

The Knuth-Yao Quadrangle-Inequality Speedup is a Consequence of Total-Monotonicity

Wolfgang W. Bein^{*} Mordecai J. Golin[†] Lawrence L. Larmore[‡] Yan Zhang[§]

Abstract

There exist several general techniques in the literature for speeding up naive implementations of dynamic programming. Two of the best known are the Knuth-Yao quadrangle inequality speedup and the SMAWK algorithm for finding the row-minima of totally monotone matrices. Although both of these techniques use a quadrangle inequality and seem similar they are actually quite different and have been used differently in the literature.

In this paper we show that the Knuth-Yao technique is actually a direct consequence of total monotonicity. As well as providing new derivations of the Knuth-Yao result, this also permits showing how to solve the Knuth-Yao problem directly using the SMAWK algorithm. Another consequence of this approach is a method for solving *online* versions of problems with the Knuth-Yao property. The online algorithms given here are asymptotically as fast as the best previously known static ones. For example the Knuth-Yao technique speeds up the standard dynamic program for finding the optimal binary search tree of n elements from $\Theta(n^3)$ down to $O(n^2)$, and the results in this paper allow construction of an optimal binary search tree in an online fashion (adding a node to the left or right of the current nodes at each step) in $O(n)$ time per step.

1 Introduction

1.1 History

The construction of optimal binary search trees is a classic optimization problem. The input is $2n+1$ weights (probabilities) $p_1, \dots, p_n, q_0, q_1, \dots, q_n$; p_l is the weight that a search is for Key_l ; such a search is called *successful*. The value q_l is the weight that the search argument is *unsuccessful* and is for an argument between Key_l and Key_{l+1} (where we set $\text{Key}_0 = -\infty$ and $\text{Key}_{n+1} = \infty$). Note that we use *weight* instead of *probability* since the p_l and q_l will not be required to add up to 1.

Our problem is to find an *optimal binary search tree* (OBST) with n internal nodes – corresponding to successful searches – and $n+1$ leaves – corresponding to unsuccessful searches – that minimizes the average search time. Let $d(p_l)$ be the depth of internal node corresponding to p_l

^{*}Department of Computer Science, University of Nevada, Las Vegas, NV 89154. Email: bein@cs.unlv.edu. Research supported by NSF grant CCR-0312093.

[†]Dept. of Computer Science, Hong Kong UST, Clear Water Bay, Kowloon, Hong Kong. Email golin@cs.ust.hk Research partially supported by Hong Kong RGC CERG grant HKUST6312/04E.

[‡]Department of Computer Science, University of Nevada, Las Vegas, NV 89154. Email: larmore@cs.unlv.edu. Research supported by NSF grant CCR-0312093.

[§]Dept. of Computer Science, Hong Kong UST, Clear Water Bay, Kowloon, Hong Kong. Email: cszy@cs.ust.hk Research partially supported by Hong Kong RGC CERG grant HKUST6312/04E.

and $d(q_l)$ the depth of leaf corresponding to q_l . Then we want to find a tree that minimizes

$$\sum_{1 \leq l \leq n} p_l(1 + d(p_l)) + \sum_{0 \leq l \leq n} q_l d(q_l).$$

It is not hard to see that this problem reduces to solving the following recurrence using an $O(n^3)$ -time dynamic program:

$$B_{i,j} = \begin{cases} 0 & \text{if } i = j \\ \sum_{l=i+1}^j p_l + \sum_{l=i}^j q_l + \min_{i < t \leq j} \{B_{i,t-1} + B_{t,j}\} & \text{if } i < j \end{cases} \quad (1)$$

where the cost of the optimal OBST is $B_{0,n}$. The naive way of calculating $B_{i,j}$ requires $\Theta(j - i)$ time, so calculating all of the $B_{i,j}$ would seem to require $\Theta(n^3)$ time. In fact, this is what was done in by Gilbert and Moore in 1956 [9]. More than a decade later, in 1971, it was noticed by Knuth [10] that, using a complicated amortization argument, the $B_{i,j}$ can all be computed using only $\Theta(n^2)$ time. Around another decade later, in the early 1980s, Yao [18, 19] simplified Knuth's proof and, in the process, showed that this *dynamic programming speedup* worked for a large class of problems satisfying a *quadrangle inequality* property.

Many other authors then used the Knuth-Yao technique, either implicitly or explicitly, to speed up different dynamic programming problems. See *e.g.*, [15, 3, 4].

In the 1980s a variety of researchers developed various related techniques for exploiting properties such as convexity and concavity to yield dynamic programming speedups; a good early survey is [8]. A high point of this strand of research was the development in the late 1980s of the linear time SMAWK algorithm [1] for finding the row-minima of totally monotone matrices. The work in [7] provides a good survey of the techniques mentioned as well as applications and later extensions. One particular extension we mention (since we will use it later) is the LARSCH algorithm of Larmore and Schieber [11] which, in some cases, permits finding row-minima even when entries of the matrix can implicitly depend upon other entries in the matrix (a case SMAWK cannot handle). Very recently [5] gives new results based on the LARSCH algorithm for certain bottleneck path problems (which extends the earlier work in [11]) and in the the same paper the LARSCH algorithm is used to find a bottleneck-shortest pyramidal traveling salesman tour in $O(n)$ time.

As we shall soon see, both the Knuth-Yao (KY) and SMAWK techniques rely on an underlying quadrangle inequality in their structure and have a similar “feel”. In spite of this, they have until now usually been thought of as being different approaches. See, *e.g.*, [13] which uses *both* KY and SMAWK to speed up different problems. In [2] Aggarwal and Park demonstrated a relationship between the KY problem and totally-monotone matrices by building a *3-D monotone matrix* based on the KY problem and then using an algorithm due to Wilber [16] to find *tube* minima in that 3-D matrix. They left as an open question the possibility of using SMAWK directly to solve the KY problem.

The main theoretical contribution of this paper is to show that the KY technique is really just a special case of the use of totally monotone matrices. We first show a direct solution to the KY problem by decomposing it into $O(n)$ totally-monotone $O(n) \times O(n)$ matrices, permitting direct application of the SMAWK algorithm to yield another $O(n^2)$ solution. After that we describe how the Knuth-Yao technique itself is actually a direct consequence of total-monotonicity of certain related matrices. Finally, we show that problems which can be solved by the KY technique statically in $O(n^2)$ time can actually be solved in an online manner using only $O(n)$ worst case time per step. This is done by using a new formulation of the problem in terms of monotone-

matrices, along with the LARSCH algorithm.¹ We conclude by discussing various extension of the standard KY speedup problem in the literature and showing that these extensions are similarly just special cases of the use of totally monotone matrices.

1.2 Definitions

Definition 1 A two dimensional upper triangular array $a(i, j)$, $0 \leq i \leq j \leq n$ satisfies a quadrangle inequality (QI) if

$$a(i, j) + a(i', j') \leq a(i', j) + a(i, j') \quad \text{for } i \leq i' \leq j \leq j'.$$

Note: In some applications we will write $a_{i,j}$ instead of $a(i, j)$.

Definition 2 A 2×2 matrix is monotone if the minimum of the upper row is not to the right of the minimum of the lower row. More formally, $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ is monotone if $b < a$ implies that $d < c$ and $b = a$ implies that $d \leq c$.

A 2-dimensional matrix M is totally monotone if every 2×2 submatrix of M is monotone.

Definition 3 A 2×2 matrix $\begin{bmatrix} a & b \\ c & d \end{bmatrix}$ is Monge if $a + d \leq b + c$.

A 2-dimensional matrix M is Monge if every 2×2 submatrix² of M is Monge.

The important observations (all of which can be found in [7]) are

Observation 1 An $(m+1) \times (n+1)$ matrix M is Monge if

$$M(i, j) + M(i+1, j+1) \leq M(i, j+1) + M(j+1, i) \quad (2)$$

for all $0 \leq i < m$ and $0 \leq j < n$,

Observation 2 Every Monge matrix is totally monotone.

Combining the above leads to the test that we will often use:

Observation 3 Let M be an $(m+1) \times (n+1)$ matrix. If

$$M(i, j) + M(i+1, j+1) \leq M(i, j+1) + M(j+1, i) \quad (3)$$

for all $0 \leq i < m$ and $0 \leq j < n$, then M is totally monotone.

1.3 Mathematical Framework

Even though both the SMAWK algorithm [1] and the Knuth-Yao (KY) speedup [10, 18, 19] use an implicit quadrangle inequality in their associated matrices, on second glance, they seem quite different from each other.

In the SMAWK technique, the quadrangle inequality is on the entries of a given $m \times n$ input matrix, which can be any totally monotone matrix.³ It is not necessary for the input matrix to

¹We should point out that, as discussed in more detail at the end of Section 3, an alternative online algorithm to the one presented here could be derived by careful deconstruction of the static Aggarwal-Park [2] method; somehow, this never seems to have been remarked before in the literature.

²In this paper, a submatrix is allowed to take non-adjacent rows/columns from the original matrix.

³Note that Monge Matrices satisfy a quadrangle inequality, but in general, a totally monotone matrix may not. However, in practice, most applications of the SMAWK algorithm make use of Monge matrices. If the input matrix is triangular, the missing entries are assigned the default value ∞ , preserving total monotonicity.

actually be given: in many applications, including those in this paper, the entries are implicit, *i.e.*, they are computed only as they are needed. All that the SMAWK algorithm requires is that, when needed, the entries can be calculated in $O(1)$ (amortized) time. The *output* of the SMAWK algorithm is a vector containing the row-minima of the input matrix. If $m \leq n$, the SMAWK algorithm outputs this vector in $O(n)$ time, an order of magnitude speedup of the naive algorithm that scans all mn matrix entries.

The KY technique, by contrast, uses a quadrangle inequality in the upper-triangular $n \times n$ matrix $B_{i,j}$. That is, it uses the QI property of its *result matrix* to speed up the evaluation, via dynamic programming, of the entries in the same result matrix.

More specifically, Yao's result [18] was formulated as follows: For $0 \leq i \leq j \leq n$ let $w(i, j)$ be a given value and

$$B_{i,j} = \begin{cases} 0 & \text{if } i = j \\ w(i, j) + \min_{i < t \leq j} \{B_{i,t-1} + B_{t,j}\} & \text{if } i < j \end{cases} \quad (4)$$

Definition 4 $w(i, j)$ is monotone in the lattice of intervals if $[i, j] \subseteq [i', j']$ implies $w(i, j) \leq w(i', j')$.

As an example, it is not difficult to see that the $w(i, j) = \sum_{l=i+1}^j p_l + \sum_{l=i}^j q_l$ of the BST recurrence (1) satisfies the quadrangle inequality and is monotone in the lattice of intervals.

Definition 5 Let

$$K_B(i, j) = \max\{t : w(i, j) + B_{i,t-1} + B_{t,j} = B_{i,j}\},$$

i.e., the largest index which achieves the minimum in (4).

Yao then proves two Lemmas (see Figure 1 for an example):

Lemma 1 (Lemma 2.1 in [18])

If $w(i, j)$ satisfies the quadrangle inequality as defined in Definition 1, and is also monotone on the lattice of intervals, then the $B_{i,j}$ defined in (4) also satisfy the quadrangle inequality.

Lemma 2 (Lemma 2.2 in [18])

If the function defined in (4) satisfies the quadrangle inequality then

$$K_B(i, j) \leq K_B(i, j+1) \leq K_B(i+1, j+1) \quad \text{for } i < j$$

Lemma 1 proves that a QI in the $w(i, j)$ implies a QI in the $B_{i,j}$. Suppose then that we evaluate the values of the $B_{i,j}$ in the order $d = 1, 2, \dots, n$, where, for each fixed d , we evaluate all of $B_{i,i+d}$, $i = 0, 1, n-d$. Then Lemma 2 says that $B_{i,i+d}$ can be evaluated in time $O(K_B(i+1, i+d) - K_B(i, i+d-1))$. Note that

$$\sum_{i=0}^{n-d} (K_B(i+1, i+d) - K_B(i, i+d-1)) \leq K_B(n-d+1, n) \leq n$$

and thus all entries for fixed d can be calculated in $O(n)$ time. Summing over all d , we see that all $B_{i,j}$ can be obtained in $O(n^2)$ time.

As mentioned, Lemma 2 and the resultant $O(n^2)$ running time have usually been viewed as unrelated to the SMAWK algorithm. While they seem somewhat similar (a QI leading to an order of magnitude speedup) they appeared not to be directly connected.

The main theoretical result of this paper is the observation that if the $w(i, j)$ satisfy the QI and are monotone in the lattice of intervals, then the $B_{i,j}$ defined by (4) can be derived as the row-minima of a sequence of $O(n)$ different totally monotone matrices, each of size $O(n) \times O(n)$,

where the entries in a matrix depend upon the row-minima of previous matrices in the sequence.⁴ In fact, we will show three totally different decomposition of the $B_{i,j}$ into $O(n)$ totally monotone matrices. In particular, our first decomposition, will permit the direct use of SMAWK.

1.4 Online Algorithms

Generally, an online problem is defined to be a problem where a stream of outputs must be generated in response to a stream of inputs, and where those responses must be given under a protocol which requires some outputs be given before all inputs are known. The performance of such an online algorithm is usually compared to the performance of an algorithm that does not have a restriction on the input stream and thus knows the entire input sequence in advance. For many optimization problems it is impossible to give an optimal solution without complete knowledge of future inputs. In such situations online algorithms are analyzed in terms of *competitiveness*, a measure of the performance that compares the decision made online with the optimal offline solution for the same problem, where the lowest possible competitiveness is best.

The online versions of the problems in which we are interested are given below, and do not involve competitiveness. Instead, our goal is to achieve the optimal result, while maintaining the same asymptotic time complexity as the offline versions.

Let $L \leq R$ be given along with values $w(i, j)$ for all $L \leq i \leq j \leq R$ that satisfy the QI and the “monotone on lattice of intervals” property. Let

$$B_{i,j} = \begin{cases} 0 & \text{if } i = j \\ w(i, j) + \min_{i < t \leq j} \{B_{i,t-1} + B_{t,j}\} & \text{if } i < j \end{cases}$$

and assume all $B_{i,j}$ for $L \leq i \leq j \leq R$ have already been calculated and stored.

The **Right-online** problem is:

Given new values $w(i, R+1)$ for $L \leq i \leq R+1$, such that $w(i, j)$ still satisfy the QI and monotone property, calculate all of the values $B_{i,R+1}$ for $L \leq i \leq R+1$.

The **Left-online** problem is:

Given new values $w(L-1, j)$ for $L-1 \leq j \leq R$, such that $w(i, j)$ still satisfy the QI and monotone property, calculate all of the values $B_{L-1,j}$ for $L-1 \leq j \leq R$.

These online problems restricted to the optimal binary search tree would be to construct the OBST for items $\text{Key}_{L+1}, \dots, \text{Key}_R$, and, at each step, add either Key_{R+1} , a new key to the right, or Key_L , a new key to the left. Every time a new element is added, we want to update the $B_{i,j}$ (dynamic programming) table and thereby construct the optimal binary search tree of the new full set of elements. (See Figure 1.) To achieve this, it is certainly possible to *recompute* the entire table; however this comes at the price of $O(n^2)$ time, where $n = R - L$ is the number of keys currently in the table. What we are interested in here is the question of how one can handle a new key where the extra computational work is neutral to the overall complexity of the problem, *i.e.*, a new key can be added in linear time. Our goal is an algorithm in which a sequence of n online key insertions will result in a worst case $O(n)$ per step to maintain an optimal tree, yielding an

⁴See [17] for a solution to a different type of problem via finding the row-minima of a sequence of dependent totally monotone matrices.

	3	4	5	6	7
3	0	91	282	499	821
4		0	169	386	686
5			0	124	348
6				0	155
7					0

	3	4	5	6	7
3	3	4	5	5	6
4		4	5	5	6
5			5	6	7
6				6	7
7					7

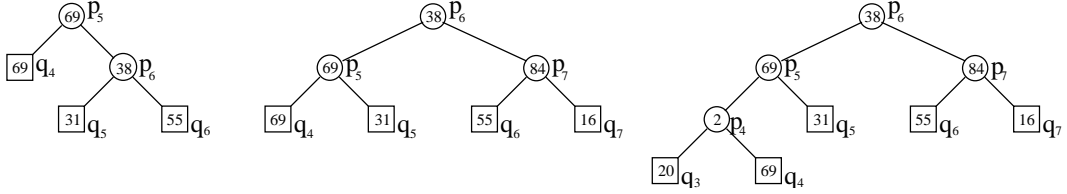


Figure 1: An example of the online case for optimal binary search trees where $(p_4, p_5, p_6, p_7) = (2, 69, 38, 84)$ and $(q_3, q_4, q_5, q_6, q_7) = (20, 69, 31, 55, 16)$. The leftmost table contains the $B_{i,j}$ values; the rightmost one, the $K_B(i, j)$ values. The unshaded entries in the table are for the problem restricted to only keys 5, 6. The dark gray cells are the entries added to the table when key 7 is added to the right. The light gray cells are the entries added when key 4 is added to the left. The corresponding optimal binary search trees are also given, where circles correspond to successful searches and squares to unsuccessful ones. The values in the nodes are the weights of the nodes (not their keys).

overall run time of $O(n^2)$.

Unfortunately, the KY speedup *cannot* be used to do this. The reason that the speedup fails is that the KY speedup is actually an amortization over the evaluation of all entries when done in a particular order. In the online case, adding a new item n to previously existing items $1, 2, \dots, n-1$, requires using (4) to compute the n new entries $B_{i,n}$, in the fixed order $i = n, n-1, \dots, 1, 0$ and it is not difficult to construct an example in which calculating these new entries in this order using (4) requires $\Theta(n^2)$ work.

We will see later that the decomposition given in section 3 permits a fully online algorithm with no penalty in performance, *i.e.*, after adding the n^{th} new key, the new $B_{i,j}$ can be calculated in $O(n)$ worst case time. Furthermore, this will be true for both the left-online and right-online case.

2 The First Decomposition

Definition 6 For $1 \leq d < n$ define the $(n-d+1) \times (n+1)$ matrix D^d by

$$D_{i,j}^d = \begin{cases} w(i, i+d) + B_{i,j-1} + B_{j,i+d} & \text{if } 0 \leq i < j \leq i+d \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (5)$$

Figure 2 illustrates the first decomposition. Note that (4) immediately implies

$$B_{i,i+d} = \min_{0 \leq j \leq n} D_{i,j}^d \quad (6)$$

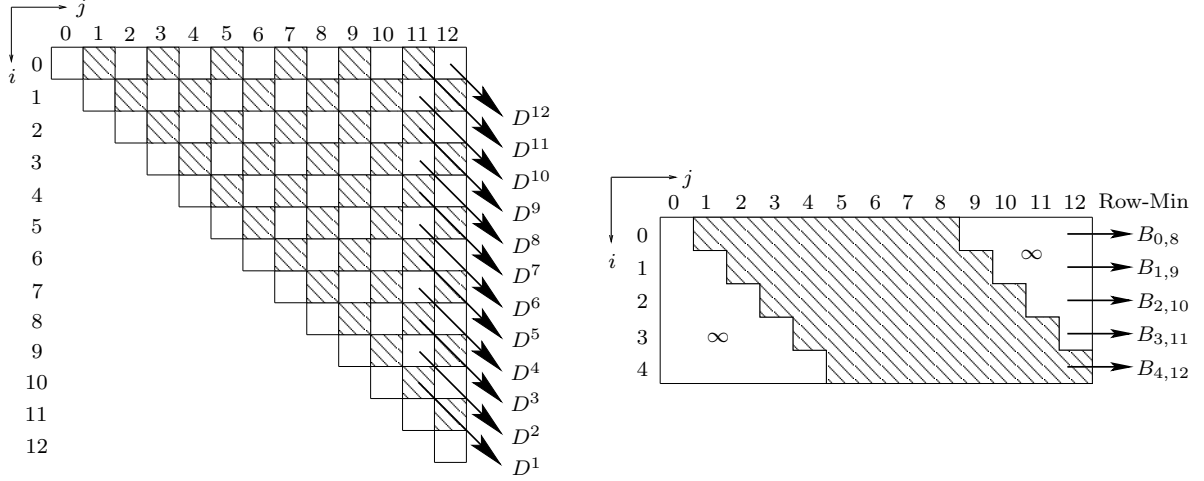


Figure 2: The lefthand figure shows the $B_{i,j}$ matrix for $n = 12$. Each diagonal, $d = i - j$, in the matrix will correspond to a totally monotone matrix D^d . The minimal item of row i in D^d will be the value $B_{i,i+d}$. The righthand figure shows D^8 .

so finding the row-minima of D^d yields $B_{i,i+d}$, $i = 0, \dots, n - d$. Put another way, the $B_{i,j}$ entries on diagonal $j - i = d$ are exactly the row-minima of matrix D^d .

Lemma 3 *If $w(i, j)$ and the function $B_{i,j}$ defined in (4) satisfies the QI then, for each $d \leq n$, D^d is a totally monotone matrix.*

Proof: From Observation 3 it suffices to prove that

$$D_{i,j}^d + D_{i+1,j+1}^d \leq D_{i+1,j}^d + D_{i,j+1}^d \quad (7)$$

Note that if $i + 1 < j < i + d$, then from Lemma 1,

$$B_{i,j-1} + B_{i+1,j} \leq B_{i+1,j-1} + B_{i,j} \quad (8)$$

and

$$B_{j,i+d} + B_{j+1,i+1+d} \leq B_{j,i+1+d} + B_{j+1,i+d}. \quad (9)$$

Thus,

$$\begin{aligned} D_{i,j}^d + D_{i+1,j+1}^d &= [w(i, i+d) + B_{i,j-1} + B_{j,i+d}] + [w(i+1, i+d+1) + B_{i+1,j} + B_{j+1,i+1+d}] \\ &= w(i, i+d) + w(i+1, i+d+1) + [B_{i,j-1} + B_{i+1,j}] + [B_{j,i+d} + B_{j+1,i+1+d}] \\ &\leq w(i, i+d) + w(i+1, i+d+1) + [B_{i+1,j-1} + B_{i,j}] + [B_{j,i+1+d} + B_{j+1,i+d}] \\ &= [w(i+1, i+d+1) + B_{i+1,j-1} + B_{j,i+1+d}] + [w(i, i+d) + B_{i,j} + B_{j+1,i+d}] \\ &= D_{i+1,j}^d + D_{i,j+1}^d \end{aligned}$$

and (7) is correct (where we note that the right hand side is ∞ if $i + 1 \not\leq j$ or $j \not\leq i + d$). \square

Lemma 4 *Assuming that all of the row-minima of D^1, D^2, \dots, D^{d-1} have already been calculated, all of the row-minima of D^d can be calculated using the SMAWK algorithm in $O(n)$ time.*

Proof: From the previous lemma, D^d is a totally monotone matrix. Also, by definition, its entries can be calculated in $O(1)$ time, using the previously calculated row-minima of $D^{d'}$ where $d' < d$. Thus SMAWK can be applied. \square

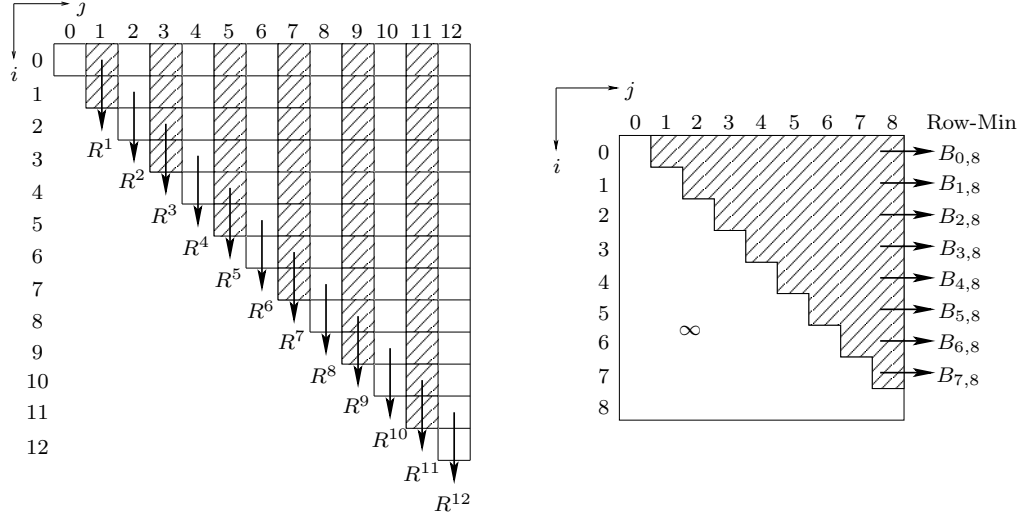


Figure 3: The lefthand figure shows the $B_{i,j}$ matrix for $n = 12$. Each column in the $B_{i,j}$ matrix will correspond to a totally monotone matrix R^m . The minimal element of row i in R^m will be the value $B_{i,m}$. The righthand figure shows R^8 .

Combined with (6) this immediately gives a new $O(n^2)$ algorithm for solving the KY problem; just run SMAWK on the D^d in the order $d = 1, 2, \dots, n - 1$ and report all of the row-minima.

We point out that this technique cannot help us solve the online problem as defined in subsection 1.4, though. To see why, suppose that items $1, \dots, n$ have previously been given, new item $n + 1$ has just been added, and we need to calculate the values $B_{i,n+1}$ for $i = 0, \dots, n + 1$. In our formulation this would correspond to *adding a new bottom row to every matrix D^d and creating a new matrix D^{n+1}* . In our formulation, we would need to find the row-minima of all of the n new bottom rows. Unfortunately, the SMAWK algorithm only works on the rows of matrices all at once and cannot help to find the row-minima of a single new row.

3 The Second & Third Decompositions

So far we have seen that it is possible to derive the KY *running time* via repeated calls to the SMAWK algorithm. We now see two more decompositions into totally-monotone matrices. These decompositions will trivially imply Lemma 2 (Lemma 2.1 in [18]), which is the basis of the KY speedup. Thus, the KY speedup is just a consequence of total-monotonicity. These new decompositions will also permit us to efficiently solve the online problem given in subsection 1.4.

The second decomposition is indexed by the *rightmost element* seen so far. See Figure 3.

Definition 7 For $1 \leq m \leq n$ define the $(m + 1) \times (m + 1)$ matrix R^m by

$$R^m_{i,j} = \begin{cases} w(i, m) + B_{i,j-1} + B_{j,m} & \text{if } 0 \leq i < j \leq m, \\ \infty & \text{otherwise.} \end{cases} \quad (10)$$

Note that (4) immediately implies

$$B_{i,m} = \min_{0 < j \leq m} R^m_{i,j} \quad (11)$$

so finding the row-minima of R^m yields $B_{i,m}$ for $i = 0, \dots, m - 1$. Put another way, the $B_{i,j}$ entries in column m are exactly the row minima of R^m .

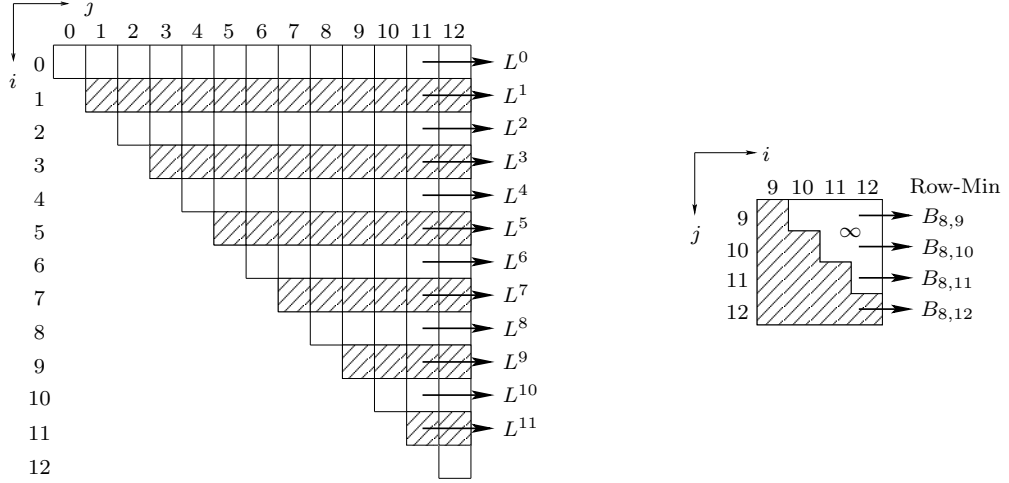


Figure 4: The lefthand figure shows the $B_{i,j}$ matrix for $n = 12$. Each row in the $B_{i,j}$ matrix will correspond to a totally monotone matrix L^m . The minimal element of row j in L^m will be the value $B_{m,j}$. The righthand figure shows L^8 .

The third decomposition is similar to the second except that it is indexed by the *leftmost element* seen so far. See Figure 4.

Definition 8 For $0 \leq m < n$ define the $(n - m) \times (n - m)$ matrix L^m by

$$L_{j,i}^m = \begin{cases} w(m, j) + B_{m,i-1} + B_{i,j} & \text{if } m < i \leq j \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (12)$$

(For convenience, we set the row and column indices to run from $(m + 1) \dots n$ and not $0 \dots (n - m - 1)$.) Note that (4) immediately implies

$$B_{m,j} = \min_{m < i \leq j} L_{j,i}^m \quad (13)$$

so finding the row-minima of L^m yields $B_{m,j}$ for $j = m + 1, \dots, n$. Put another way, the $B_{i,j}$ entries in row m are exactly the row minima of matrix L^m .

Lemma 5 If the function defined in (4) satisfies the QI then R^m and L^m are totally monotone matrices.

Proof: The proofs are very similar to that of Lemma 3. Note that if $i + 1 < j \leq m$, we can again use (8); writing the entries from (8) in boldface gives

$$\begin{aligned} R_{i,j}^m + R_{i+1,j+1}^m &= [w(i, m) + \mathbf{B}_{i,j-1} + B_{j,m}] + [w(i + 1, m) + \mathbf{B}_{i+1,j} + B_{j+1,m}] \\ &\leq [w(i + 1, m) + \mathbf{B}_{i+1,j-1} + B_{j,m}] + [w(i, m) + \mathbf{B}_{i,j} + B_{j+1,m}] \\ &= R_{i+1,j}^m + R_{i,j+1}^m \end{aligned}$$

and thus R^m is Monge (where we note that the right hand side is ∞ if $i + 1 \not\leq j$) and thus totally monotone. If $m < i < j$ then we again use (8) (with j replaced by $j + 1$) to get

$$\begin{aligned} L_{j,i}^m + L_{j+1,i+1}^m &= [w(m, j) + B_{m,i-1} + \mathbf{B}_{i,j}] + [w(m, j + 1) + B_{m,i} + \mathbf{B}_{i+1,j+1}] \\ &\leq [w(m, j + 1) + B_{m,i-1} + \mathbf{B}_{i,j+1}] + [w(m, j) + B_{m,i} + \mathbf{B}_{i+1,j}] \\ &= L_{j+1,i}^m + L_{j,i+1}^m \end{aligned}$$

and thus L^m is Monge (where we note that the right hand side is ∞ if $i \not\leq j$) and thus totally monotone. \square

We point out these two decompositions immediately imply a new proof of Lemma 2 (Lemma 2.1 in [18]) which states that

$$K_B(i, j) \leq K_B(i, j+1) \leq K_B(i+1, j+1). \quad (14)$$

To see this note that $K_B(i, j+1)$ is the location of the rightmost row-minimum of row i in matrix R^{j+1} , while $K_B(i+1, j+1)$ is the location of the rightmost row-minimum of row $i+1$ in matrix R^{j+1} . Thus, the definition of total monotonicity (Definition 2) immediately gives

$$K_B(i, j+1) \leq K_B(i+1, j+1). \quad (15)$$

Similarly, $K_B(i, j)$ is the rightmost row-minimum of row j in L^i while $K_B(i, j+1)$ is the location of the rightmost row-minimum of row $j+1$ in L^i . Thus

$$K_B(i, j) \leq K_B(i, j+1). \quad (16)$$

Combining (15) and (16) yields (14), which is what we want. Since the actual speedup in the KY technique comes from an amortization argument based on (14), we have just seen that the original KY-speedup itself is also a consequence of total monotonicity.

We have still not seen how to actually calculate the $B_{i,j}$ using the R^m and L^m . Before continuing, we point out that even though the R^m are totally monotone, their row minima *cannot* be calculated using the SMAWK algorithm. This is because, for $0 \leq i < j \leq m$, the value of entry $R_{i,j}^m = w(i, m) + B_{i,j-1} + B_{j,m}$, which is dependent upon $B_{j,m}$ which is itself the row-minimum of row j in the same matrix R^m . Thus, the values of the entries of R^m depend upon the other entries in R^m which is something that SMAWK does not allow. The same problem occurs with the L^m .

We will now see that, despite this dependence, we can still use the LARSCH algorithm to find the row-minima of the R^m . This will have the added advantage of solving the online problem as well.

At this point we should note that our decompositions L^m could also be derived by careful cutting of the 3-D monotone matrices of Aggarwal and Park [2] along particular planes. Aggarwal and Park used an algorithm of Wilber [16] (derived for finding the maxima of certain concave-sequences) to find various tube maxima of their matrices, leading to another $O(n^2)$ algorithm for solving the KY-problem. In fact, even though their algorithm was presented as a static algorithm, careful decomposition of what they do permits using it to solve what we call the left-online KY-problem. A symmetry argument could then yield a right-online algorithm. This never seems to have been noted in the literature, though. In the next section, we present a different online algorithm, based on our decompositions and the LARSCH algorithm.

4 Online Algorithms Without Losing the KY Speedup

To execute the LARSCH algorithm, as defined in Section 3 of [11] we need only that X satisfy the following conditions:

1. X is a totally monotone $n \times m$ monotone matrix.
2. For each row index i of X , there is a column index C_i such that for $j > C_i$, $X_{i,j} = \infty$. Furthermore, $C_i \leq C_{i+1}$.

3. If $j \leq C_i$, then $X_{i,j}$ can be evaluated in $O(1)$ time provided that the row minima of the first $i - 1$ rows are already known.

If these conditions are satisfied, the LARSCH algorithm then calculates all of the row minima of X in $O(n + m)$ time. We can now use this algorithm to derive

Lemma 6

- Given that all values $B_{i,j}$, $m < i \leq j \leq n$ have already been calculated, all of the row-minima of L^m can be calculated in $O(n - m)$ time.
- Given that all values $B_{i,j}$, $0 \leq i \leq j < m$ have already been calculated, all of the row-minima of R^m can be calculated in $O(m)$ time.

Proof: For the first part, it is easy to see that L^m satisfies the first two conditions required by the LARSCH algorithm with $C_j = j$. For the third condition, note that, for $m < i \leq j$, $L_{j,i}^m = w(m, j) + B_{m,i-1} + B_{i,j}$. The values $w(m, j)$ and $B_{i,j}$ are already known and can be retrieved in $O(1)$ time. $B_{m,i-1}$ is the row minima of row $i - 1$ of L^m but, since we are assuming $i \leq j$ this means that $B_{m,i-1}$ is the row minima of an earlier row in L^m and the third LARSCH condition is satisfied. Thus, all of the row-minima of the $(n - m) \times (n - m)$ matrix L^m can be calculated in $O(n - m)$ time. For the second part set X to be the $(m + 1) \times (m + 1)$ matrix defined by $X_{i,j} = R_{m-i,m-j}^m$. Then X satisfies the first two LARSCH conditions with $C_i = i - 1$. For the third condition note that $X_{i,j} = R_{m-i,m-j}^m = w(m - i, m) + B_{m-i,m-j-1} + B_{m-j,m}$. The values $w(m - i, m)$ and $B_{m-i,m-j-1}$ are already known and can be calculated in $O(1)$ time. $B_{m-j,m}$ is the row minima of row j of X ; but, since we are assuming $j \leq C_i = i - 1$ this means that $B_{m-j,m}$ is the row minima of an earlier row in X so the third LARSCH condition is satisfied. Thus, all of the row-minima of X and equivalently R^m can be calculated in $O(m)$ time. \square

Note that Lemma 6 immediately solves the “right-online” and “left-online” problems described in subsection 1.4; Given the new values $w(i, R + 1)$ for $L \leq i \leq R + 1$, simply find the row minima of R^{R+1} in time $O(R - L)$. Given the new values $w(L - 1, j)$ for $L - 1 \leq j \leq R$, simply find the row minima of L^{L-1} .

We have therefore just shown that *any* dynamic programming problem for which the KY speedup can statically improve run time from $\Theta(n^3)$ to $O(n^2)$ time can be solved in an online fashion in $O(n)$ time per step. That is, online processing incurs no penalty compared to static processing. In particular, the optimum binary search tree (as illustrated in 1.4), can be maintained in $O(n)$ time per step as nodes are added to both its left and right.

5 Further Applications

In this section, we consider two extensions of the Knuth-Yao quadrangle inequality; The first was due to Wachs [14] in 1989 and the second to Borchers and Gupta (BG) [6] in 1994.

In our presentation we will first quickly describe the Wachs and BG extensions. We then sketch how our various results, *i.e.*, the D^d , R^m and L^m Monge matrix decompositions and their consequences, can be generalized to work for the Wachs and BG extensions.

Note: In order to maintain the consistency of our presentation we sometimes slightly modify the statements of the theorems in [14] and [6]. After our presentations we will note the various modifications made (with the exception of trivial renaming of variables).

5.1 The Wachs Extension

In [14] Wachs was interested in solving the *system* of dynamic programming recurrences

$$B_{i,j} = \begin{cases} 0 & \text{if } i = j \\ v(i,j) + \min_{i < t \leq j} \{u(i,t-1)w(i,j) + \bar{B}_{i,t-1} + B_{t,j}\} & \text{if } i < j \end{cases} \quad (17)$$

$$\bar{B}_{i,j} = \begin{cases} 0 & \text{if } i = j \\ v(i,j) + \min_{i < t \leq j} \{u(t,j)w(i,j) + \bar{B}_{i,t-1} + B_{t,j}\} & \text{if } i < j \end{cases} \quad (18)$$

where $v(i,j)$, $u(i,j)$ and $w(i,j)$ were functions satisfying the QI and other special properties.

Her motivation was to calculate the binary search tree corresponding to the optimal comparison search procedure on a tape. A tape can only be accessed sequentially, either from left to right, or from right to left (with no cost imposed for changing the direction). The n records on the tape are sorted in increasing order by their keys $\text{Key}_1, \dots, \text{Key}_n$. As in Knuth's original problem, denote by p_l the weight that a search is for Key_l , and by q_l the weight that a search is for an argument between Key_l and Key_{l+1} . $\text{Key}_0 = -\infty$ and $\text{Key}_{n+1} = +\infty$. Denote by x_l the location of Key_l on the tape. Let $x_0 = x_1$ and $x_{n+1} = x_n$. The cost of moving the tape from Key_{l_1} to Key_{l_2} ($l_1 \leq l_2$) is the same as moving from Key_{l_2} to Key_{l_1} , which is $a(x_{l_2} - x_{l_1}) + b$ where a and b are nonnegative constants. The binary search tree constructed by Wachs lies in the random access memory but is only used to *model* the search procedure on the tape. This means, the node that *represents* Key_l in the binary search tree does not contain the key value of Key_l , but only contains x_l and the two child pointers. So, in the search step at Key_l , we need to move the tape from current location to x_l , then compare with Key_l , and then decide whether to choose the left or right branch in the binary search tree.

Define

$$u(i,j) = \begin{cases} a(x_{j+1} - x_i) + b & \text{if } j \geq i - 1 \\ \infty & \text{otherwise} \end{cases} \quad (19)$$

which is the cost of moving the tape from Key_i to Key_{j+1} when $j \geq i - 1$. Define

$$w(i,j) = \begin{cases} \sum_{l=i+1}^j p_l + \sum_{l=i}^j q_l & \text{if } i \leq j \\ -p_i & \text{if } j = i - 1 \\ \infty & \text{otherwise} \end{cases} \quad (20)$$

which is the weight of the subtree from Key_{i+1} to Key_j when $i \leq j$, as Knuth's original problem. Both $u(i,j)$ and $w(i,j)$ satisfy the QI. It is important to note in this section (Section 5.1) that $u(i,j)$ and $w(i,j)$ satisfy the QI *as equality* on their finite elements. That is, for $i \leq i'$, $j \leq j'$ and $j \geq i' - 1$,

$$u(i,j) + u(i',j') = u(i',j) + u(i,j') \quad (21)$$

$$w(i,j) + w(i',j') = w(i',j) + w(i,j') \quad (22)$$

We call (21) and (22) the quadrangle *equality* (QE), and we will say $u(i,j)$ and $w(i,j)$ “satisfy the QE” instead of saying “satisfy the QE on finite elements” since the infinite elements can be defined such that the QE is satisfied on all elements.

Let $B_{i,j}$ (resp. $\bar{B}_{i,j}$) be the optimal cost of searching the subtree from Key_{i+1} to Key_j , where the tape is initially at x_i (resp. x_{j+1}). Wachs showed that $B_{i,j}$ and $\bar{B}_{i,j}$ satisfy (17) and (18), when $u(i,j)$ and $w(i,j)$ are defined as (19) and (20), and $v(i,j) \equiv 0$.

The naive method of evaluating all of the $B_{i,j}$ and $\bar{B}_{i,j}$ requires $\Theta(n^3)$ time. Using a generalization of the KY speedup Wachs was able to reduce this down to $O(n^2)$. In our notation, her main results were

Lemma 7 (Theorem 3.1 in [14])

If (i) $v(i, j)$ satisfies the QI, (ii) $u(i, j)$ and $w(i, j)$ satisfy the QE, (iii) all three functions are monotone on the lattice of intervals and, furthermore, (iv) if $u(i, i-1) = b$ is a nonnegative constant independent of i , and $w(i, j) \geq 0$ for all $i \leq j$, then $B_{i,j}$ and $\bar{B}_{i,j}$ as defined by (17) and (18) satisfy the following stronger version of the QI: For all $0 \leq i \leq i' \leq j \leq j' \leq n$,

$$B_{i',j} + B_{i,j'} - B_{i,j} - B_{i',j'} \geq [u(i, i'-1) - u(i', i'-1)][w(j+1, j') - w(j+1, j)] \geq 0 \quad (23)$$

$$\bar{B}_{i',j} + \bar{B}_{i,j'} - \bar{B}_{i,j} - \bar{B}_{i',j'} \geq [u(j+1, j') - u(j+1, j)][w(i, i'-1) - w(i', i'-1)] \geq 0 \quad (24)$$

Lemma 8 (Theorem 3.2 in [14])

If $B_{i,j,r}$ and $\bar{B}(i, j, r)$ satisfy the QI, then

$$K_B(i, j) \leq K_B(i, j+1) \leq K_B(i+1, j+1) \quad (25)$$

$$K_{\bar{B}}(i, j) \leq K_{\bar{B}}(i, j+1) \leq K_{\bar{B}}(i+1, j+1) \quad (26)$$

where $K_B(i, j)$ and $K_{\bar{B}}(i, j)$ are the maximum splitting points at which $B_{i,j}$ and $\bar{B}_{i,j}$, respectively, attain their minimum values.

Lemma 8 is then used in exactly the same fashion as was Lemma 2 by Knuth and Yao, to speed up the solution of the DP recurrence from $\Theta(n^3)$ to $O(n^2)$. Since the $u(i, j)$ and $w(i, j)$, as well as $(v(i, j) \equiv 0)$ in Wachs' tape searching problem satisfy all of the conditions of Lemma 7, this solves Wachs' motivating problem in $O(n^2)$ time.

Setting $u(i, j) \equiv 0$ or $w(i, j) \equiv 0$ gives $B(i, j) = \bar{B}_{i,j}$ for all i, j . Further setting $v(i, j)$ be the $w(i, j)$ in (4) collapses Wachs' results down to the standard KY speedup. Thus, Wachs' results can be seen as an extension of KY.

Note: The indices here are slightly shifted from those in Wachs' [14]. The statement of Theorem 3.1 in [14] assumes that $v(i, j) \equiv 0$. The extension to arbitrary $v(i, j)$ satisfying the QI and monotonicity is noted in the last paragraph of [14].

We now apply our schemes to the system of recurrences of Wachs' problem. We first provide the analogue to our old D^d matrices.

Definition 9 For $1 \leq d \leq n$, define the $(n-d+1) \times (n+1)$ matrix D^d and \bar{D}^d by

$$D_{i,j}^d = \begin{cases} v(i, i+d) + u(i, j-1)w(i, i+d) + \bar{B}_{i,j-1} + B_{j,i+d} & \text{if } 0 \leq i < j \leq i+d \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (27)$$

$$\bar{D}_{i,j}^d = \begin{cases} v(i, i+d) + u(j, i+d)w(i, i+d) + \bar{B}_{i,j-1} + B_{j,i+d} & \text{if } 0 \leq i < j \leq i+d \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (28)$$

Lemma 9 If $B_{i,j}$ and $\bar{B}_{i,j}$ both satisfy the stronger QI given by (23) and (24) in Lemma 7, $u(i, j)$ and $w(i, j)$ both satisfy the QE and are monotone, and $u(i, i-1) = b$ is a nonnegative constant independent of i , then D^d and \bar{D}^d as defined by (27) and (28) are Monge, i.e., for all $1 \leq d \leq n$, $0 \leq i < n-d$ and $0 \leq j < n$,

$$D_{i,j}^d + D_{i+1,j+1}^d \leq D_{i+1,j}^d + D_{i,j+1}^d \quad (29)$$

$$\bar{D}_{i,j}^d + \bar{D}_{i+1,j+1}^d \leq \bar{D}_{i+1,j}^d + \bar{D}_{i,j+1}^d \quad (30)$$

Proof: Since (29) and (30) are symmetric, we will only show the proof of (29). If $i + 1 \not\leq j$ or $j \not\leq i + d$, (29) is trivially true since the right hand side is ∞ . So we assume $i + 1 < j < i + d$. To save space, we write $f(i, j)$ as $f_{i,j}$, where f is v , u or w . Define

$$H_{i,j} = v_{i,i+d} + u_{i,j-1}w_{i,i+d}. \quad (31)$$

Then, $D_{i,j}^d = H_{i,j} + \bar{B}_{i,j-1} + B_{j,i+d}$. From Lemma 7,

$$\bar{B}_{i+1,j-1} + \bar{B}_{i,j} - \bar{B}_{i,j-1} - \bar{B}_{i+1,j} \geq (u_{j,j} - b)(w_{i,i} - w_{i+1,i}) \quad (32)$$

$$B_{j,i+d+1} + B_{j+1,i+d} - B_{j,i+d} - B_{j+1,i+d+1} \geq (u_{j,j} - b)(w_{i+d+1,i+d+1} - w_{i+d+1,i+d}) \quad (33)$$

Denote by $\text{QI}(f; i, i', j, j')$ the QI that $f_{i,j} + f_{i',j'} \leq f_{i',j} + f_{i,j'}$ where $i \leq i'$, $j \leq j'$ and $j \geq i - 1$, and by $\text{QE}(f; i, i', j, j')$ the corresponding QE (QI as equality).

$$\begin{aligned} & H_{i+1,j} + H_{i,j+1} - H_{i,j} - H_{i+1,j+1} \\ &= (u_{i,j} - u_{i,j-1})w_{i,i+d} - (u_{i+1,j} - u_{i+1,j-1})w_{i+1,i+d+1} \\ &\geq (u_{i,j} - u_{i,j-1})(w_{i,i+d} - w_{i+1,i+d+1}) \quad [\text{QI}(u; i, i+1, j-1, j)] \\ &= (u_{j,j} - u_{j,j-1})(w_{i,i+d} - w_{i+1,i+d+1}) \quad [\text{QE}(u; i, j, j-1, j)] \\ &= (u_{j,j} - b)(w_{i,i+d} - w_{i+1,i+d+1}) \quad [u_{j,j-1} = b] \end{aligned} \quad (34)$$

Combine (32) to (34),

$$\begin{aligned} & D_{i+1,j}^d + D_{i,j+1}^d - D_{i,j}^d - D_{i+1,j+1}^d \\ &= (\bar{B}_{i+1,j-1} + \bar{B}_{i,j} - \bar{B}_{i,j-1} - \bar{B}_{i+1,j}) + (B_{j,i+d+1} + B_{j+1,i+d} - B_{j,i+d} - B_{j+1,i+d+1}) + \\ & \quad (H_{i+1,j} + H_{i,j+1} - H_{i,j} - H_{i+1,j+1}) \\ &\geq (u_{j,j} - b)[w_{i,i} - w_{i+1,i} + w_{i,i+d} + (w_{i+d+1,i+d+1} - w_{i+d+1,i+d} - w_{i+1,i+d+1})] \\ &= (u_{j,j} - b)(w_{i,i} - w_{i+1,i} + w_{i,i+d} - w_{i+1,i+d}) \quad [\text{QE}(w; i+1, i+d+1, i+d, i+d+1)] \\ &\geq u_{j,j} - b \quad [\text{monotonicity of } w] \\ &\geq 0 \quad [\text{monotonicity of } u] \end{aligned} \quad (35)$$

which yields the lemma. \square

Note that

$$B_{i,i+d} = \min_{0 \leq j \leq n} D_{i,j}^d \quad \text{and} \quad \bar{B}_{i,i+d} = \min_{0 \leq j \leq n} \bar{D}_{i,j}^d. \quad (36)$$

Thus, as in Section 2, we can use the SMAWK algorithm to evaluate all of the $B_{i,j}$ and $\bar{B}_{i,j}$ in $O(n^2)$ time.

We now generalize the R^m and L^m matrices.

Definition 10 For $1 \leq m \leq n$ define the $(m+1) \times (m+1)$ matrix R^m and \bar{R}^m by

$$R_{i,j}^m = \begin{cases} v(i, m) + u(i, j-1)w(i, m) + \bar{B}_{i,j-1} + B_{j,m} & \text{if } 0 \leq i < j \leq m, \\ \infty & \text{otherwise.} \end{cases} \quad (37)$$

$$\bar{R}_{i,j}^m = \begin{cases} v(i, m) + u(j, m)w(i, m) + \bar{B}_{i,j-1} + B_{j,m} & \text{if } 0 \leq i < j \leq m, \\ \infty & \text{otherwise.} \end{cases} \quad (38)$$

For $0 \leq m < n$ define the $(n-m) \times (n-m)$ matrix L^m and \bar{L}^m by

$$L_{i,j}^m = \begin{cases} v(m, j) + u(m, i-1)w(m, j) + \bar{B}_{m,i-1} + B_{i,j} & \text{if } m < i \leq j \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (39)$$

$$\bar{L}_{i,j}^m = \begin{cases} v(m,j) + u(i,j)w(m,j) + \bar{B}_{m,i-1} + B_{i,j} & \text{if } m < i \leq j \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (40)$$

Lemma 10 *If $B_{i,j}$ and $\bar{B}_{i,j}$ both satisfy the stronger QI given by (23) and (24) in Lemma 7, $u(i,j)$ and $w(i,j)$ both satisfy the QE and are monotone, and $u(i,i-1) = b$ is a nonnegative constant independent of i , then the four matrices R^m , \bar{R}^m , L^m and \bar{L}^m as defined by (37) to (40) are all Monge matrices.*

Proof: We will only show the proof for R and \bar{R} ; the proof for L and \bar{L} is symmetric. We will show for all $1 \leq m \leq n$, $0 \leq i < m$ and $0 \leq j < m$,

$$R_{i,j}^m + R_{i+1,j+1}^m \leq R_{i+1,j}^m + R_{i,j+1}^m \quad (41)$$

$$\bar{R}_{i,j}^m + \bar{R}_{i+1,j+1}^m \leq \bar{R}_{i+1,j}^m + \bar{R}_{i,j+1}^m \quad (42)$$

If $i+1 \not\leq j$, (41) and (42) are trivially true since the right hand side is ∞ . So we assume $i+1 < j$. The proof of (41) is easier than that of (42).

$$\begin{aligned} & R_{i+1,j}^m + R_{i,j+1}^m - R_{i,j}^m - R_{i+1,j+1}^m \\ &= (\bar{B}_{i+1,j-1} + \bar{B}_{i,j} - \bar{B}_{i,j-1} - \bar{B}_{i+1,j}) + (v_{i+1,m} + v_{i,m} - v_{i,m} - v_{i+1,m}) + \\ & \quad (u_{i+1,j-1}w_{i+1,m} + u_{i,j}w_{i,m} - u_{i,j-1}w_{i,m} - u_{i+1,j}w_{i+1,m}) \\ &\geq u_{i+1,j-1}w_{i+1,m} + u_{i,j}w_{i,m} - u_{i,j-1}w_{i,m} - u_{i+1,j}w_{i+1,m} \quad [\text{QI}(\bar{B}; i, i+1, j-1, j)] \\ &\geq u_{i+1,j-1}w_{i+1,m} + (u_{i,j-1} + u_{i+1,j} - u_{i+1,j-1})w_{i,m} - u_{i,j-1}w_{i,m} - u_{i+1,j}w_{i+1,m} \\ & \quad [\text{QI}(u; i, i+1, j-1, j)] \\ &= (u_{i+1,j} - u_{i+1,j-1})(w_{i,m} - w_{i+1,m}) \\ &\geq 0 \quad [\text{monotonicity of } u, w] \end{aligned} \quad (43)$$

For (42),

$$\begin{aligned} & \bar{R}_{i+1,j}^m + \bar{R}_{i,j+1}^m - \bar{R}_{i,j}^m - \bar{R}_{i+1,j+1}^m \\ &= (\bar{B}_{i+1,j-1} + \bar{B}_{i,j} - \bar{B}_{i,j-1} - \bar{B}_{i+1,j}) + (v_{i+1,m} + v_{i,m} - v_{i,m} - v_{i+1,m}) + \\ & \quad (u_{j+1,m} - u_{j,m})(w_{i,m} - w_{i+1,m}) \\ &\geq (u_{j,j} - u_{j,j-1})(w_{i,i} - w_{i+1,i}) + (u_{j+1,m} - u_{j,m})(w_{i,m} - w_{i+1,m}) \\ & \quad [\text{stronger QI}(\bar{B}; i, i+1, j-1, j)] \\ &= (u_{j,j} - u_{j,j-1})(w_{i,i} - w_{i+1,i}) + (u_{j+1,j} - u_{j,j})(w_{i,m} - w_{i+1,m}) \quad [\text{QE}(u; j, j+1, j, m)] \\ &= (u_{j,j} - b)(w_{i,i} - w_{i+1,i} - w_{i,m} + w_{i+1,m}) \quad [u_{j,j-1} = u_{j+1,j} = b] \\ &= u_{j,j} - b \quad [\text{QE}(w; i, i+1, i, m)] \\ &\geq 0 \quad [\text{monotonicity of } u] \end{aligned} \quad (44)$$

which yields the lemma. \square

In the same way as Lemma 5 from Section 3 implied Lemma 2, our new Lemma 10 implies Lemma 8.

Recall too that in section Section 4, we saw how Lemma 5 implied that the LARSCH algorithm could be used to solve the two-sided KY online problem in $O(n)$ time per step, and in $O(n^2)$ time in total. Similarly, Lemma 10 implies that LARSCH algorithm could be used to solve the two-sided Wachs online problem in $O(n)$ time per step, and in $O(n^2)$ time in total.

5.2 The Borchers and Gupta Extension

In [6], motivated by various problems, Borchers and Gupta address the following dynamic programming recurrence: For $0 \leq i \leq j \leq n$ and $0 < r \leq k$,

$$B_{i,j,r} = \begin{cases} c_i & \text{if } i = j \\ \min_{i < t \leq j} \{w(i, t, j) + aB_{i,t-1,f(r)} + bB_{t,j,g(r)}\} & \text{if } i < j \end{cases} \quad (45)$$

In comparing this to (4) one notes many differences. As far as the analysis is concerned the *major* difference is that $w(i, j)$ is replaced by $w(i, t, j)$, which is dependent upon the *splitting-point* t and therefore needs to be moved *inside* the “min”. Our previous definitions of the “quadrangle-inequality” and being “monotone on the lattice of intervals” can not apply to a function of three variables so we need to extend them as follows:

Definition 11 $w(i, t, j)$ satisfies the quadrangle inequality (QI) if for all $i \leq i' < t \leq t' \leq j'$ and $t \leq j \leq j'$,

$$w(i, t, j) + w(i', t', j') \leq w(i', t, j) + w(i, t', j'); \quad (46)$$

and, for all $i < t \leq t' \leq j \leq j'$ and $i \leq i' < t'$,

$$w(i', t', j') + w(i, t, j) \leq w(i', t', j) + w(i, t, j'). \quad (47)$$

Definition 12 $w(i, t, j)$ is monotone (with respect to the lattice of intervals) if for all $[i, j] \subseteq [i', j']$ and $i < t \leq j$, $w(i, t, j) \leq w(i', t, j')$.

Note that if for all t and $i < t \leq j$, $w(i, t, j) = w(i, j)$, then our new definitions of the QI and monotonicity collapse down to Definitions 1 and 4.

The straightforward approach to compute all of the $B_{i,j,r}$ would use $\Theta(kn^3)$ time, where k is the maximum value of the nonnegative integer r , and nonnegative integer functions $f(r) \leq r$, $g(r) \leq r$. Borchers and Gupta [6] showed they can be computed in $O(kn^2)$ time if $w(i, t, j)$ satisfies the generalized QI and monotone property:

Lemma 11 (Lemma 1 in [6])

For $0 \leq i \leq j \leq n$ and $0 < r \leq k$, let $B_{i,j,r}$ be defined by (45). Furthermore let (i) $a, b \geq 1$, (ii) for all $r_1 < r_2$, $f(r_1) \leq f(r_2)$, $g(r_1) \leq g(r_2)$, and (iii) $c_i \geq 0$ and $w(i, j) \geq 0$. If $w(i, t, j)$ satisfies the generalized QI and monotone property, then, for every fixed r , $B_{i,j,r}$ satisfies the QI, i.e., for all $i \leq i' \leq j \leq j'$ and all r ,

$$B_{i,j,r} + B_{i',j',r} \leq B_{i',j,r} + B_{i,j',r}. \quad (48)$$

Lemma 12 (Lemma 2 in [6])

If $B_{i,j,r}$ satisfies the QI for every fixed r , then

$$K_B(i, j, r) \leq K_B(i, j + 1, r) \leq K_B(i + 1, j + 1, r), \quad (49)$$

where $K_B(i, j, r)$ is the maximum splitting point at which $B_{i,j,r}$ attains its minimum value.

As in the KY case, this last lemma provides a $\Theta(n)$ speedup, i.e., from $\Theta(kn^3)$ to $O(kn^2)$.

Setting $k = 1$, $a = b = 1$ and for all t , $w(i, t, j) = w(i, j)$ for some function $w(i, j)$, collapses (45) down to (4) and the BG result collapses down to the standard KY speedup.

Note: Our Lemma 11 is essentially the same as Lemma 1 in [6]. A reader of both papers would note that our statement looks different. The reason for this is that our lemma collects the various conditions required by their Lemma 1 in one place and then lists them in such a way as to easily compare their Lemma with the KY and Wachs results.

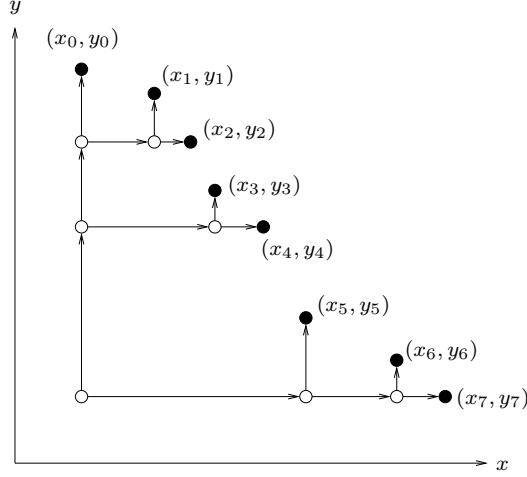


Figure 5: The Rectilinear Steiner Arborescence connecting slide-points $(x_0, y_0), \dots, (x_7, y_7)$. The slide-points are terminals and denoted by solid circles, empty circles denotes the Steiner points. The directed edges can only go up or right.

One interesting immediate application of this result pointed out in the BG paper [6] is finding an optimal Rectilinear Steiner Minimal Arborescence (RSMA) of a *slide*. A slide is a set of points (x_i, y_i) such that, if $i < j$, then $x_i < x_j$ and $y_i > y_j$. See Fig. 5. A Rectilinear Steiner Arborescence is a directed tree in which each edge either goes up or to the right. In [12] it was shown that the minimum cost Rectilinear Steiner Arborescence connecting slide-points $(x_i, y_i), (x_{i+1}, y_{i+1}), \dots, (x_j, y_j)$ satisfies

$$B_{i,j} = \min_{i < t \leq j} \{(x_t - x_i + y_{t-1} - y_j) + B_{i,t-1} + B_{t,j}\}. \quad (50)$$

[12] solved this recurrence in $O(n^3)$ time. As noted in [6] it is actually easy to see that in the RSMA problem, $w(i, t, j) = x_t - x_i + y_{t-1} - y_j$ satisfies the QI (as equality) and is monotone on the lattice of intervals so the BG extension automatically speeds this up to $O(n^2)$ time.

We now generalize our decompositions to the BG recurrence.

Definition 13 For $1 \leq d \leq n$ and $0 < r \leq k$ define the $(n - d + 1) \times (n + 1)$ matrix $D^{d,r}$ by

$$D_{i,j}^{d,r} = \begin{cases} w(i, j, i + d) + aB_{i,j-1,f(r)} + bB_{j,i+d,g(r)} & \text{if } 0 \leq i < j \leq i + d \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (51)$$

Before proving that these matrices are Monge we must first prove the following utility lemma:

Lemma 13 If $w(i, t, j)$ satisfies the QI as defined by (46) and (47), then, for all $i \leq i' < t \leq t' \leq j \leq j'$,

$$w(i, t, j) + w(i', t', j') \leq w(i', t, j') + w(i, t', j). \quad (52)$$

Proof: From (46),

$$\begin{aligned} w(i, t, j) + w(i', t', j) &\leq w(i', t, j) + w(i, t', j), \\ w(i, t, j') + w(i', t', j') &\leq w(i', t, j') + w(i, t', j'). \end{aligned}$$

From (47),

$$w(i', t', j') + w(i', t, j) \leq w(i', t', j) + w(i', t, j'),$$

$$w(i, t', j') + w(i, t, j) \leq w(i, t', j) + w(i, t, j').$$

Summing up these four inequalities, subtracting equal parts from both sides gives

$$2[w(i, t, j) + w(i', t', j')] \leq 2[w(i', t, j') + w(i, t', j)],$$

which yields the lemma. \square

We now continue and show:

Lemma 14 *If the function $B_{i,j,r}$ defined in (45) satisfies the QI for fixed r , and $w(i, t, j)$ satisfies the generalized QI, then, for each $1 \leq d \leq n$ and $0 \leq r \leq k$, $D^{d,r}$ as defined by (51) is a Monge matrix, i.e., for all $0 \leq i < n - d$ and $0 \leq j < n$,*

$$D_{i,j}^{d,r} + D_{i+1,j+1}^{d,r} \leq D_{i+1,j}^{d,r} + D_{i,j+1}^{d,r}. \quad (53)$$

Proof: If $i+1 \not\prec j$ or $j \not\prec i+d$, (53) is trivially true since the right hand side is ∞ . So we assume $i+1 < j < i+d$. Since $B_{i,j,r}$ satisfies the QI,

$$a(B_{i,j-1,f(r)} + B_{i+1,j,f(r)}) \leq a(B_{i+1,j-1,f(r)} + B_{i,j,f(r)}), \quad (54)$$

$$b(B_{j,i+d,g(r)} + B_{j+1,i+1+d,g(r)}) \leq b(B_{j,i+1+d,g(r)} + B_{j+1,i+d,g(r)}). \quad (55)$$

From Lemma 13,

$$w(i, j, i+d) + w(i+1, j+1, i+1+d) \leq w(i+1, j, i+1+d) + w(i, j+1, i+d). \quad (56)$$

Summing up these three inequalities yields (53). \square

Thus, as in Section 2, we can use the SMAWK algorithm to evaluate all of the $B_{i,j}$ in $O(n^2)$ time.

We now again generalize the R^m and L^m matrices:

Definition 14

$$R_{i,j}^{m,r} = \begin{cases} w(i, j, m) + aB_{i,j-1,f(r)} + bB_{j,m,g(r)} & \text{if } 0 \leq i < j \leq m, \\ \infty & \text{otherwise.} \end{cases} \quad (57)$$

$$L_{j,i}^{m,r} = \begin{cases} w(m, i, j) + aB_{m,i-1,f(r)} + bB_{i,j,g(r)} & \text{if } m < i \leq j \leq n, \\ \infty & \text{otherwise.} \end{cases} \quad (58)$$

Lemma 15 *If the function $B_{i,j,r}$ defined in (45) satisfies the QI for fixed r , and $w(i, t, j)$ satisfies the generalized QI, then for each $1 \leq m \leq n$ and $0 \leq r \leq k$, $R^{m,r}$ and $L^{m,r}$ as defined by (57) and (58) are Monge matrices, i.e., for all $0 \leq i < m$ and $0 \leq j \leq m$,*

$$R_{i,j}^{m,r} + R_{i+1,j+1}^{m,r} \leq R_{i+1,j}^{m,r} + R_{i,j+1}^{m,r}, \quad (59)$$

and for all $m < i < n$ and $m < j < n$,

$$L_{i,j}^{m,r} + L_{i+1,j+1}^{m,r} \leq L_{i+1,j}^{m,r} + L_{i,j+1}^{m,r}. \quad (60)$$

Proof: We will only show the proof of $R^{m,r}$, the proof of $L^{m,r}$ is symmetric. If $i+1 \not\prec j$, (59) is trivially true since the right hand side is ∞ . So we assume $i+1 < j$. Since $B_{i,j,r}$ satisfies the QI,

$$a(B_{i,j-1,f(r)} + B_{i+1,j,f(r)}) \leq a(B_{i+1,j-1,f(r)} + B_{i,j,f(r)}). \quad (61)$$

From (46) of Def. 11,

$$w(i, j, m) + w(i + 1, j + 1, m) \leq w(i + 1, j, m) + w(i, j + 1, m). \quad (62)$$

Finally, it is trivially true that

$$b(B_{j,m,g(r)} + B_{j+1,m,g(r)}) = b(B_{j,m,g(r)} + B_{j+1,m,g(r)}). \quad (63)$$

Summing up these three inequalities yields (59). \square

Again just the same way as Lemma 5 from Section 3 implied Lemma 2, our new Lemma 15 implies Lemma 12.

Also, again as before, Lemma 15 implies that the two sided online BG problem could be solved in $O(kn)$ time per step using the LARSCH algorithm. As an example, this implies that the two-sided online minimum cost Rectilinear Steiner Arborescence problem (in which points could be added to the slide one at a time from the left and the right) can be solved in $O(n)$ worst-case time per step.

References

- [1] Alok Aggarwal, Maria M. Klawe, Shlomo Moran, Peter W. Shor, and Robert E. Wilber. Geometric applications of a matrix-searching algorithm. *Algorithmica*, 2:195–208, 1987.
- [2] Alok Aggarwal and James Park. Notes on searching in multidimensional monotone arrays. In *Proceedings of the 29th Annual Symposium on Foundations of Computer Science*, pages 497–512, 1988.
- [3] M. J. Atallah, S. R. Kosaraju, L. L. Larmore, G. L. Miller, and S.-H. Teng. Constructing trees in parallel. In *Proceedings of the First Annual ACM Symposium on Parallel Algorithms and Architectures*, pages 421–431, 1989.
- [4] Amotz Bar-Noy and Richard E. Ladner. Efficient algorithms for optimal stream merging for media-on-demand. *SIAM Journal on Computing*, 33(5):1011–1034, 2004.
- [5] Wolfgang W. Bein, Peter Brucker, Lawrence L. Larmore, and James K. Park. The algebraic Monge property and path problems. *Discrete Applied Mathematics*, 145(3):455–464, 2005.
- [6] Al Borchers and Prosenjit Gupta. Extending the quadrangle inequality to speed-up dynamic programming. *Information Processing Letters*, 49(6):287–290, 1994.
- [7] Rainer E. Burkard, Bettina Klinz, and Rudiger Rudolf. Perspectives of Monge properties in optimization. *Discrete Applied Mathematics*, 70(2):95–161, 1996.
- [8] Zvi Galil and Kunsoo Park. Dynamic programming with convexity, concavity and sparsity. *Theoretical Computer Science*, 92(1):49–76, 1992.
- [9] E. N. Gilbert and E. F. Moore. Variable length encodings. *Bell System Technical Journal*, 38:933–967, 1959.
- [10] Donald E. Knuth. Optimum binary search trees. *Acta Informatica*, 1:14–25, 1971.
- [11] Lawrence L. Larmore and Baruch Schieber. On-line dynamic programming with applications to the prediction of RNA secondary structure. *Journal of Algorithms*, 12(3):490–515, 1991.
- [12] Sailesh K. Rao, P. Sadayappan, Frank K. Hwang, and Peter W. Shor. The rectilinear steiner arborescence problem. *Algorithmica*, 7(2-3):277–288, 1992.
- [13] Amir Said. Efficient alphabet partitioning algorithms for low-complexity entropy coding. In *Proceedings of the 2005 Data Compression Conference*, pages 183–192, 2005.

- [14] Michelle L. Wachs. On an efficient dynamic programming technique of F. F. Yao. *Journal of Algorithms*, 10(4):518–530, 1989.
- [15] Russell L. Wessner. Optimal alphabetic search trees with restricted maximal height. *Information Processing Letters*, 4(4):90–94, 1976.
- [16] Robert Wilber. The concave least-weight subsequence problem revisited. *Journal of Algorithms*, 9(3):418–425, 1988.
- [17] Gerhard J. Woeginger. Monge strikes again: Optimal placement of web proxies in the Internet. *Operations Research Letters*, 27(3):93–96, 2000.
- [18] F. Frances Yao. Efficient dynamic programming using quadrangle inequalities. In *Proceedings of the Twelfth Annual ACM Symposium on Theory of Computing*, pages 429–435, 1980.
- [19] F. Frances Yao. Speed-up in dynamic programming. *SIAM Journal on Matrix Analysis and Applications (formerly SIAM Journal on Algebraic and Discrete Methods)*, 3(4):532–540, 1982.

Last Corrections:

qi03.tex
 mlg: Tue Apr 5 15:59:17 HKT 2005
 qi04.tex
 larry: Tue Apr 5 06:10:11 PDT 2005
 qi04.tex
 mlg:Wed Apr 6 00:49:41 HKT 2005
 wolf: qi06.tex
 Tue Apr 5 15:40:12 PST 2005
 qi07
 larry: Tue Apr 5 20:09:57 PDT 2005
 qi08
 mlg: Wed Apr 6 18:20:04 HKT 2005
 qi09
 larry: Wed Apr 6 05:11:31 PDT 2005
 qi10
 mlg:Wed Apr 6 22:44:38 HKT 2005
 qi11
 larry: Wed Apr 6 20:19:56 PDT 2005
 qi12
 mlg:Thu Apr 7 18:44:57 HKT 2005
 qi13
 Thu Apr 7 22:36:54 HKT 2005
 qi14
 larry,wolf: Thu Apr 7 20:41:33 PDT 2005
 qi15.tex
 wolf: Thu Apr 7 21:17:20 PST 2005
 qi16.tex
 mlg:Fri Apr 8 15:34:41 HKT 2005
 qi16.tex

 qi18.tex
 mlg:Sun Apr 10 23:11:06 HKT 2005
 qi19.tex
 larry: Sun Apr 10 13:06:13 PDT 2005
 qi20.tex
 mlg:Mon Apr 11 13:59:11 HKT 2005
 qi21.tex
 wolf: Mon Apr 11 17:48:29 PST 2005
 qi22.tex
 wolf: Mon Apr 11 22:11:23 PST 2005
 qi23.tex
 mlg:Tue Apr 12 22:23:03 HKT 2005
 qi24.tex, qi25.tex
 mlg:Fri May 6 13:13:20 HKT 2005
 qi26.tex
 yan:Tue May 10 19:07:xx HKT 2005
 qi27.tex
 yan:Fri May 13 10:30:xx HKT 2005
 qi29.tex
 mlg:Mon May 16 16:20:14 HKT 2005