence of nucleotides: e.g. AGTCCGATTACAGAGA...
  - genetic code for a particular species (human, gorilla, etc)
- Phylogenetic tree
  - Binary tree with helices (species) at the nodes and leaves



## Simple User-Defined Datatypes

OCaml lets programmers define *new* datatypes

type day =
 Sunday
 Monday
 Monday
 Tuesday
 Wednesday
 Thursday
 Friday
 Saturday



(*must* be capitalized)

The *constructors* are the values of the datatype

- e.g. A : nucleotide
 [A; G; C] : nucleotide list

## Pattern Matching on Simple Datatypes

Datatype values can be analyzed by pattern matching:

let string\_of\_n (n:nucleotide) : string =
 begin match n with
 I A -> "adenine"
 I C -> "cytosine"
 I G -> "guanine"
 I T -> "thymine"
 end

• One case per constructor

- you will get a warning if you leave out a case or list one twice

• As with lists, the pattern syntax follows that of the datatype values (*i.e.*, the constructors)

#### A Point About Abstraction

- We *could* represent weekdays by using integers:
  - Sunday = 0, Monday = 1, Tuesday = 2, *etc.*
- But...!
  - Integers support different operations than days do:
     Wednesday Monday = Tuesday (?!?)
     Wednesday \* Tuesday = Saturday (?!?)
  - There are *more* integers than days (What day is 17? -3?)
- Confusing integers with days can lead to bugs
  - Many "scripting" languages (PHP, Javascript, Perl, Python,...) do confuse values of different types (true == 1 == "1"), leading to much misery when debugging...
- For these reasons, most modern languages (Java, C#, C++, Rust, Swift,...) provide *user-defined types*

## Type Abbreviations II

Abbreviations let us give *names* to complex types but do not introduce new abstractions



- *i.e.*, a helix is the same as a list of nucleotides
   let x : helix = [A;C;C] in length x
- Can make code easier to read & write, but does not provide all the benefits of user-defined types

#### **Data-Carrying Constructors**

• Datatype constructors can also carry values



 Values of type 'measurement' include: Missing NucCount (A, 3) CodonCount ((A,G,T), 17)

#### Pattern Matching on Datatypes

Pattern matching notation combines syntax of tuples and simple datatype constructors:

- Defining a datatype also defines its patterns
- Datatype patterns *bind* identifiers (*e.g.*, 'n') just like for lists and tuples

#### 5: What is the type of this expression?



nucleotide	
	<b>0</b> %
nucleotide list	
	0%
helix	
	0%
nucleotide * string	
	0%
string * string	
	0%
none (expression is ill typed)	
	0%

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

## type nucleotide = | A | C | G | T type helix = nucleotide list

What is the type of this expression?

(A, "A")

- 1. nucleotide
- 2. nucleotide list
- 3. helix
- 4. nucleotide \* string
- 5. string \* string
- 6. none (expression is ill typed)



#### 5. What is the type of the expression [A;C]?

nucleotide	
	0%
helix	
	0%
nucleotide list	
	0%
string * string	
	0%
nucleotide * nucleotide	
	0%
none (expression is ill typed)	
	0%

Start the presentation to see live content. For screen share software, share the entire screen. Get help at pollev.com/app

## type nucleotide = | A | C | G | T type helix = nucleotide list

What is the type of this expression?

[A;C]

- 1. nucleotide
- 2. helix
- 3. nucleotide list
- 4. string \* string
- 5. nucleotide \* nucleotide
- 6. none (expression is ill typed)

Answer: both 2 and 3

Defining a *datatype* adds a fresh abstraction as a firstclass concept to your program.

- *Constructors* explain the shape/structure of the values.
- *Patterns* explain how to inspect/name the components of those values.
- *Abstraction* means that the type can't be confused with other, existing types.

#### Trees

- We now know how to define types for nucleotides, codons, DNA helices, *etc.*
- What about the evolutionary tree itself?



## **Recursive User-defined Datatypes**

Datatype definitions can mention themselves *recursively*:



#### **Tree Values**

```
type labeled_tree =
    LLeaf of helix
    LNode of labeled_tree * helix * labeled_tree
```

Example values of type tree:







## Trees are everywhere

#### Family trees



## Organizational charts



#### Filesystem Folder Structure



#### **Domain Name Hierarchy**



#### Game trees



#### Natural-Language Parse Trees



### COVID evolutionary tree



The Economist

#### **Binary Trees**

A particular form of tree-structured data

#### **Binary Trees**



A binary tree is either *empty*, or a *node* with at most two children, both of which are also binary trees.

A *leaf* is a node whose children are both empty.

## Trees are Drawn Upside Down



#### **Another Tree**



CIS1200

#### **Binary Trees in OCaml**



#### **Representing trees**



#### Working with binary trees

see tree.ml treeExamples.ml

#### Structural Recursion Over Trees

Structural recursion builds an answer from smaller components:

The branch for Empty calculates the value (f Empty) directly.

- this is the *base case* of the recursion

The branch for Node(1,x,r) calculates

- (f (Node(l,x,r)) given X and (f l) and (f r).
- this is the *inductive case* of the recursion

#### Tree vs. List Recursion

#### Trees as Containers

- Like lists, trees aggregate ordered data
- As we did for lists, we can write a function to determine whether a tree *contains* a particular element...

#### Searching for Data in a Tree

- This function searches through the tree, looking for n
- In the worst case, it might have to traverse the *entire tree*

# Search during (contains t 8)



#### Searching for Data in a Tree



5 = 7

II contains (Node(Node (Empty, 0, Empty), 1, Node(Empty, 3, Empty))) 7
II contains (Node (Empty, 7, Empty)) 7

<pre>((0 = 7    contains Empty 7    contains Empty 7)</pre>	
<pre>Il contains (Node (Empty, 7, Empty)) 7</pre>	Eliding some steps
<pre>contains (Node(Empty, 3, Empty)) 7</pre>	
<pre>   contains (Node (Empty, 7, Empty)) 7</pre>	Eliding some steps
contains (Node (Empty, 7, Empty)) 7	