

Data Structures II

Advanced Functional Programming

Andres Löh (andres@cs.uu.nl)

Universiteit Utrecht

24 May 2005



Overview

Introduction (library overview)

Views

Binary search trees

Discussion



Overview

Introduction (library overview)

Views

Binary search trees

Discussion



Excerpt from the Haskell hierarchical libraries I

Data.Array	standard immutable arrays
Data.Array.*	mutable and unboxed arrays
Data.Bits	class for bit operations
Data.Bool	standard bool type and logical operations
Data.Char	standard characters and character classes
Data.Complex	standard complex numbers
Data.Dynamic	dynamic types
Data.Either	standard binary sum type
Data.FiniteMap	deprecated, see Data.Map
Data.FunctorM	monadic functor class
Data.Generics	“scrap your boilerplate” combinators
Data.Graph	easy-to-use graph library
Data.Graph.Inductive	“functional graph library”



Excerpt from the Haskell hierarchical libraries II

Data.HashTable	ephemeral hash table implementation (IO)
Data.Int	standard integers
Data.IntMap	efficient finite maps with integers keys
Data.IntSet	efficient sets of integers
Data.IRef	mutable variables (IO)
Data.Ix	class of array index types
Data.Lists	standard lists
Data.Map	“generic” finite maps
Data.Maybe	standard option/exception type
Data.Monoid	monoid class
Data.PackedString	packed (space-efficient) strings



Excerpt from the Haskell hierarchical libraries III

Data.Queue	single-ended queues
Data.Ratio	standard rational numbers
Data.Set	“generic” sets
Data.STRef	mutable variables (ST)
Data.Tree	(limited) tree operations
Data.Tuple	standard tuples
Data.Typeable	class for dynamic type information
Data.Unique	unique identities (IO)
Data.Version	very simplistic version numbers
Data.Word	words of different bit-length



Overview

Introduction (library overview)

Views

Binary search trees

Discussion



Using data types

- ▶ High-level recursion operators.
- ▶ Special syntax.
- ▶ Pattern matching.



How to use pattern matching on dequeues?

Deque implementation:

```
| data Deque a = D ! Int [a] ! Int [a]
```

Other possible implementation:

```
| data Deque a = D [a] [a] [a]
```

Pattern matching on the implementation type is bad, because

- ▶ it breaks the abstraction
- ▶ the implementation might respect invariants that are not obvious from the representation type



Using functions to destruct deques is tedious

Example: remove the first and the last element of a deque, return their sum and the resulting deque.

```
removefl :: Deque Int → (Int, Deque Int)  
removefl q = (head q + last q, init (tail q))
```

Imagine we could write

```
removefl :: Deque Int → (Int, Deque Int)  
removefl (f <q> l) = (f + l, q)
```

instead.



Using functions to destruct deques is tedious

Example: remove the first and the last element of a deque, return their sum and the resulting deque.

```
removefl :: Deque Int → (Int, Deque Int)  
removefl q = (head q + last q, init (tail q))
```

Imagine we could write

```
removefl :: Deque Int → (Int, Deque Int)  
removefl (f ◁ q ▷ l) = (f + l, q)
```

instead.



Datatypes and functions

```
data Front a = Nil | a ◁ Deque a  
data Back a = Lin | Deque a ▷ a
```

```
front :: Deque a → Front a  
front q = if isEmpty q then Nil  
          else head q ◁ tail q  
  
back :: Deque a → Back a  
back q = if isEmpty q then Lin  
          else init q ▷ last q
```



But then?

```
removefl :: Deque Int → (Int, Deque Int)
removefl q = case back q of
    (q' ▷ l) → case front q' of
        (f ◁ q'') → (f + l, q'')
```

This gets even worse if we want *removefl* to return $(0, \text{empty})$ on queues with less than two elements.



But then?

```
removefl :: Deque Int → (Int, Deque Int)
removefl q = case back q of
    (q' ▷ l) → case front q' of
        (f ◁ q'') → (f + l, q'')
```

This gets even worse if we want *removefl* to return $(0, \text{empty})$ on queues with less than two elements.



But then?

```
removefl :: Deque Int → (Int, Deque Int)
removefl q = case back q of
    (q' ▷ l) → case front q' of
        (f ◁ q'') → (f + 1, q'')
        -         → (0, q)
    - → (0, q)
```

This gets even worse if we want *removefl* to return $(0, \text{empty})$ on queues with less than two elements.



Pattern guards

```
removefl :: Deque Int → (Int, Deque Int)
removefl q
  | q' ▷ l ← back q, f ◁ q'' ← front q' = (f + l, q'')
```

Pattern guards

- ▶ are implemented in GHC;
- ▶ allow “list comprehension”-syntax in guards
- ▶ are a conservative extension of normal guards



Pattern guards

```
removefl :: Deque Int → (Int, Deque Int)
removefl q
  |  $q' \triangleright l \leftarrow \text{back } q, f \triangleleft q'' \leftarrow \text{front } q = (f + l, q'')$ 
  | otherwise = (0, q)
```

Pattern guards

- ▶ are implemented in GHC;
- ▶ allow “list comprehension”-syntax in guards
- ▶ are a conservative extension of normal guards



Views

view Front a of Deque $a = Nil \mid a \triangleright$ Deque a

where

$front\ q =$ **if** $isEmpty\ q$ **then** Nil
else $head\ q \triangleright tail\ q$

view Back a of Deque $a = Lin \mid$ Deque $a \triangleleft a$

where

$back\ q =$ **if** $isEmpty\ q$ **then** Lin
else $init\ q \triangleleft last\ q$

$removefl ::$ Deque Int \rightarrow (Int, Deque Int)

$removefl\ (f \triangleright q \triangleleft l) = (f + l, q)$



Uni-directional views

view Front a of Deque $a = Nil \mid a \triangleright$ Deque a

where

$front\ q =$ **if** $isEmpty\ q$ **then** Nil
else $head\ q \triangleright tail\ q$

- ▶ View constructors such as Nil and (\triangleright) must not appear on the right hand side of functions, except in the view transformation. Why?
- ▶ The view type must not be recursive. Why?
- ▶ However, the view transformation may use the view recursively ...



Recursive view definitions

view Ord $a \Rightarrow$ Minimum a of $[a]$ = *Empty* | *Min a [a]*

where

$min [] = Empty$

$min (x : Empty) = Min x []$

$min (x : Min y ys) = \mathbf{if} \ x \leq y \ \mathbf{then} \ Min \ x \ (y : ys)$
 $\mathbf{else} \ Min \ y \ (x : xs)$

$sort :: Ord \ a \Rightarrow [a] \rightarrow [a]$

$sort \ Empty = []$

$sort (Min \ x \ xs) = x : sort \ xs$



Views on classes

```
class ListLike l where
```

```
  null :: l a → Bool
```

```
  nil   :: l a
```

```
  cons :: a → l a → l a
```

```
  head :: l a → a
```

```
  tail :: l a → l a
```

```
view (ListLike l) ⇒ List l a of l a = Nil | Cons a l
```

```
where
```

```
  list xs = if null xs then Nil
```

```
           else Cons (head xs) (tail xs)
```

Can views be instances of classes?



Views in Haskell?

- ▶ Views are not implemented in any Haskell implementation.
- ▶ If they apply in both directions, isomorphism has to be checked manually.
- ▶ It is difficult to estimate the efficiency of pattern matching in the presence of views.
- ▶ It is said that pattern guards are enough.
- ▶ Nevertheless, I think that views would make a useful addition to Haskell.



Overview

Introduction (library overview)

Views

Binary search trees

Discussion



Binary search trees

| **data** BinTree $a = \text{Tip} \mid \text{Node} (\text{BinTree } a) a (\text{BinTree } a)$

Binary search tree (BST) property: Invariant for Node $l x r$:

| $\text{all } (\leq x) (\text{toList } l) \wedge \text{all } (x \leq) (\text{toList } r)$



Binary search trees

data BinTree $a = \text{Tip} \mid \text{Node} (\text{BinTree } a) a (\text{BinTree } a)$

$\text{toList} :: \text{BinTree } a \rightarrow [a]$

$\text{toList Tip} = []$

$\text{toList Node } l \ x \ r = \text{toList } l \ ++ [x] \ ++ \text{toList } r$

Binary search tree (BST) property: Invariant for $\text{Node } l \ x \ r$:

$\text{all } (\leq x) (\text{toList } l) \wedge \text{all } (x \leq) (\text{toList } r)$



Binary search trees

data BinTree $a = \text{Tip} \mid \text{Node} (\text{BinTree } a) a (\text{BinTree } a)$

$\text{toList} :: \text{BinTree } a \rightarrow [a]$

$\text{toList} = \text{toList}' []$

$\text{toList}' :: \text{BinTree } a \rightarrow [a] \rightarrow [a]$

$\text{toList}' \text{Tip} = \text{id}$

$\text{toList}' \text{Node } l \ x \ r = (\text{toList}' l++) \cdot (x:) \cdot (\text{toList}' r++)$

Binary search tree (BST) property: Invariant for $\text{Node } l \ x \ r$:

$\text{all } (\leq x) (\text{toList } l) \wedge \text{all } (x \leq) (\text{toList } r)$



Binary search trees

data BinTree $a = \text{Tip} \mid \text{Node} (\text{BinTree } a) a (\text{BinTree } a)$

$\text{toList} :: \text{BinTree } a \rightarrow [a]$

$\text{toList} = \text{toList}' []$

$\text{toList}' :: \text{BinTree } a \rightarrow [a] \rightarrow [a]$

$\text{toList}' \text{Tip} = \text{id}$

$\text{toList}' \text{Node } l \ x \ r = (\text{toList}' l++) \cdot (x:) \cdot (\text{toList}' r++)$

Binary search tree (BST) property: Invariant for $\text{Node } l \ x \ r$:

$\text{all } (\leq x) (\text{toList } l) \wedge \text{all } (x \leq) (\text{toList } r)$



Searching an element in a BST

elem :: Ord a ⇒ a → BinTree a → Bool

elem x Tip = False

elem x (Node l y r)

| x == y = True

| x < y = *elem* x l

| x > y = *elem* x r



Inserting an element in a BST

insert :: Ord a ⇒ a → BinTree a → BinTree a

insert x Tip = Node Tip x Tip

insert x (Node l y r)

| x ≤ y = Node (insert x l) y r

| x > y = Node l y (insert x r)

Observations:

- ▶ Biased insertion.
- ▶ It would be easy to disallow duplicates.
- ▶ Can lead to unbalanced trees.



Inserting an element in a BST

```
insert :: Ord a => a -> BinTree a -> BinTree a
insert x Tip = Node Tip x Tip
insert x (Node l y r)
  | x <= y = Node (insert x l) y r
  | x > y = Node l y (insert x r)
```

Observations:

- ▶ Biased insertion.
- ▶ It would be easy to disallow duplicates.
- ▶ Can lead to unbalanced trees.



Sorting using a BST

$sort :: [a] \rightarrow [a]$
 $sort = toList \cdot foldr\ insert\ Tip$

Performance?

Quadratic in the worst-case, unless BST is balanced



Sorting using a BST

$sort :: [a] \rightarrow [a]$
 $sort = toList \cdot foldr insert Tip$

Performance?

Quadratic in the worst-case, unless BST is balanced



Sorting using a BST

```
| sort :: [a] → [a]  
| sort = toList · foldr insert Tip
```

Performance?

Quadratic in the worst-case, unless BST is balanced



Balancing schemes

There are multiple balancing schemes known:

- ▶ AVL trees
- ▶ Red-black trees
- ▶ ...

It turns out to be more efficient to balance relatively rarely, because when used with random elements, sufficient balancing is often achieved on its own.



Balancing schemes

There are multiple balancing schemes known:

- ▶ AVL trees
- ▶ Red-black trees
- ▶ ...

It turns out to be more efficient to balance relatively rarely, because when used with random elements, sufficient balancing is often achieved on its own.



Data.Map and Data.Set

- ▶ From Daan Leijen's DData library.
- ▶ Since ghc-6.4 the standard finite map and set types.
- ▶ Implemented as balanced BSTs.
- ▶ Based on *Efficient sets: a balancing act* by Stephen Adams, JFP 3(4), pages 553–562, October 1993.



Rotations

Balancing is based on **rotations** (drawing).

$$\begin{aligned} & \text{singleL, singleR} :: \text{BinTree } a \rightarrow \text{BinTree } a \\ & \text{singleL } (\text{Node } l \ x \ (\text{Node } m \ y \ r)) = \text{Node } (\text{Node } l \ x \ m) \ y \ r \\ & \text{singleR } (\text{Node } (\text{Node } l \ x \ m) \ y \ r) = \text{Node } l \ x \ (\text{Node } m \ y \ r) \end{aligned}$$
$$\begin{aligned} & \text{doubleL, doubleR} :: \text{BinTree } a \rightarrow \text{BinTree } a \\ & \text{doubleL } (\text{Node } l \ x \ (\text{Node } (\text{Node } m \ y \ n) \ z \ r)) \\ & \quad = \text{Node } (\text{Node } l \ x \ m) \ y \ (\text{Node } n \ z \ r) \\ & \text{doubleR } (\text{Node } (\text{Node } l \ x \ (\text{Node } m \ y \ n)) \ z \ r) \\ & \quad = \text{Node } (\text{Node } l \ x \ m) \ y \ (\text{Node } n \ z \ r) \end{aligned}$$


Rotations

Balancing is based on **rotations** (drawing).

singleL, singleR :: BinTree a → BinTree a

singleL (Node l x (Node m y r)) = Node (Node l x m) y r

singleR (Node (Node l x m) y r) = Node l x (Node m y r)

doubleL, doubleR :: BinTree a → BinTree a

*doubleL (Node l x (Node (Node m y n) z r))
= Node (Node l x m) y (Node n z r)*

*doubleR (Node (Node l x (Node m y n)) z r)
= Node (Node l x m) y (Node n z r)*



Rotations

Balancing is based on **rotations** (drawing).

$singleL, singleR :: \text{BinTree } a \rightarrow \text{BinTree } a$

$singleL (\text{Node } l \ x \ (\text{Node } m \ y \ r)) = \text{Node } (\text{Node } l \ x \ m) \ y \ r$

$singleR (\text{Node } (\text{Node } l \ x \ m) \ y \ r) = \text{Node } l \ x \ (\text{Node } m \ y \ r)$

$doubleL, doubleR :: \text{BinTree } a \rightarrow \text{BinTree } a$

$doubleL (\text{Node } l \ x \ (\text{Node } (\text{Node } m \ y \ n) \ z \ r))$

$= \text{Node } (\text{Node } l \ x \ m) \ y \ (\text{Node } n \ z \ r)$

$doubleR (\text{Node } (\text{Node } l \ x \ (\text{Node } m \ y \ n)) \ z \ r)$

$= \text{Node } (\text{Node } l \ x \ m) \ y \ (\text{Node } n \ z \ r)$



Smart constructor

A **smart constructor** is a function that takes the role of a constructor, but performs additional operations such as to establish invariants.

Recall *makeQ*.

Now *node* :: Ord a ⇒ BinTree a → a → BinTree a → BinTree a.



Smart constructor

A **smart constructor** is a function that takes the role of a constructor, but performs additional operations such as to establish invariants.

Recall *makeQ*.

Now $node :: Ord\ a \Rightarrow BinTree\ a \rightarrow a \rightarrow BinTree\ a \rightarrow BinTree\ a.$



Smart constructor

A **smart constructor** is a function that takes the role of a constructor, but performs additional operations such as to establish invariants.

Recall *makeQ*.

Now $node :: Ord\ a \Rightarrow BinTree\ a \rightarrow a \rightarrow BinTree\ a \rightarrow BinTree\ a$.



Keeping the tree balanced

To be able to check the balance of a tree efficiently, we change the representation:

| **data** BinTree $a = \text{Tip} \mid \text{Node} ! \text{Int} (\text{BinTree } a) a (\text{BinTree})$

New invariant for $\text{Node } s \ l \ x \ r$:

| $s == \text{length } (\text{toList } l) + 1 + \text{length } (\text{toList } r)$

| $\text{size} :: \text{BinTree } a \rightarrow \text{Int}$
| $\text{size } \text{Tip} = 0$
| $\text{size } (\text{Node } s \ _ \ _ \ _) = s$



Keeping the tree balanced

To be able to check the balance of a tree efficiently, we change the representation:

| **data** BinTree $a = \text{Tip} \mid \text{Node} ! \text{Int} (\text{BinTree } a) a (\text{BinTree})$

New invariant for $\text{Node } s \ l \ x \ r$:

| $s == \text{length } (\text{toList } l) + 1 + \text{length } (\text{toList } r)$

| $\text{size} :: \text{BinTree } a \rightarrow \text{Int}$
| $\text{size } \text{Tip} = 0$
| $\text{size } (\text{Node } s \ _ \ _ \ _) = s$



Keeping the tree balanced

Smart constructor, assumes that the tree was originally balanced at that only one of the two trees has been changed by one element.

```
node :: Ord a => BinTree a -> a -> BinTree a -> BinTree a
node l x r
  | sl + sr <= 1 = Node s l x r
  | sr >= delta*sl = rotateL l x r
  | sl >= delta*sr = rotateR l x r
  | otherwise    = Node s l x r
where sl = size l
        sr = size r
        s  = size l + 1 + size r
```

The constant δ can be chosen within certain parameters and is 5 in Data.Map.



Insertion and deletion

During insertion and deletion, the smart constructor is used to maintain the balance.

```
insert :: Ord a => a -> BinTree a -> BinTree a
```

```
insert x Tip = Node Tip x Tip
```

```
insert x (Node s l y r)
```

```
  | x < y  = node (insert x l) y r
```

```
  | x > y  = node l y (insert x r)
```

```
  | x == y = Node s l x r
```

```
delete :: Ord k => a -> BinTree a -> BinTree a
```

```
delete x Tip = Tip
```

```
delete x (Node s l y r)
```

```
  | x < y  = node (delete x l) y r
```

```
  | x > y  = node l y (delete x r)
```

```
  | x == y = glue l r
```



Insertion and deletion

During insertion and deletion, the smart constructor is used to maintain the balance.

```
insert :: Ord a => a -> BinTree a -> BinTree a
```

```
insert x Tip = Node Tip x Tip
```

```
insert x (Node s l y r)
```

```
  | x < y   = node (insert x l) y r
```

```
  | x > y   = node l y (insert x r)
```

```
  | x == y  = Node s l x r
```

```
delete :: Ord k => a -> BinTree a -> BinTree a
```

```
delete x Tip = Tip
```

```
delete x (Node s l y r)
```

```
  | x < y   = node (delete x l) y r
```

```
  | x > y   = node l y (delete x r)
```

```
  | x == y  = glue l r
```



Rotating once or twice

```
rotateL l x r@(Node _ lr _ rr)
  |  $\alpha * \text{size } lr < \text{size } rr = \text{singleL } l \ x \ r$ 
  | otherwise                       = doubleL l x r
rotateR l@(Node _ ll _ rl) x r
  |  $\alpha * \text{size } rl < \text{size } ll = \text{singleR } l \ x \ r$ 
  | otherwise                       = doubleR l x r
```

Again, α is a constant that can be chosen, and is 0.5 in Data.Map.

The functions *singleL*, *doubleL*, *singleR*, *doubleR* need to be adapted to the correct type, but are otherwise the same as mentioned before.



Glueing two balanced trees together

glue :: BinTree *a* → BinTree *a* → BinTree *a*

glue Tip *r* = *r*

glue *l* Tip = *l*

glue *l* *r*

| *size l* > *size r* = **let** (*x*, *l'*) = *deleteFindMax l* **in** *node l' x r*

| *otherwise* = **let** (*x*, *r'*) = *deleteFindMin r* **in** *node l x r'*

deleteFindMax :: BinTree *a* → (*a*, BinTree *a*)

deleteFindMax (Node *l* *x* Tip) = (*x*, *l*)

deleteFindMax (Node *l* *x* *r*) =

let (*y*, *r'*) = *deleteFindMax r* **in** (*y*, *node l x r'*)



Glueing two balanced trees together

$glue :: \text{BinTree } a \rightarrow \text{BinTree } a \rightarrow \text{BinTree } a$

$glue \text{ Tip } r = r$

$glue \text{ l Tip} = l$

$glue \text{ l } r$

| $\text{size } l > \text{size } r = \mathbf{let} (x, l') = \text{deleteFindMax } l \mathbf{in} \text{ node } l' \ x \ r$

| $\text{otherwise} = \mathbf{let} (x, r') = \text{deleteFindMin } r \mathbf{in} \text{ node } l \ x \ r'$

$\text{deleteFindMax} :: \text{BinTree } a \rightarrow (a, \text{BinTree } a)$

$\text{deleteFindMax} (\text{Node } l \ x \ \text{Tip}) = (x, l)$

$\text{deleteFindMax} (\text{Node } l \ x \ r) =$

$\mathbf{let} (y, r') = \text{deleteFindMax } r \mathbf{in} (\text{node } l \ x \ r')$



Sets

- ▶ The operations discussed correspond almost directly to the implementation of sets in `Data.Set`.
- ▶ All operations (lookup, insertion, deletion) are in $O(\log n)$.
- ▶ BSTs can be used persistently. When modified, a part of the tree must be copied.
- ▶ The `Data.Set` module supports more operations: update (logarithmic), union (linear), difference (linear), intersection (linear).
- ▶ What about a map on sets?



Sets

- ▶ The operations discussed correspond almost directly to the implementation of sets in `Data.Set`.
- ▶ All operations (lookup, insertion, deletion) are in $O(\log n)$.
- ▶ BSTs can be used persistently. When modified, a part of the tree must be copied.
- ▶ The `Data.Set` module supports more operations: update (logarithmic), union (linear), difference (linear), intersection (linear).
- ▶ What about a map on sets? It's $O(n \log n)$ in general, and linear only for **monotonic** functions.



Sets

- ▶ The operations discussed correspond almost directly to the implementation of sets in Data.Set.
- ▶ All operations (lookup, insertion, deletion) are in $O(\log n)$.
- ▶ BSTs can be used persistently. When modified, a part of the tree must be copied.
- ▶ The Data.Set module supports more operations: update (logarithmic), union (linear), difference (linear), intersection (linear).
- ▶ What about a map on sets? It's $O(n \log n)$ in general, and linear only for **monotonic** functions.



Finite maps

The step from sets to finite maps is very small:
We use set elements of the form $(key, value)$, where the order is determined only by the keys.

In practice, we use a specialized datatype

```
| data Map k a = Tip | Node ! Int (Map k a) k a (Map k a)
```

and adapt all the operations.



Finite maps

The step from sets to finite maps is very small:
We use set elements of the form $(key, value)$, where the order is determined only by the keys.

In practice, we use a specialized datatype

```
| data Map k a = Tip | Node ! Int (Map k a) k a (Map k a)
```

and adapt all the operations.



What if operations on the key types are expensive?

...for example, if we use strings as keys.

Use a **trie** (aka **digital search tree**).



What if operations on the key types are expensive?

...for example, if we use strings as keys.

Use a **trie** (aka **digital search tree**).



What if operations on the key types are expensive?

...for example, if we use strings as keys.

Use a **trie** (aka **digital search tree**).



Tries

Trie representation for keys of type $[k]$:

```
data Trie k a = Node (Maybe a) (Map k (Trie k a))
```

- ▶ Can be generalized to other structures of keys than lists.
- ▶ Can be implemented as a type-indexed type.
- ▶ Are currently not available as a standard GHC library.



Overview

Introduction (library overview)

Views

Binary search trees

Discussion



Other data structures

- ▶ Heaps/Priority queues.
- ▶ Hybrid structures: priority search queues, finger trees.
- ▶ Functional graphs.



Papers to read

- ▶ Okasaki
- ▶ Hinze



Practical advice

Be bold enough to use a non-list data structure once in a while.

At the very least, use finite maps when random-access is desired, and avoid arrays when multiple updates occur.



Observations

- ▶ Functional languages are suitable to express complex data structures clearly.
- ▶ Persistence is not always expensive.
- ▶ Laziness can sometimes be helpful in the context of persistence.
- ▶ There are still few, but nevertheless usable libraries for data structures available in Haskell.
- ▶ Views are useful.



Haskell as a language for data structures

Clear advantages, but also problems:

- ▶ lazy evaluation
- ▶ not a real module system
- ▶ no views
- ▶ at least multi-parameter type classes with fundeps needed

