

Bootstrapping One-sided Flexible Arrays

RALF HINZE

Institut für Informatik III, Universität Bonn, Römerstraße 164
53117 Bonn, Germany

Email: `ralf@informatik.uni-bonn.de`

Homepage: `http://www.informatik.uni-bonn.de/~ralf`

February, 2002

(Pick the slides at `.../~ralf/talks.html#T31.`)

Motivation

One-sided flexible arrays support look-up and update of elements and can grow and shrink at one end.

A variety of tree-based implementations is available.

- ✘ Braun trees:
all operations in $\Theta(\log n)$ time.
- ✘ Binary random-access lists:
 $\Theta(\log n)$ access and $\Theta(1)$ list operations (amortized).
- ✘ Skew binary random-access list:
 $\Theta(\log n)$ access and $\Theta(1)$ list operations (worst-case).

A common characteristic of the tree-based implementations is the *logarithmic* time bound for the look-up operation.

Motivation—continued

Assume that you have an application that uses indexing a lot but also updates or extends occasionally (ruling out 'real' arrays).

There is no data structure available that fits these needs.

Idea: Improve look-up by using *fat* multiway trees trading the running time of look-up operations for the running time of update operations (this idea is due to Chris Okasaki).

Signature

infixl 9 !

class *Array* *a* **where**

-- array-like operations

(!) :: $a\ x \rightarrow Int \rightarrow x$

update :: $(x \rightarrow x) \rightarrow Int \rightarrow a\ x \rightarrow a\ x$

-- list-like operations

empty :: $a\ x \rightarrow Bool$

size :: $a\ x \rightarrow Int$

nil :: $a\ x$

copy :: $Int \rightarrow x \rightarrow a\ x$

cons :: $x \rightarrow a\ x \rightarrow a\ x$

Signature—continued

head :: $a\ x \rightarrow x$

tail :: $a\ x \rightarrow a\ x$

-- mapping functions

map :: $(x \rightarrow y) \rightarrow (a\ x \rightarrow a\ y)$

zip :: $(x \rightarrow y \rightarrow z) \rightarrow (a\ x \rightarrow a\ y \rightarrow a\ z)$

-- conversion functions

list :: $a\ x \rightarrow [x]$

array :: $[x] \rightarrow a\ x$

Notational convenience: we write both *map f* and *zip f* simply as f^* .

Multiway trees

Idea: bootstrap an implementation based on multiway trees from a standard implementation of flexible arrays.

data $Tree\ a\ x = \langle a\ x, a\ (Tree\ a\ x) \rangle$

A node $\langle xs, ts \rangle$ is a pair consisting of an array xs of elements, called the *prefix*, and an array ts of subtrees.

We will show how to turn $Tree\ a$ into an instance of $Array$ given that a is already an instance.

instance $(Array\ a) \Rightarrow Array\ (Tree\ a)$ **where**

List-like operations

Let us start with the *cons* operation since this operation will determine the way indexing is done.

Idea: fill up the root node; if it is full up, distribute the elements evenly among the subtrees and start afresh.

$$\begin{array}{l} \text{cons } x \langle xs, ts \rangle \\ | \text{ size } xs < \text{size } ts \quad = \quad \langle \text{cons } x \ xs, \ ts \rangle \\ | \text{ otherwise} \quad \quad \quad = \quad \langle \text{cons } x \ \text{nil}, \ \text{cons}^* \ xs \ ts \rangle \end{array}$$

To make this algorithm work, we have to maintain some invariants.

Invariants

For all nodes $t = \langle \text{array } [x_1, \dots, x_m], \text{array } [t_1, \dots, t_n] \rangle$:

- ❶ $m \leq n$ and $1 \leq n$,
- ❷ $|t_1| = \dots = |t_n|$,
- ❸ $|t| = 0 \iff m = 0$.

where $|t|$ denotes the size of a tree (the total number of elements).

NB. The third invariant is necessary so that we can effectively check whether a given tree is empty.

List-like operations—continued

$empty \langle xs, ts \rangle = empty \ xs$

$head \langle xs, ts \rangle$

| $empty \ xs = error \ "head: \ empty \ array"$

| $otherwise = head \ xs$

$tail \langle xs, ts \rangle$

| $empty \ xs = error \ "tail: \ empty \ array"$

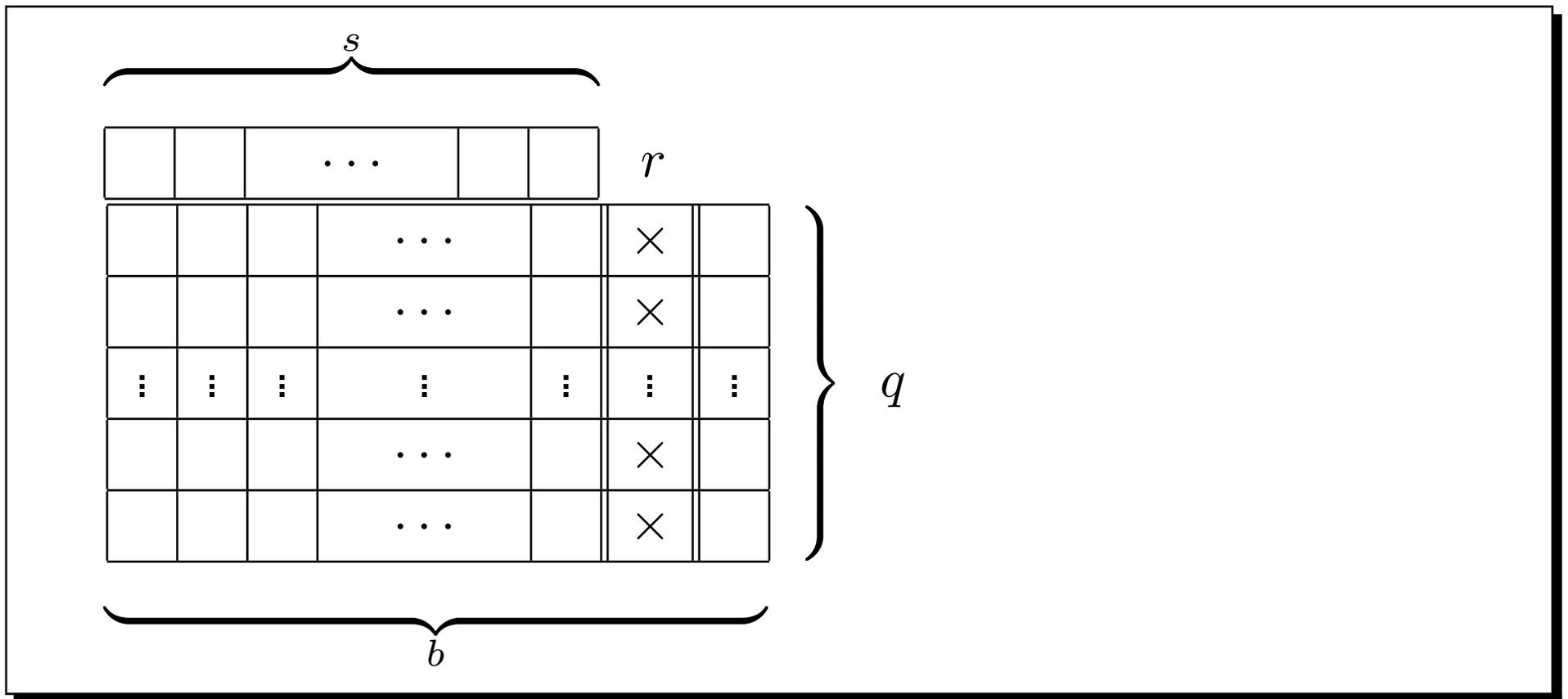
| $size \ xs > 1 = \langle tail \ xs, ts \rangle$

| $all \ empty \ ts = \langle nil, ts \rangle$

| $otherwise = \langle head^* \ ts, tail^* \ ts \rangle$

Array-like operations

After q overflows the r -th subtree comprises the q elements at positions $s + 0 * b + r$, $s + 1 * b + r$, \dots , $s + (q - 1) * b + r$, where s is the size of the prefix and b is the total number of subtrees.



Array creation

The shape of the multiway trees is solely determined by the array creation functions *nil*, *copy*, and *array*.

We may think of an initial array as an infinite tree, whose branching structure is fixed and which will be populated through repeated applications of *cons*.

☞ To be able to analyze the running times reasonably well, we make one further assumption: *we require that nodes of the same level have the same size*.

Given this assumption the structure of trees can be described using a special number system, the so-called *mixed-radix number system*

Mixed-radix number systems

A *mixed-radix numeral* is given by a sequence of digits d_0, d_1, d_2, \dots (determining the size of the element arrays) and a sequence of bases b_0, b_1, b_2, \dots (determining the size of the subtree arrays).

$$\begin{bmatrix} d_0, d_1, d_2, \dots \\ b_0, b_1, b_2, \dots \end{bmatrix} = \sum_{i=0} d_i \cdot w_i \text{ where } w_i = b_{i-1} \cdot \dots \cdot b_1 \cdot b_0$$

The bases are positive numbers $1 \leq b_i$ and we require the digits to lie in the range $0 \leq d_i \leq b_i$ (cf Invariant 1). Furthermore, we require $d_i = 0 \implies d_{i+1} = 0$ (cf Invariant 3).

 Each natural number has a unique representation in this system.

Converting from the mixed-radix number system

```
type Bases = [Int]
type Mix    = [(Int, Int)]
```

Converting a mixed-radix number to a natural number is straightforward (using the Horner's rule).

```
decode      :: Mix → Int
decode ((d, b) : σ)
  | d == 0   = 0
  | otherwise = d + b * decode σ
```

Converting to the mixed-radix number system

Given a list of bases we can easily convert a natural number into a mixed-radix number.

```
encode           :: Bases → (Int → Mix)  
encode (b : bs) n  
  | n == 0       = (0, b) : zip (repeat 0) bs  
  | otherwise    = (r + 1, b) : encode bs q  
where (q, r) = divMod (n - 1) b
```

Generic array creation functions

$$\begin{aligned} gnil &:: (\text{Array } a) \Rightarrow \text{Bases} \rightarrow \text{Tree } a \ x \\ gnil \ (b : bs) &= \langle nil, \text{copy } b \ (gnil \ bs) \rangle \\ gcopy &:: (\text{Array } a) \Rightarrow \text{Mix} \rightarrow x \rightarrow \text{Tree } a \ x \\ gcopy \ ((d, b) : \sigma) \ x &= \langle \text{copy } d \ x, \text{copy } b \ (gcopy \ \sigma \ x) \rangle \end{aligned}$$

 Now, all we have to do is to come up with interesting bases.

Analysis of running times

Let $H(n)$ be the height of the *tallest* tree with size n .

The running time of ' $!$ ' and *update* is

$$\begin{aligned}\mathcal{T}_!(n) &= \sum_{i=0}^{H(n)-1} \bar{\mathcal{T}}_!(b_i) \\ \mathcal{T}_{update}(n) &= \sum_{i=0}^{H(n)-1} \bar{\mathcal{T}}_{update}(b_i),\end{aligned}$$

where $\bar{\mathcal{T}}_{op}$ is the running time of *op* on base arrays.

Analysis of running times—continued

The *amortized* running-time of *cons* is given by

$$\mathcal{T}_{cons}(n) = \frac{1}{n} \sum_{i=0}^{H(n)-1} \frac{n}{w_i} w_i \bar{\mathcal{T}}_{cons}(b_i) = \sum_{i=0}^{H(n)-1} \bar{\mathcal{T}}_{cons}(b_i).$$

The sum calculates the costs of n successive *cons* operations. If we divide the result by n , we obtain the amortized running-time. Each summand describes the total costs at level i : we have a carry every n/w_i steps; if a carry occurs w_i nodes must be rearranged; and the rearrangement of one node takes $\bar{\mathcal{T}}_{cons}(b_i)$ time.

Variant 1: b -ary trees

Mixed-radix numeral:

$$\left[\begin{array}{c} d_0, d_1, d_2, \dots, d_n, \dots \\ b, b, b, \dots, b, \dots \end{array} \right]$$

The radices are constant.

$bary \quad :: \quad Int \rightarrow Bases$

$bary \ b \ = \ repeat \ b$

b -ary trees—running times

Performance:

base array	bootstrapped array
$\Theta(1)$	$\Theta(\log n)$
$\Theta(\log n)$	$\Theta(\log n)$
$\Theta(n)$	$\Theta(\log n)$

Of course, the constants hidden in the Θ notation differ widely.

$$\mathcal{T}_{op}(n) \approx \lg n \cdot \bar{\mathcal{T}}_{op}(b) / \lg b.$$

Arithmetic progression trees—running times

If we fix $\alpha = \beta = 1$, we obtain the so-called *factorial number system*.

$$\begin{bmatrix} d_0, d_1, d_2, \dots, d_n, & \dots \\ 1, 2, 3, \dots, n+1, \dots \end{bmatrix}$$

Performance ($\alpha = \beta = 1$):

base array	bootstrapped array
$\Theta(1)$	$\Theta(\log n / \log \log n)$
$\Theta(\log n)$	$\Theta(\log n)$
$\Theta(n)$	$\Theta((\log n)^2 / (\log \log n)^2)$

An excerpt of the asymptotic hierarchy

$$\log \log n$$
$$\prec \sqrt{\log n}$$
$$\prec \log n / \log \log n$$
$$\prec \log n$$
$$\prec (\log n)^2 / (\log \log n)^2$$
$$\prec 2^{\sqrt{\log n}}$$
$$\prec n$$

Arithmetic progression trees—running times

Performance ($\alpha = 1$ and $\beta = 2$):

base array	bootstrapped array
$\Theta(1)$	$\Theta(\sqrt{\log n})$
$\Theta(\log n)$	$\Theta(\log n)$
$\Theta(n)$	$\Theta(2^{\sqrt{\log n}})$

Conversion functions: *list*

A straightforward recursive implementation of *list*.

$$\begin{aligned} \textit{list} &:: (\textit{Array } a) \Rightarrow \textit{Tree } a \ x \rightarrow [x] \\ \textit{list } t & \\ &| \textit{empty } t = [] \\ &| \textit{otherwise} = \textit{elements } t \ ++ \ \textit{riffle } (\textit{list}^* (\textit{subtrees } t)) \\ \textit{riffle} &:: [[x]] \rightarrow [x] \\ \textit{riffle } x & \\ &| \textit{all empty } x = [] \\ &| \textit{otherwise} = \textit{head}^* x \ ++ \ \textit{riffle } (\textit{tail}^* x) \end{aligned}$$

NB. $\textit{riffle} = \textit{concat} \cdot \textit{transpose}$.

Conversion functions: *array*

A generic version of *array*.

```
garray :: (Array a) ⇒ Mix → [x] → Tree a x  
garray ((d, b) :  $\sigma$ ) xs  
  | d == 0           = gnil (b : map snd  $\sigma$ )  
  | otherwise       = node ys ((garray  $\sigma$ )* (unriffle b zs))  
  where (ys, zs) = splitAt d xs  
unriffle :: Int → [x] → [[x]]  
unriffle n xs  
  | empty xs       = replicate n []  
  | otherwise     = cons* ys (unriffle n zs)  
  where (ys, zs) = splitAt n xs
```

Programming challenge

Give linear-time implementations of *list* and *garray*.

Conclusion

- ✘ Bootstrapped arrays are simple to implement.
- ✘ Again, number systems have proven their worth in designing purely functional data structures.
- ✘ One can nicely trade the running time of look-up operations for the running time of update operations.
- ✘ Sensible choices for the base arrays are ‘real’ arrays or lists (‘logarithmic base arrays’ don’t lead to improvements).
- ✘ Preliminary measurements show that bootstrapped arrays perform well for random access, the main factor being the underlying base array.