

Data Structures I

Advanced Functional Programming

Andres Löh (andres@cs.uu.nl)

Universiteit Utrecht

19 May 2005

Universiteit Utrecht



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



Question

What is the most frequently used data structure in Haskell?

Clearly, lists ...



Question

What is the most frequently used data structure in Haskell?

Clearly, lists ...



What are lists good for?

$head :: [a] \rightarrow a$

$tail :: [a] \rightarrow [a]$

$(:) :: a \rightarrow [a] \rightarrow [a]$

- ▶ These are efficient operations on lists.
- ▶ These are the **stack** operations.



What are lists good for?

$head :: [a] \rightarrow a$

$tail :: [a] \rightarrow [a]$

$(:) :: a \rightarrow [a] \rightarrow [a]$

- ▶ These are efficient operations on lists.
- ▶ These are the **stack** operations.



What are lists good for?

$head :: [a] \rightarrow a \quad -- O(1)$

$tail :: [a] \rightarrow [a] \quad -- O(1)$

$(:) :: a \rightarrow [a] \rightarrow [a] \quad -- O(1)$

- ▶ These are efficient operations on lists.
- ▶ These are the **stack** operations.



What are lists good for?

$top :: [a] \rightarrow a \quad -- O(1)$

$pop :: [a] \rightarrow [a] \quad -- O(1)$

$push :: a \rightarrow [a] \rightarrow [a] \quad -- O(1)$

- ▶ These are efficient operations on lists.
- ▶ These are the **stack** operations.



Haskell stacks are persistent

A data structure is called **persistent** if after an operation both the original and the resulting version of the data structure are available.

If not persistent, a data structure is called **ephemeral**.

- ▶ Functional data structures are naturally persistent.
- ▶ Imperative data structures are usually ephemeral.
- ▶ Persistent data structures are often, but not always, less efficient than ephemeral data structures.



Haskell stacks are persistent

A data structure is called **persistent** if after an operation both the original and the resulting version of the data structure are available.

If not persistent, a data structure is called **ephemeral**.

- ▶ Functional data structures are naturally persistent.
- ▶ Imperative data structures are usually ephemeral.
- ▶ Persistent data structures are often, but not always, less efficient than ephemeral data structures.



Haskell stacks are persistent

A data structure is called **persistent** if after an operation both the original and the resulting version of the data structure are available.

If not persistent, a data structure is called **ephemeral**.

- ▶ Functional data structures are naturally persistent.
- ▶ Imperative data structures are usually ephemeral.
- ▶ Persistent data structures are often, but not always, less efficient than ephemeral data structures.



Other operations on lists

<i>snoc</i>	$:: [a] \rightarrow a \rightarrow [a]$	$-- O(n)$
<i>snoc</i>	$= \lambda xs x \rightarrow xs \# [x]$	
$(\#)$	$:: [a] \rightarrow [a] \rightarrow [a]$	$-- O(n)$
<i>reverse</i>	$:: [a] \rightarrow [a]$	$-- O(n), \text{ naively: } O(n^2)$
<i>union</i>	$:: \text{Eq } a \Rightarrow [a] \rightarrow [a] \rightarrow [a]$	$-- O(mn)$
<i>elem</i>	$:: \text{Eq } a \Rightarrow a \rightarrow [a] \rightarrow \text{Bool}$	$-- O(n)$

Although not efficient for these purposes, Haskell lists are frequently used as

- ▶ arrays
- ▶ queues, dequeues, catenable queues
- ▶ sets
- ▶ lookup tables, association lists, finite maps
- ▶ ...

Why?



Other operations on lists

<i>snoc</i>	$:: [a] \rightarrow a \rightarrow [a]$	$-- O(n)$
<i>snoc</i>	$= \lambda xs x \rightarrow xs \# [x]$	
$(\#)$	$:: [a] \rightarrow [a] \rightarrow [a]$	$-- O(n)$
<i>reverse</i>	$:: [a] \rightarrow [a]$	$-- O(n)$, naively: $O(n^2)$
<i>union</i>	$:: \text{Eq } a \Rightarrow [a] \rightarrow [a] \rightarrow [a]$	$-- O(mn)$
<i>elem</i>	$:: \text{Eq } a \Rightarrow a \rightarrow [a] \rightarrow \text{Bool}$	$-- O(n)$

Although not efficient for these purposes, Haskell lists are frequently used as

- ▶ arrays
- ▶ queues, dequeues, catenable queues
- ▶ sets
- ▶ lookup tables, association lists, finite maps
- ▶ ...

Why?



Other operations on lists

$snoc$	$:: [a] \rightarrow a \rightarrow [a]$	$-- O(n)$
$snoc$	$= \lambda xs x \rightarrow xs \# [x]$	
$(\#)$	$:: [a] \rightarrow [a] \rightarrow [a]$	$-- O(n)$
$reverse$	$:: [a] \rightarrow [a]$	$-- O(n), \text{ naively: } O(n^2)$
$union$	$:: Eq\ a \Rightarrow [a] \rightarrow [a] \rightarrow [a]$	$-- O(mn)$
$elem$	$:: Eq\ a \Rightarrow a \rightarrow [a] \rightarrow Bool$	$-- O(n)$

Although not efficient for these purposes, Haskell lists are frequently used as

- ▶ arrays
- ▶ queues, dequeues, catenable queues
- ▶ sets
- ▶ lookup tables, association lists, finite maps
- ▶ ...

Why?



Lists are everywhere, because ...

- ▶ There is a convenient built-in notation for lists.
- ▶ There are even list comprehensions in Haskell.
- ▶ Lots of library functions on lists.
- ▶ Pattern matching!
- ▶ Haskell strings are lists.
- ▶ Other data structures not widely known.
- ▶ Arrays are often worse.
- ▶ Not enough standard libraries for data structures.

We are going to change **this** ...

Unfortunately, the remaining reasons are valid.



Lists are everywhere, because ...

- ▶ There is a convenient built-in notation for lists.
- ▶ There are even list comprehensions in Haskell.
- ▶ Lots of library functions on lists.
- ▶ Pattern matching!
- ▶ Haskell strings are lists.
- ▶ Other data structures not widely known.
- ▶ Arrays are often worse.
- ▶ Not enough standard libraries for data structures.

We are going to change **this** ...

Unfortunately, the remaining reasons are valid.



Lists are everywhere, because ...

- ▶ There is a convenient built-in notation for lists.
- ▶ There are even list comprehensions in Haskell.
- ▶ Lots of library functions on lists.
- ▶ Pattern matching!
- ▶ Haskell strings are lists.
- ▶ **Other data structures not widely known.**
- ▶ Arrays are often worse.
- ▶ Not enough standard libraries for data structures.

We are going to change **this** ...

Unfortunately, the remaining reasons are valid.



Lists are everywhere, because ...

- ▶ There is a convenient built-in notation for lists.
- ▶ There are even list comprehensions in Haskell.
- ▶ Lots of library functions on lists.
- ▶ Pattern matching!
- ▶ Haskell strings are lists.
- ▶ Other data structures not widely known.
- ▶ Arrays are often worse.
- ▶ Not enough standard libraries for data structures.

We are going to change **this** ...

Unfortunately, the remaining reasons are valid.



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



About arrays

Imperative arrays feature

- ▶ constant-time lookup
- ▶ constant-time update

Update is usually at least as important as lookup.

Functional arrays do

- ▶ lookup in $O(1)$; yay!
- ▶ update in $O(n)$! Why? Persistence!

Array update is even worse than list update.

- ▶ To update the n th element of a list, $n - 1$ elements are copied.
- ▶ To update any element of an array, the **whole** array is copied.



About arrays

Imperative arrays feature

- ▶ constant-time lookup
- ▶ constant-time update

Update is usually at least as important as lookup.

Functional arrays do

- ▶ lookup in $O(1)$; yay!
- ▶ update in $O(n)$! Why? Persistence!

Array update is even worse than list update.

- ▶ To update the n th element of a list, $n - 1$ elements are copied.
- ▶ To update any element of an array, the **whole** array is copied.



About arrays

Imperative arrays feature

- ▶ constant-time lookup
- ▶ constant-time update

Update is usually at least as important as lookup.

Functional arrays do

- ▶ lookup in $O(1)$; yay!
- ▶ update in $O(n)$! Why? Persistence!

Array update is even worse than list update.

- ▶ To update the n th element of a list, $n - 1$ elements are copied.
- ▶ **To update any element of an array, the whole array is copied.**



Space efficiency vs. space leaks

Arrays can be stored in a compact way.
Lists require lots of pointers.

If arrays are updated frequently and used persistently,
space leaks will occur!



Space efficiency vs. space leaks

Arrays can be stored in a compact way.
Lists require lots of pointers.

If arrays are updated frequently and used persistently,
space leaks will occur!



Mutable arrays

- ▶ Are like imperative arrays.
- ▶ Defined in `Data.Array.MArray` and `Data.Array.IO`.
- ▶ All operations in a state monad (possibly IO monad).
- ▶ Often awkward to use in a functional setting.



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



Boxed vs. unboxed types

Haskell data structures are **boxed**.

- ▶ Each value is behind an additional indirection.
- ▶ This allows polymorphic datastructures (because the size of a pointer is always the same).
- ▶ This allows laziness, because the pointer can be to a computation as well as to evaluated data.

GHC offers **unboxed** datatypes, too. Naturally, they

- ▶ are slightly more efficient (in both space and time),
- ▶ are strict,
- ▶ cannot be used in polymorphic data structures.



Boxed vs. unboxed types

Haskell data structures are **boxed**.

- ▶ Each value is behind an additional indirection.
- ▶ This allows polymorphic datastructures (because the size of a pointer is always the same).
- ▶ This allows laziness, because the pointer can be to a computation as well as to evaluated data.

GHC offers **unboxed** datatypes, too. Naturally, they

- ▶ are slightly more efficient (in both space and time),
- ▶ are strict,
- ▶ cannot be used in polymorphic data structures.



Unboxed types

- ▶ Defined in `GHC.Base`.
- ▶ For example, `Int#`, `Char#`, `Double#`.
- ▶ Have kind `#`, not `*`.
- ▶ Use specialized operations such as

| `(+#) :: Int# → Int# → Int#`

- ▶ Cannot be used in polymorphic functions or datatypes.
- ▶ Are used by GHC internally to define the usual datatypes:

| `data Int = I# Int#`



Packed strings

- ▶ Defined in `Data.PackedString`.
- ▶ Implemented as immutable, unboxed arrays.
- ▶ Can be more space-efficient than standard strings.
- ▶ Manipulating packed strings can be expensive.



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



Queues

- ▶ Stacks are LIFO (last-in-first-out).
- ▶ Queues are FIFO (first-in-first-out).
- ▶ A list is not very suitable to represent a queue, because efficient access to both ends is desired.

The standard trick is:

```
| data Queue a = Q [a] [a]
```

The first list is the **front**, the second the **back** of the queue, in reversed order.



Queues

- ▶ Stacks are LIFO (last-in-first-out).
- ▶ Queues are FIFO (first-in-first-out).
- ▶ A list is not very suitable to represent a queue, because efficient access to both ends is desired.

The standard trick is:

data Queue $a = Q [a] [a]$

The first list is the **front**, the second the **back** of the queue, in reversed order.



Queues

- ▶ Stacks are LIFO (last-in-first-out).
- ▶ Queues are FIFO (first-in-first-out).
- ▶ A list is not very suitable to represent a queue, because efficient access to both ends is desired.

The standard trick is:

| **data** Queue $a = Q [a] [a]$

The first list is the **front**, the second the **back** of the queue, in reversed order.



Queue operations

This is what we want for a queue:

```
empty    :: Queue a                -- produce an empty queue
snoc     :: a → Queue a → Queue a -- insert at the back
head     :: Queue a → a           -- get first element
tail     :: Queue a → Queue a     -- remove first element
```

```
toList   :: Queue a → [a]         -- queue to list
fromList :: [a] → Queue a         -- list to queue
```



Queue operations

This is what we want for a queue:

<i>empty</i>	:: Queue <i>a</i>	-- produce an empty queue
<i>snoc</i>	:: <i>a</i> → Queue <i>a</i> → Queue <i>a</i>	-- insert at the back
<i>head</i>	:: Queue <i>a</i> → <i>a</i>	-- get first element
<i>tail</i>	:: Queue <i>a</i> → Queue <i>a</i>	-- remove first element
<i>toList</i>	:: Queue <i>a</i> → [<i>a</i>]	-- queue to list
<i>fromList</i>	:: [<i>a</i>] → Queue <i>a</i>	-- list to queue



Implementing queue operations

$empty :: Queue\ a$
 $empty = Q\ []\ []$

$snoc :: a \rightarrow Queue\ a \rightarrow Queue\ a$
 $snoc\ x\ (Q\ fs\ bs) = Q\ fs\ (x : bs)$



Implementing queue operations

$empty :: Queue\ a$
 $empty = Q\ []\ []$

$snoc :: a \rightarrow Queue\ a \rightarrow Queue\ a$
 $snoc\ x\ (Q\ fs\ bs) = Q\ fs\ (x : bs)$



Invocations of *reverse*

$head :: Queue\ a \rightarrow a$

$head\ (Q\ (f : fs)\ bs) = f$

$head\ (Q\ []\ bs) = head\ (reverse\ bs)$

$tail :: Queue\ a \rightarrow Queue\ a$

$tail\ (Q\ (f : fs)\ bs) = Q\ fs\ bs$

$tail\ (Q\ []\ bs) = Q\ (tail\ (reverse\ bs))\ []$

Persistence spoils the fun here; without persistence, all operations would be in $O(1)$ amortized time.



Invocations of *reverse*

$head :: Queue\ a \rightarrow a$

$head\ (Q\ (f : fs)\ bs) = f$

$head\ (Q\ []\ bs) = head\ (reverse\ bs)$

$tail :: Queue\ a \rightarrow Queue\ a$

$tail\ (Q\ (f : fs)\ bs) = Q\ fs\ bs$

$tail\ (Q\ []\ bs) = Q\ (tail\ (reverse\ bs))\ []$

Persistence spoils the fun here; without persistence, all operations would be in $O(1)$ amortized time.



Invocations of *reverse*

$head :: Queue\ a \rightarrow a$

$head\ (Q\ (f : fs)\ bs) = f$

$head\ (Q\ []\ bs) = head\ (reverse\ bs)$

$tail :: Queue\ a \rightarrow Queue\ a$

$tail\ (Q\ (f : fs)\ bs) = Q\ fs\ bs$

$tail\ (Q\ []\ bs) = Q\ (tail\ (reverse\ bs))\ []$

Persistence spoils the fun here; without persistence, all operations would be in $O(1)$ amortized time.



Amortized analysis

Amortized complexity can be better than **worst-case complexity** if the worst-case cannot occur that often in practice.

In an amortized analysis, we look at the cost of multiple operations rather than single operations.



Idea

- ▶ Distribute the work that *reverse* causes over multiple operations in such a way that the amortized cost of each operation is constant.
- ▶ Use laziness (and memoization) to ensure that expensive operations are not performed too early or too often.



Memoization

A suspended expression in a lazy language is evaluated only once. The suspension is then updated with the result. Whenever the same expression is needed again, the result can be used immediately. This is called **memoization**.



Memoization

A suspended expression in a lazy language is evaluated only once. The suspension is then updated with the result. Whenever the same expression is needed again, the result can be used immediately. This is called **memoization**.



Efficient queues

Recall the queue representation:

| **data** Queue $a = Q \ [a] \ [a]$

As we will see, the work of reversing the list can be distributed well by choosing the following invariant for $Q \ fs \ bs$:

| $length \ fs \geq length \ bs$

In particular, $length \ fs == 0$ if and only if the queue is empty.

We need the lengths of both lists available in constant time.



Efficient queues

Recall the queue representation:

| **data** Queue $a = Q [a] [a]$

As we will see, the work of reversing the list can be distributed well by choosing the following invariant for $Q fs bs$:

| $length fs \geq length bs$

In particular, $length fs == 0$ if and only if the queue is empty.

We need the lengths of both lists available in constant time.



Efficient queues

Recall the queue representation:

| **data** Queue $a = Q \text{ !Int } [a] \text{ !Int } [a]$

As we will see, the work of reversing the list can be distributed well by choosing the following invariant for $Q \text{ !Int } [a] \text{ !Int } [a]$ *lf fs lb bs*:

| *lf* \geq *lb*

In particular, *length fs* == 0 if and only if the queue is empty.

We need the lengths of both lists available in constant time.



empty and *head* are simple due to the invariant

```
empty :: Queue a  
empty = Q 0 [] 0 []
```

```
head :: Queue a      → a  
head (Q _ (f:fs) bs) = f  
head (Q _ []      _) = error "empty queue"
```



empty and *head* are simple due to the invariant

empty :: Queue *a*
empty = Q 0 [] 0 []

head :: Queue *a* → *a*
head (Q _ (f:fs) bs) = f
head (Q _ [] _) = error "empty queue"



What about *tail* and *snoc*?

```
tail :: Queue a          → a
tail (Q lf (f : fs) lb b) = makeQ (lf - 1) fs lb b
tail (Q - [] - -) = error "empty queue"
```

```
snoc :: a → Queue a     → Queue a
snoc x (Q lf fs lb bs) = makeQ lf fs (lb + 1) (x : bs)
```

In both cases, we have to make a new queue using a call `makeQ lf f lb f'`, where we may need to re-establish the invariant.



What about *tail* and *snoc*?

$$\begin{aligned} \text{tail} &:: \text{Queue } a && \rightarrow a \\ \text{tail} \quad (\text{Q } lf \ (f : fs) \ lb \ b) &= \text{makeQ } (lf - 1) \ fs \ lb \ b \\ \text{tail} \quad (\text{Q } - \ [] \quad - \ -) &= \text{error "empty queue"} \end{aligned}$$
$$\begin{aligned} \text{snoc} &:: a \rightarrow \text{Queue } a && \rightarrow \text{Queue } a \\ \text{snoc} \quad x \quad (\text{Q } lf \ fs \ lb \ bs) &= \text{makeQ } lf \ fs \ (lb + 1) \ (x : bs) \end{aligned}$$

In both cases, we have to make a new queue using a call $\text{makeQ } lf \ f \ lb \ f'$, where we may need to re-establish the invariant.



What about *tail* and *snoc*?

```
tail :: Queue a          → a
tail (Q lf (f : fs) lb b) = makeQ (lf - 1) fs lb b
tail (Q - [] - -) = error "empty queue"
```

```
snoc :: a → Queue a     → Queue a
snoc x (Q lf fs lb bs) = makeQ lf fs (lb + 1) (x : bs)
```

In both cases, we have to make a new queue using a call *makeQ lf f lb f'*, where we may need to re-establish the invariant.



How to make a queue

```
makeQ :: Int → [a] → Int → [a]
        → Queue a
```

```
makeQ lf fs lb bs
```

```
  | lf ≥ lb = Q lf fs lb bs
```

```
  | otherwise = Q (lf + lb) (fs ++ reverse bs) 0 []
```



Why is this implementation “better”?

(drawing and lots of handwaving)

Read Okasaki's book for a proof.



Why is this implementation “better”?

(drawing and lots of handwaving)

Read Okasaki's book for a proof.



Queues in GHC

- ▶ Available in `Data.Queue` (ghc-6.4).
- ▶ Based on a slight variation of the implementation described here, allowing operations in constant worst-case time (“real-time queues”).
- ▶ Representation of queues is then

```
| data Queue a = Q [a] [a] [a]
```

where the third list is used to maintain the unevaluated part of the front queue.

- ▶ Described in the paper *Simple and efficient functional queues and dequeues*, JFP 5(4), pages 583–592, October 1995, by Chris Okasaki.
- ▶ Also described in Okasaki’s book.



Queues in GHC

- ▶ Available in `Data.Queue` (ghc-6.4).
- ▶ Based on a slight variation of the implementation described here, allowing operations in constant worst-case time (“real-time queues”).
- ▶ Representation of queues is then

| data `Queue a = Q [a] [a] [a]`

where the third list is used to maintain the unevaluated part of the front queue.

- ▶ Described in the paper *Simple and efficient functional queues and dequeues*, JFP 5(4), pages 583–592, October 1995, by Chris Okasaki.
- ▶ Also described in Okasaki’s book.



Queues in GHC

- ▶ Available in `Data.Queue` (ghc-6.4).
- ▶ Based on a slight variation of the implementation described here, allowing operations in constant worst-case time (“real-time queues”).
- ▶ Representation of queues is then

| **data** Queue $a = Q [a] [a] [a]$

where the third list is used to maintain the unevaluated part of the front queue.

- ▶ Described in the paper *Simple and efficient functional queues and dequeues*, JFP 5(4), pages 583–592, October 1995, by Chris Okasaki.
- ▶ Also described in Okasaki’s book.



Dequeues

A **deque** is a double-ended queue.

Operations like queue operations:

```
empty    :: Deque a           -- produce an empty queue
snoc     :: a → Deque a → Deque a -- insert at the back
head     :: Deque a → a       -- get first element
tail     :: Deque a → Deque a  -- remove first element
toList   :: Deque a → [a]     -- queue to list
fromList :: [a] → Deque a     -- list to queue
```

Additionally (also in constant time):

```
cons     :: a → Deque a → Deque a -- insert at the front
init     :: Deque a → Deque a     -- remove last element
last     :: Deque a → a           -- get last element
```



Dequeues

A **deque** is a double-ended queue.

Operations like queue operations:

<i>empty</i>	:: Deque <i>a</i>	-- produce an empty queue
<i>snoc</i>	:: <i>a</i> → Deque <i>a</i> → Deque <i>a</i>	-- insert at the back
<i>head</i>	:: Deque <i>a</i> → <i>a</i>	-- get first element
<i>tail</i>	:: Deque <i>a</i> → Deque <i>a</i>	-- remove first element
<i>toList</i>	:: Deque <i>a</i> → [<i>a</i>]	-- queue to list
<i>fromList</i>	:: [<i>a</i>] → Deque <i>a</i>	-- list to queue

Additionally (also in constant time):

<i>cons</i>	:: <i>a</i> → Deque <i>a</i> → Deque <i>a</i>	-- insert at the front
<i>init</i>	:: Deque <i>a</i> → Deque <i>a</i>	-- remove last element
<i>last</i>	:: Deque <i>a</i> → <i>a</i>	-- get last element



Efficient deques

The previous implementation can easily be extended to work for deques:

| **data** Deque $a = D ! \text{Int } [a] ! \text{Int } [a]$

Of course, we have to make the representation more symmetric. The invariant for $D \text{ } lf \text{ } fs \text{ } lb \text{ } bs$ becomes:

| $lf \leq c * lb + 1 \wedge lr \leq c * lf + 1$

(for some constant $c > 1$).



Implementation of deques

- ▶ The implementation of *makeQ* must be adapted to maintain this invariant.
- ▶ The other operations are straight-forward, we only have to pay attention to the one-element queue.
- ▶ How much time does it cost to reverse a deque?
- ▶ Unfortunately, there currently is no standard Haskell library for deques.



Catenable queues or dequeues

- ▶ Queues or dequeues that support efficient concatenation are called **catenable**.
- ▶ It is possible to support concatenation in $O(\log n)$ and even in $O(1)$ amortized time, but this requires a completely different implementation of queues/deques.
- ▶ Again, there currently are no standard Haskell libraries for catenable queues and dequeues.



Overview

Introduction (Lists)

Arrays

Unboxed types

Queues and dequeues

Summary and next lecture



Summary

- ▶ Lists are everywhere in Haskell, for a lot of good reasons.
- ▶ Functional data structures are persistent.
- ▶ Persistence and efficiency and evaluation order all interact.
- ▶ Array updates are inherently inefficient in a functional language.
- ▶ Queues and deques support many operations efficiently that normal lists do not.
- ▶ In a persistent setting, queue and deque operations can be implemented with the same complexity bounds as in an ephemeral setting.
- ▶ GHC has a standard library that supports many (but not all desirable) datastructures, for instance lists, queues, arrays in all flavors, but also unboxed types and packed strings.



Next lecture

- ▶ Pattern matching, abstract datatypes, views.
- ▶ Trees, finite maps and sets.
- ▶ ...

