

FUNCTIONAL PEARL

On merging and selection

RICHARD S. BIRD

*Programming Research Group, Oxford University,
Wolfson Building, Parks Road, Oxford OX1 3QD, UK*

1 Introduction

Given two ascending lists xs and ys of combined length greater than n , consider the computation of

$$\text{merge } (xs, ys) !! n$$

The standard function *merge* merges two ascending sequences and $(!!)$ denotes list indexing. With a lazy evaluator the computation takes $O(n)$ steps; with an eager one it takes $O(p + q)$ steps, where $p = \text{length } xs$ and $q = \text{length } ys$. Now in functional programming it is more efficient to index a tree than a list, so the question arises: can we find a faster solution if xs and ys are each represented by a tree? Somewhat surprisingly the answer is yes: if xs and ys are each represented by balanced binary search trees, then the computation can be reduced to $O(\log p + \log q)$ steps. This is despite the fact that there is no known method for merging two binary search trees in better than linear time. The details, presented below, depend on a subtle relationship between merging and indexing.

2 Specification

We begin by specifying the problem more precisely. Consider the type

$$\mathbf{data} \text{ Tree } a = \text{Null} \mid \text{Node } \text{Int } (\text{Tree } a) a (\text{Tree } a)$$

of binary trees under the restriction that each element xt of $\text{Tree } a$ satisfies the following datatype invariants:

1. The sequence *flatten* xt is ascending, where

$$\begin{aligned} \text{flatten } \text{Null} &= [] \\ \text{flatten } (\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= \text{flatten } xt \ ++ \ [x] \ ++ \ \text{flatten } yt \end{aligned}$$

In words, xt is a binary search tree.

2. If $xt = \text{Node } p \text{ } yt \text{ } y \text{ } zt$, then $p = \text{size } yt$, where

$$\begin{aligned} \text{size } \text{Null} &= 0 \\ \text{size } (\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= \text{size } xt \ + \ 1 \ + \ \text{size } yt. \end{aligned}$$

In words, each node is labelled with the size of its left subtree.

In particular, the following equations hold, referred to subsequently by the hint ‘datatype invariant’:

$$\begin{aligned} \text{take } p (\text{flatten } (\text{Node } p \text{ } xt \text{ } x \text{ } yt)) &= \text{flatten } xt \\ \text{flatten } (\text{Node } p \text{ } xt \text{ } x \text{ } yt) !! p &= x \\ \text{drop } (p + 1) (\text{flatten } (\text{Node } p \text{ } xt \text{ } x \text{ } yt)) &= \text{flatten } yt \end{aligned}$$

Now define

$$\text{select } (xt, yt) n = \text{merge } (\text{flatten } xt, \text{flatten } yt) !! n$$

Our aim is to show that $\text{select } (xt, yt) n$ can be computed in $O(\text{height } xt + \text{height } yt)$ steps, where height is defined by

$$\begin{aligned} \text{height } \text{Null} &= 0 \\ \text{height } (\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= 1 + \max (\text{height } xt, \text{height } yt) \end{aligned}$$

3 Derivation

It is instructive to consider first a simpler problem. Define index by

$$\text{index } xt \ n = (\text{flatten } xt) !! n$$

To synthesise a more efficient program for index we need the following relationship between concatenation and indexing:

$$\begin{aligned} (xs \ ++ \ ys) !! n &= xs !! n, && \text{if } n < p \\ &= ys !! (n - p), && \text{otherwise} \\ &\text{where } p = \text{length } xs \end{aligned}$$

It is now an easy task to synthesise the following alternative program for index :

$$\begin{aligned} \text{index } (\text{Node } p \text{ } xt \text{ } x \text{ } yt) \ n &= \text{index } xt \ n, && \text{if } n < p \\ &= x, && \text{if } n = p \\ &= \text{index } yt \ (n - p - 1), && \text{if } n > p \end{aligned}$$

Evaluation of $\text{index } xt \ n$ takes $O(\text{height } xt)$ steps. In particular, if xt is balanced, then the cost is $O(\log(\text{size } xt))$ steps.

The efficient program for select is derived in an analogous fashion, and depends on the relationship between merging and indexing. Suppose xs and ys are two ascending sequences of combined length greater than n , and let p and q be two natural numbers satisfying $0 \leq p < \text{length } xs$ and $0 \leq q < \text{length } ys$. The facts we need are:

1. If $n \leq p + q$, then

$$\begin{aligned} \text{merge } (xs, ys) !! n &= \text{merge } (xs, \text{take } q \ ys) !! n, && \text{if } xs !! p \leq ys !! q \\ &= \text{merge } (\text{take } p \ xs, ys) !! n, && \text{if } ys !! q \leq xs !! p \end{aligned}$$

2. If $p + q < n$, then

$$\begin{aligned} & \text{merge}(xs, ys) !! n \\ &= \text{merge}(\text{drop}(p + 1) xs, ys) !! (n - p - 1), \quad \text{if } xs !! p \leq ys !! q \\ &= \text{merge}(xs, \text{drop}(q + 1) ys) !! (n - q - 1), \quad \text{if } ys !! q \leq xs !! p \end{aligned}$$

Given these facts, the synthesis of the program for *select* is straightforward. We give the details for just one case. For brevity, introduce the functions

$$\begin{aligned} \text{label}(\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= p \\ \text{left}(\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= xt \\ \text{value}(\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= x \\ \text{right}(\text{Node } p \text{ } xt \text{ } x \text{ } yt) &= yt \end{aligned}$$

Now, with $p = \text{label } xt$ and $q = \text{label } yt$ we argue:

$$\begin{aligned} & \text{select}(xt, yt) n \\ &= \{\text{definition}\} \\ & \text{merge}(\text{flatten } xt, \text{flatten } yt) !! n \\ &= \{\text{property (1), assuming } n \leq p + q \text{ and } \text{value } xt \leq \text{value } yt\} \\ & \text{merge}(\text{flatten } xt, \text{take } q(\text{flatten } yt)) !! n \\ &= \{\text{datatype invariant}\} \\ & \text{merge}(\text{flatten } xt, \text{flatten}(\text{left } yt)) !! n \\ &= \{\text{definition of } \text{select}\} \\ & \text{select}(xt, \text{left } yt) n \end{aligned}$$

In full, the improved program for *select* is:

$$\begin{aligned} \text{select}(xt, yt) n &= \text{index } yt \ n, && \text{if } xt = \text{Null} \\ &= \text{index } xt \ n, && \text{if } yt = \text{Null} \\ &= \text{select}(xt, \text{left } yt) \ n, && \text{if } n \leq p + q \wedge x \leq y \\ &= \text{select}(\text{left } xt, yt) \ n, && \text{if } n \leq p + q \wedge y \leq x \\ &= \text{select}(\text{right } xt, yt) (n - p - 1), && \text{if } p + q < n \wedge x \leq y \\ &= \text{select}(xt, \text{right } yt) (n - q - 1), && \text{if } p + q < n \wedge y \leq x \\ &\text{where } p = \text{label } xt \\ & \quad q = \text{label } yt \\ & \quad x = \text{value } xt \\ & \quad y = \text{value } yt \end{aligned}$$

It is clear that *select* has the desired time complexity: at each step, one or other of the two trees is reduced in height by at least one.

4 Merging and indexing

The relationship between merging and indexing on which the fast algorithm for *select* depends is none too obvious I would say. Certainly the proof seems quite tricky.

In what follows it is convenient to imagine that all ascending sequences are extended on the left by $-\infty$ values, and on the right by ∞ values. Thus, we assume that

$$\begin{aligned} xs !! k &= -\infty, & \text{if } k < 0 \\ &= \infty, & \text{if } \text{length } xs \leq k \end{aligned}$$

The following lemma will be useful.

Lemma 1

Let xs and ys be two ascending sequences and n a natural number. Then there exist unique natural numbers i and j with $i + j = n$ and

$$xs !! (i - 1) \leq ys !! j \quad \text{and} \quad ys !! (j - 1) < xs !! i$$

Furthermore, $\text{merge } (xs, ys) !! n = \min (xs !! i, ys !! j)$.

Proof

The proof is by induction on n . For the base case the unique assignment is $(i, j) = (0, 0)$. For the induction step, suppose (i, j) are the values associated with case n . If $xs !! i \leq ys !! j$, then $(i + 1, j)$ is the unique assignment in case $n + 1$, while if $ys !! j < xs !! i$, the assignment is $(i, j + 1)$.

With the definition

$$\begin{aligned} \text{merge } ([], ys) &= ys \\ \text{merge } (x : xs, []) &= x : xs \\ \text{merge } (x : xs, y : ys) &= x : \text{merge } (xs, y : ys), & \text{if } x \leq y \\ &= y : \text{merge } (x : xs, ys), & \text{otherwise} \end{aligned}$$

of merge it is easy to show, with i and j as given above, that

$$\text{drop } n (\text{merge } (xs, ys)) = \text{merge } (\text{drop } i \text{ } xs, \text{drop } j \text{ } ys)$$

Hence we can argue:

$$\begin{aligned} &\text{merge } (xs, ys) !! n \\ &= \{ \text{since } zs !! k = \text{head } (\text{drop } k \text{ } zs) \} \\ &\quad \text{head } (\text{drop } n (\text{merge } (xs, ys))) \\ &= \{ \text{above} \} \\ &\quad \text{head } (\text{merge } (\text{drop } i \text{ } xs, \text{drop } j \text{ } ys)) \\ &= \{ \text{definition of } \text{merge} \} \\ &\quad \min (xs !! i, ys !! j) \end{aligned}$$

□

Since $\text{merge } (xs, ys) = \text{merge } (ys, xs)$ the lemma has a dual version in which the roles of xs and ys are interchanged. Thus, for property (1) it is sufficient to show that if $n \leq p + q$ and $xs !! p \leq ys !! q$, then

$$\text{merge } (xs, ys) !! n = \text{merge } (xs, \text{take } q \text{ } ys) !! n$$

Let i and j be the numbers associated with xs and ys as specified in the lemma. If

$q < j$, then $i = n - j < n - q \leq p$, and so

$$xs !! i \leq xs !! p \leq ys !! q \leq ys !! (j - 1) \tag{1}$$

This contradicts the definition of i and j , so $j \leq q$. Furthermore, since

$$\begin{aligned} xs !! (i - 1) &\leq ys !! j \leq (take\ q\ ys) !! j \\ (take\ q\ ys) !! (j - 1) &= ys !! (j - 1) < xs !! i \end{aligned}$$

the numbers i and j are also the numbers associated with xs and $take\ q\ ys$.

We now need a case analysis. In the case $j = q$ we have $i \leq p$, and so

$$\begin{aligned} &merge\ (xs,\ ys)\ !!\ n \\ = &\quad \{lemma,\ and\ assumption\ j = q\} \\ &min\ (xs\ !!\ i,\ ys\ !!\ q) \\ = &\quad \{since\ xs\ !!\ i \leq xs\ !!\ p \leq ys\ !!\ q\} \\ &xs\ !!\ i \\ = &\quad \{since\ (take\ q\ ys)\ !!\ q = \infty\} \\ &min\ (xs\ !!\ i,\ (take\ q\ ys)\ !!\ q) \\ = &\quad \{lemma\} \\ &merge\ (xs,\ take\ q\ ys)\ !!\ n \end{aligned}$$

In the case $j < q$ we reason

$$\begin{aligned} &merge\ (xs,\ ys)\ !!\ n \\ = &\quad \{lemma\} \\ &min\ (xs\ !!\ i,\ ys\ !!\ j) \\ = &\quad \{assumption\ j < q\} \\ &min\ (xs\ !!\ i,\ (take\ q\ ys)\ !!\ j) \\ = &\quad \{lemma\} \\ &merge\ (xs,\ take\ q\ ys)\ !!\ n \end{aligned}$$

For property (2) it is sufficient to prove that if $p + q < n$ and $xs !! p \leq ys !! q$, then

$$merge\ (xs,\ ys)\ !!\ n = merge\ (drop\ (p + 1)\ xs,\ ys)\ !!\ (n - p - 1)$$

Again, let i and j be the numbers associated with xs and ys . If $i \leq p$, then $j = n - i \geq n - p > q$ and (1) holds, a contradiction. Thus $i > p$. Since

$$\begin{aligned} (drop\ (p + 1)\ xs)\ !!\ (i - p - 2) &\leq xs !! (i - 1) \leq ys !! j \\ ys !! (j - 1) &< xs !! i = (drop\ (p + 1)\ xs)\ !!\ (i - p - 1) \end{aligned}$$

the unique numbers associated with the sequences $drop\ (p + 1)\ xs$ and ys are $i - p - 1$ and j . Hence

$$\begin{aligned} &merge\ (xs,\ ys)\ !!\ n \\ = &\quad \{lemma\} \end{aligned}$$

$$\begin{aligned}
& \min (xs !! i, ys !! j) \\
= & \quad \{\text{since } i > p\} \\
& \min ((drop (p + 1) xs) !! (i - p - 1), ys !! j) \\
= & \quad \{\text{lemma}\} \\
& \text{merge } (drop (p + 1) xs, ys) !! (n - p - 1)
\end{aligned}$$

This concludes the proof.

5 Postscript

The problem treated in this pearl arose out of a tutorial exercise on divide and conquer set by my colleague, Bill McColl. Students were asked for an algorithm to find the median element of a set represented by two sorted lists, each of length $n > 0$, in $O(\log n)$ steps. Equivalently, the problem is to compute $\text{merge } (xs, ys) !! n$ in $O(\log n)$ steps, assuming that list indexing takes constant time and xs and ys are both strictly increasing and have no elements in common. This variation is left as an exercise. Another variation, also left as an exercise, is to compute the same expression in the same time, assuming constant-time list indexing and that xs and ys are infinite ascending lists (not necessarily increasing, nor necessarily disjoint).

Acknowledgement

I would like to thank two anonymous referees for pointing out errors in previous drafts of this pearl.