

Parallel Wavelet Tree Construction ^{*†}

Julian Shun[‡]

Carnegie Mellon University
jshun@cs.cmu.edu

We present parallel algorithms for wavelet tree construction with polylogarithmic depth, improving upon the linear depth of the recent parallel algorithms by Fuentes-Sepulveda et al. We experimentally show on a 40-core machine with two-way hyper-threading that we outperform the existing parallel algorithms by 1.3–5.6x and achieve up to 27x speedup over the sequential algorithm on a variety of real-world and artificial inputs. Our algorithms show good scalability with increasing thread count, input size and alphabet size. We also discuss extensions to variants of the standard wavelet tree.

1 Introduction

The *wavelet tree* was first described by Grossi et al. [15], where it was used in compressed suffix arrays. It is a space-efficient data structure that supports access, rank and select queries on a sequence in $O(\log \sigma)$ work, where σ is the alphabet size of the sequence. Since its initial use, wavelet trees have found many other applications, for example in compressed representations of sequences, permutations, grids, graphs, self-indexes based on the Burrows-Wheeler transform [4], images, two-dimensional range queries [19], among many others (see [20, 22] for surveys of applications). While applications of wavelet trees have attracted significant attention, wavelet tree construction has not been widely studied. This is not surprising, as the standard sequential algorithm for wavelet tree construction is very straightforward. The algorithm requires $O(n \log \sigma)$ work for a sequence of length n . However, constructing the wavelet tree of large sequences (with large alphabets) can be time-consuming, and hence parallelizing the construction is important. A step in this direction was taken recently by Fuentes-Sepulveda et al. [12], who describe parallel algorithms for constructing wavelet trees that require $O(n)$ depth (number of parallel time steps).

In this paper, we describe parallel algorithms for wavelet tree construction that exhibit much more parallelism (in particular, polylogarithmic depth). We first describe an algorithm that constructs the tree level-by-level, and requires $O(n \log \sigma)$ work and $O(\log n \log \sigma)$ depth. We then describe a second algorithm that requires $O(W_{\text{sort}}(n) \log \sigma)$ work and $O(D_{\text{sort}}(n) + \log n)$ depth, where $W_{\text{sort}}(n)$ and $D_{\text{sort}}(n)$ are the work and depth, respectively, of the parallel stable integer sorting routine used in the algorithm. Using a linear-work integer sort [23], we obtain a work bound of $O(\frac{n}{\epsilon} \log \sigma)$ and depth bound of $O(\frac{1}{\epsilon}(\sigma^\epsilon + \log n))$ for some constant $0 < \epsilon < 1$, which is sub-linear. For alphabets of polylogarithmic size, this gives an algorithm with $O(\log n)$ depth. Using a super-linear work integer sort [3, 24], we can obtain a work bound of $O(n \log \log n \log \sigma)$ and depth bound of $O(\log n)$ for all alphabets. In addition to having good theoretical bounds, our algorithms are also efficient in practice. We implement our algorithms using Cilk Plus and show experiments on a 40-core shared-memory machine (with two-way hyper-threading) indicating that they outperform the existing parallel algorithms for wavelet tree construction by 1.3–5.6x and achieve up

^{*}This is a longer version of the paper that appears in the *Proceedings of the IEEE Data Compression Conference, 2015*. The levelWT algorithm has been slightly simplified.

[†]We thank Simon Gog and Matthias Petri for discussions, helping with the code in [14], and providing the *trac8* data set, Leo Ferres and Jose Fuentes-Sepulveda for discussions and providing the code from [12], and Guy Blelloch for discussions.

[‡]Supported by the National Science Foundation under grant number CCF-1314590.

to 27x speedup over the sequential algorithm. We show that our implementations scale well with increasing thread count, input size and alphabet size. We then discuss the parallel construction of rank/select structures on binary sequences, which are an essential component to wavelet trees. Finally, we discuss how to adapt our algorithms to variants of wavelet trees—Huffman-shaped wavelet trees [11], multiary wavelet trees [10], and wavelet matrices [7].

2 Preliminaries

We state complexity bounds of algorithms in the work-depth model, where the *work* W is the number of operations required and the *depth* D is the number of time steps required. Then if P processors are available, using Brent’s scheduling theorem [18], we can bound the running time by $O(W/P + D)$. The parallelism of an algorithm is equal to $O(W/D)$. We allow for concurrent reading and writing in the model.

We denote a sequence by S , where $S[i]$ is the i ’th symbol of S and n is its length. We denote an alphabet by $\Sigma = [0, \dots, \sigma - 1]$, where σ is the alphabet size. $\text{access}(S, i)$ returns the symbol at position i of S , $\text{rank}_c(S, i)$ returns the number of times c appears in S from positions 0 to i , and $\text{select}_c(S, i)$ returns the position of the i ’th occurrence of c in S .

A *wavelet tree* is a data structure that supports access, rank and select operations on a sequence in $O(\log \sigma)$ work (we use $\log x$ to mean the base 2 logarithm of x , unless specified otherwise). The standard wavelet tree is a binary tree where each node represents a range of the symbols in Σ using a bitmap (binary sequence). We assume $\sigma \leq n$ as the symbols can be mapped to a contiguous range otherwise. The structure of the wavelet tree is defined recursively as follows: The root represents the symbols $[0, \dots, 2^{\lceil \log \sigma \rceil} - 1]$. A node v which represents the symbols $[a, \dots, b]$ stores a bitmap which has a 0 in position i if the i ’th symbol in the range $[a, \dots, b]$ is in the range $[a, \dots, \frac{(a+b+1)}{2} - 1]$, and 1 otherwise. It will have a left child that represents the symbols $[a, \dots, \frac{(a+b+1)}{2} - 1]$ and a right child that represents the symbols $[\frac{(a+b+1)}{2}, \dots, b]$. The recursion stops when the size of the range is 2 or less or if a node has no symbols to represent. We note that the original wavelet tree description in [15] uses a root whose range is not necessarily a power of 2. However, the definition we use gives the same query complexities and leads to a simpler description of our algorithms.

Along with the bitmaps, each node stores a *succinct* rank/select structure (whose size is sub-linear in the bitmap length) to allow for constant work rank and select queries. The structure of a wavelet tree requires $n \lceil \log \sigma \rceil + o(n \log \sigma)$ bits (the lower order term is for the rank/select structures). The tree topology (parent and child pointers) requires $O(\sigma \log n)$ bits, though this can be reduced or removed by modifying the queries accordingly [6, 19]. The standard sequential algorithm for wavelet tree construction takes $O(n \log \sigma)$ work.

In this paper, we will use the basic parallel primitives, prefix sum and filter [18]. *Prefix sum* takes an array X of length n , an associative binary operator \oplus , and an identity element \perp such that $\perp \oplus x = x$ for any x , and returns the array $(\perp, \perp \oplus X[0], \perp \oplus X[0] \oplus X[1], \dots, \perp \oplus X[0] \oplus X[1] \oplus \dots \oplus X[n-2])$, as well as the overall sum $\perp \oplus X[0] \oplus X[1] \oplus \dots \oplus X[n-1]$. We will use \oplus to be the $+$ operator on integers. *Filter* takes an array X of length n , a predicate function f and returns an array X' of length $n' \leq n$ containing the elements in $x \in X$ such that $f(x)$ returns true, in the same order that they appear in X . Filter can be implemented using prefix sum, and both require $O(n)$ work and $O(\log n)$ depth [18].

3 Related Work

Fuentes-Sepulveda et al. [12] describe a parallel algorithm for constructing a wavelet tree. They observe that for an alphabet where the symbols are contiguous in $[0, \sigma - 1]$, the node at which a symbol s is represented at level i of the wavelet tree can be computed as $s \gg \lceil \log \sigma \rceil - i$, requiring constant work. With this observation they can compute the bitmaps of each level independently. Each level is computed sequentially, requiring $O(n)$ work and depth. Thus, their algorithm requires an overall work of $O(n \log \sigma)$ and $O(n)$ depth. They describe a second algorithm which splits the input sequence into P sub-sequences, where P is the number of processors available. In the first step, the wavelet tree for each sub-sequence is computed sequentially and independently. Then in the second step, the partial wavelet trees are merged. The merging step is non-trivial and requires $O(n)$ depth. Thus the algorithm again requires $O(n \log \sigma)$ work and $O(n)$ depth. This algorithm was shown to perform better than the first algorithm due to the high parallelism in the first step.

Multiple queries on the wavelet tree can be answered in parallel since they do not modify the tree. Furthermore, they can be batched to take advantage of cache locality [12].

Arroyuelo et al. [1] explore the use of wavelet trees in distributed search engines. They do not construct the wavelet tree for the entire text in parallel, but instead sequentially construct the wavelet tree for parts of the text on each machine.

Tischler [26] and Claude et al. [8] discuss how to reduce the space usage of sequential wavelet tree construction. Foschini et al. [11] describe an improved algorithm for sequentially constructing the wavelet tree in compressed format, requiring $O(n + \min(n, nH_h) \log \sigma)$ work, where H_h is the h 'th order entropy of the input. The approach only works if the object produced is the wavelet tree compressed using run-length encoding. Very recently, Babenko et al. [2] and Munro et al. [21] describe sequential wavelet tree construction algorithms that require $O(n \log \sigma / \sqrt{\log n})$ work. The algorithms pack small integers into words, and require extensive bit manipulation. As far as we know, there are no implementations of the algorithms available.

4 Parallel Wavelet Tree Construction

We now describe our algorithms for wavelet tree construction. The construction requires a rank/select data structure for binary sequences. For now we assume that such structures can be created in linear work and logarithmic depth, and defer the discussion to Section 6.

Our first algorithm, *levelWT*, constructs the wavelet tree level-by-level. On each level, the nodes and their bitmaps are constructed in parallel in $O(n)$ work and $O(\log n)$ depth, which gives an overall complexity of $O(n \log \sigma)$ work and $O(\log n \log \sigma)$ depth since there are $O(\log \sigma)$ levels in the tree. The pseudo-code for *levelWT* is shown in Figure 1. We maintain a bitmap B of length n shared by all nodes on each level (Line 4). Nodes will simply store its starting point in this array along with its bitmap length. To keep track of which nodes need to be constructed on each level, we maintain an array A of information for nodes to be added at the next level. Each entry in A stores the starting point (*start*) in the level bitmap (also the starting point in the sequence for the level) for the node, bitmap length (*len*), node identifier (*id*) and length of the alphabet range that it represents (*range*). We represent an entry of A as a 4-tuple (*start*, *len*, *id*, *range*).

Initially, the array A contains just the root node, with a starting index of 0, length of n (it represents all elements), node ID of 0, and a range of $2^{\lceil \log \sigma \rceil}$ (Line 2). The array A' is used as

```

1: procedure LEVELWT( $S, n, \sigma$ )
2:   Nodes = {},  $L = \lceil \log \sigma \rceil$ ,  $S' = S$ ,  $A = \{(0, n, 0, 2^L)\}$ ,  $A' = \{\}$ 
3:   for  $l = 0$  to  $L - 1$  do ▷ Process level-by-level
4:     mask =  $2^{L-(l+1)}$ ,  $B =$  bitmap of length  $n$  ▷ Mask and bitmap for this level
5:     parfor  $j = 0$  to  $|A| - 1$  do
6:       start =  $A[j].start$ , len =  $A[j].len$ , id =  $A[j].id$ ,  $r = A[j].range$ 
7:       Nodes[id].bitmap =  $B + start$ , Nodes[id].len = len
8:       if  $r \leq 2$  then ▷ Node has no children
9:         parfor  $i = 0$  to len - 1 do
10:          if ( $S[start + i] \& mask \neq 0$ ) then  $B[start + i] = 1$  else  $B[start + i] = 0$ 
11:        else
12:           $X = \{\}$  ▷ Array used to store target positions into  $S'$ 
13:          parfor  $i = 0$  to len - 1 do { if ( $S[start + i] \& mask = 0$ )  $X[i] = 1$  else  $X[i] = 0$  }
14:          Perform prefix sum on  $X$  to get offsets of “left” characters ( $X_s$  is the total sum, i.e. number of “left” characters)
15:          parfor  $i = 0$  to len - 1 do
16:            if ( $S[start + i] \& mask \neq 0$ ) then  $B[start + i] = 1$ ,  $S'[start + X_s + i - X[i]] = S[start + i]$ 
17:            else  $B[start + i] = 0$ ,  $S'[start + X[i]] = S[start + i]$ 
18:            if  $X_s > 0$  then  $A'[2 * j] = (start, X_s, 2 * id + 1, r/2)$  ▷ Left child
19:            if  $(len - X_s) > 0$  then  $A'[2 * j + 1] = (start + X_s, len - X_s, 2 * id + 2, r/2)$  ▷ Right child
20:          Filter out empty  $A'$  entries and store into  $A$ 
21:          swap( $S, S'$ )
22:   return Nodes

```

Figure 1: levelWT: Level-by-level parallel algorithm for wavelet tree construction.

the output array for the level. The algorithm proceeds one level at a time for $\lceil \log \sigma \rceil$ levels (Line 3). The bitmaps on each level are determined by the $l + 1$ 'st highest bit in the symbols, so we use a mask to determine the sign of this bit in the symbols (Line 4). The algorithm then loops through all the nodes on the current level in parallel (Lines 5–21). For each node, it sets its bitmap pointer and length in the Nodes array (Line 7). If the alphabet range of the node is 2 or less, then it has no children (Line 8). The algorithm then just loops over the node's symbols in S in parallel, and sets each bit in the bitmap according to the sign of the symbol's $l + 1$ 'st highest bit (Lines 9–10).¹ Otherwise, the symbols in S are rearranged (and stored into S') so that they are in the correct order on the next level. To do this in parallel, we first count the number of symbols that go to the left child ($l + 1$ 'st highest bit is 0), and the offset of each such symbol using a prefix sum (Lines 12–14). With the prefix sum array X and result X_s , the symbols that go to the right child can be computed as well using the formula $X_s + i - X[i]$ (the number of symbols with an $l + 1$ 'st highest bit of 0 is X_s , so this is the number of symbols to offset, and then the number of symbols with an $l + 1$ 'st highest bit of 1 up to index i is $i - X[i]$). In Lines 15–17, the bits in the bitmap and positions in S' are set in parallel. Children nodes are placed into A' on Lines 18–19 if the number of symbols represented for the child is greater than 0. The children's node IDs are computed as twice the current node ID plus one for the left child and twice the current node ID plus two for the right child in order to give all nodes unique IDs. The starting point in the bitmap of the next level is the same as the current starting point for the left child, and is the current starting point plus the number of elements on the left (X_s) for the right child. The length is stored from the computation before. The range of each child is half of the current range. After each level, the non-empty entries of A' are filtered out and stored into A (Line 20). The roles of S and S' are swapped for the next level (Line 21).

We note that setting the bits in B (Lines 10, 16 and 17) must be done atomically since multiple processors may write to the same word concurrently. This can be implemented using a loop with a compare-and-swap until successful. An optimization we use is to only perform atomic writes if a word is shared between two nodes (there can be at most two shared words per node, as the bits for each node are contiguous in B), and for the remaining

¹The actual code requires bit arithmetic to access the appropriate bit in the word, which we omit for simplicity.

```

1: procedure sortWT(S, n,  $\sigma$ )
2:   Nodes = {}, L =  $\lceil \log \sigma \rceil$ 
3:   parfor l = 0 to L - 1 do
4:     mask =  $2^{L-(l+1)}$ , B = bitmap of length n ▷ Mask and bitmap for this level
5:     if l = 0 then ▷ No sorting required for first level
6:       parfor i = 0 to n - 1 do { if (S[i] & mask  $\neq$  0) B[i] = 1 else B[i] = 0 }
7:       Nodes[0].bitmap = 0, Nodes[0].len = n
8:     else
9:       S' = S stably sorted by top l bits
10:      parfor i = 0 to n - 1 do { if (S'[i] & mask  $\neq$  0) B[i] = 1 else B[i] = 0 }
11:      O = indices i such that (S'[i]  $\gg$  L - l)  $\neq$  (S'[i - 1]  $\gg$  L - l) ▷ Using a filter
12:      parfor i = 0 to |O| - 1 do
13:        id =  $2^l - 1 + (S'[O[i]] \gg L - l)$ , Nodes[id].bitmap = O[i], Nodes[id].len = O[i + 1] - O[i]
14:   return Nodes

```

Figure 2: sortWT: Sorting-based parallel algorithm for wavelet tree construction.

words we parallelize at the granularity of a word. Inside each word, the updates are done sequentially. This allows us to use regular writes for all but at most two words per node.

We also note that we are able to stop early when the range of a node is 2 or less since the construction of the previous level provides this information. The algorithm of [12] (and also our next algorithm) processes all levels independently, so it is not easy to stop early.

We now analyze the complexity of levelWT. For each level, there is a total of $O(n)$ work performed, since each symbol is processed a constant number of times. The prefix sum on Lines 14 and the filter on Line 20 require linear work and $O(\log n)$ depth per level. The parallel for-loops on Lines 9–10, 13 and 15–17 require linear work and $O(1)$ depth per level. There are $O(\log \sigma)$ levels so the total work is $O(n \log \sigma)$ and depth is $O(\log n \log \sigma)$. Since $\sigma \leq n$, the depth is polylogarithmic in n . We obtain the following theorem:

Theorem 4.1. *levelWT requires $O(n \log \sigma)$ work and $O(\log n \log \sigma)$ depth.*

We now describe our second wavelet tree construction algorithm, which constructs all levels of the wavelet tree in parallel. We refer to this algorithm as *sortWT*. Since preceding levels cannot provide information to later levels, we must do independent computation per level to obtain the correct ordering of the sequence for the level. We will make use of the observation of Fuentes-Sepulveda et al. [12] that the node at which a symbol s is represented at level l of the wavelet tree ($l = 0$ for the root) is encoded in the top l bits. For level l , our algorithm sorts S using the top l bits as the key, which gives the correct ordering of the sequence for the level. We note that the sort must be stable since the relative ordering of nodes with the same top l bits must be preserved in the wavelet tree.

The pseudo-code for sortWT is shown in Figure 2. As in levelWT, we use a mask and a bitmap B for each level (Line 4). For the first level ($l = 0$), no sorting of S is required, so we simply fill the bitmap according to the highest bit of each symbol (Lines 5–6). For each subsequent level, we first stably sort S by the top l bits to obtain the symbols in the correct order, and store it in S' (Line 9). The bitmap is filled according to the $l + 1$ 'st highest bit of each symbol (Line 10). To compute the length of each node's bitmap we use a filter to find all the indices where the symbol's top l bits differ from the previous symbol's top l bits in S' (Line 11). These mark the bitmaps of each node since S' is sorted by the top l bits. The length of each bitmap can be computed by the difference in indices. On lines 12–13, we set the bitmap pointers and lengths for the nodes on the current level. The IDs of the nodes start at $2^l - 1$, since there are up to $2^l - 1$ nodes in previous levels, and each node ID is offset by the top l bits of the symbols that it represents, as this determines the node's position in the level. As in levelWT, updates to the bitmaps are done in parallel at word granularity.

We now discuss the algorithm's complexity. Let $W_{\text{sort}}(n)$ and $D_{\text{sort}}(n)$ be the work and

depth, respectively, of the stable sort on Line 9. The filter on Line 11 requires $O(n)$ work and $O(\log n)$ depth. The parallel for-loops on Lines 6, 10 and 12–13 require $O(n)$ work and $O(1)$ depth. The overall work is $O(W_{\text{sort}}(n) \log \sigma)$ and since all levels can be computed in parallel, the overall depth is $O(D_{\text{sort}}(n) + \log n)$. This gives the following theorem:

Theorem 4.2. *sortWT requires $O(W_{\text{sort}}(n) \log \sigma)$ work and $O(D_{\text{sort}}(n) + \log n)$ depth.*

Using $\frac{1}{\epsilon}$ rounds of linear-work parallel stable integer sorting [23] ($W_{\text{sort}}(n) = O(\frac{n}{\epsilon})$ and $D_{\text{sort}}(n) = O(\frac{1}{\epsilon}(\sigma^\epsilon + \log n))$) for some constant $0 < \epsilon < 1$, we obtain a work bound of $O(\frac{n}{\epsilon} \log \sigma)$ and depth bound of $O(\frac{1}{\epsilon}(\sigma^\epsilon + \log n))$, which is sub-linear. For $\sigma = O(\log^c n)$ for any constant c , this gives $O(\log n)$ depth (by setting ϵ appropriately). We can alternatively use a stable integer sorting algorithm with $W_{\text{sort}}(n) = O(n \log \log n)$ and $D_{\text{sort}} = O(\frac{\log n}{\log \log n})$ (either using randomization [24] or using super-linear space [3]) to obtain an overall work of $O(n \log n \log \sigma)$ and depth of $O(\log n)$ for any alphabet size.

Space usage. The input and output of the algorithms is $O(n \log \sigma)$ bits. levelWT requires two auxiliary arrays for the prefix sum on each level (that can be reused per level), which takes $O(n \log n)$ bits. For sortWT, since all levels are processed in parallel, and each level requires $O(n \log n)$ bits for the integer sort, the total space usage is $O(n \log n \log \sigma)$ bits.

Due to the high space usage and hence memory footprint of sortWT, we found that processing the levels one-by-one gives better performance in practice, as will be discussed in Section 5 (although this increases the depth by a factor of $O(\log \sigma)$). We call this modified version *msortWT*. On each level, msortWT sorts the sequence from the previous level. Since the levels are processed one-by-one, msortWT has a space usage of $O(n \log n)$ bits.

5 Experiments

Implementations. We compare our implementations of levelWT, sortWT and msortWT with existing parallel implementations as well as a sequential implementation. The implementations all use the levelwise representation of the bitmaps, where one bitmap of length n is stored per level, and nodes have pointers into the bitmaps. sortWT and msortWT use linear-work parallel stable integer sorting. We use the parallel prefix sum, filter and integer sorting routines from the Problem Based Benchmark Suite [25]. We compare with the implementations of Fuentes-Sepulveda et al. [12], one which computes each level of the wavelet tree independently in parallel (*FEFS*), and one which computes a partial wavelet tree for each thread, and then merges them together (*FEFS2*). Both implementations require the alphabet size to be a power of 2, so we only report times for the inputs with such an alphabet size. We implement a sequential version of wavelet tree construction (*serialWT*), and found its performance to be competitive with the times of the sequential algorithm reported in [16]. We also tried the serial implementation in SDSL [14] for constructing a balanced wavelet tree, but found it to be slower than serialWT on our inputs. However, the SDSL implementation is more space-efficient, and sometimes faster on the Burrows-Wheeler transformed inputs. A comparison with SDSL is presented in Section A.1 of the Appendix.

Experimental Setup. We run our experiments on a 40-core (with two-way hyper-threading) machine with 4×2.4 GHz Intel 10-core E7-8870 Xeon processors (with a 1066MHz bus and 30MB L3 cache), and 256GB of main memory. The parallel codes use Cilk Plus, and we compile our code with the g++ compiler version 4.8.0 (which supports Cilk Plus) with the `-O2` flag. The times reported are based on a median of 3 trials.

We use a variety of real-world and artificial sequences. The real-world sequences

Text	n	σ	serialWT (T_1)	levelWT (T_1)	levelWT (T_{40})	sortWT (T_1)	sortWT (T_{40})	msortWT (T_1)	msortWT (T_{40})	FEFS (T_1)	FEFS (T_{40})	FEFS2 (T_1)	FEFS2 (T_{40})
chr22	$3.35 \cdot 10^7$	4	0.486	0.611	0.018	0.786	0.046	0.768	0.029	1.03	0.53	0.98	0.1
etext99	$1.05 \cdot 10^8$	146	4.12	6.99	0.28	10.5	0.393	9.79	0.364	–	–	–	–
HG18	$2.83 \cdot 10^9$	4	35.5	45.5	1.32	57.7	1.88	56.9	1.67	76.6	39.1	72.2	2.55
howto	$3.94 \cdot 10^7$	197	1.65	2.69	0.105	4.07	0.161	3.86	0.144	–	–	–	–
jdk13c	$6.97 \cdot 10^7$	113	2.51	3.9	0.159	6.13	0.234	5.63	0.22	–	–	–	–
proteins	$1.18 \cdot 10^9$	27	31.3	52.2	1.82	75.3	2.59	71.3	2.28	–	–	–	–
rectail96	$1.15 \cdot 10^8$	93	3.54	5.88	0.231	10.2	0.373	9.39	0.34	–	–	–	–
rfc	$1.16 \cdot 10^8$	120	3.8	6.51	0.261	10.1	0.37	9.28	0.348	–	–	–	–
sprot34	$1.1 \cdot 10^8$	66	3.69	6.26	0.248	9.48	0.36	8.8	0.328	–	–	–	–
trec8	$2.43 \cdot 10^8$	528155	33.3	50.4	2.08	138	5.5	104	4.55	–	–	–	–
w3c2	$1.04 \cdot 10^8$	256	3.82	6.66	0.275	10.6	0.388	9.78	0.357	11.1	2.0	10.6	0.51
wikisamp	10^8	204	3.52	6.16	0.264	9.78	0.374	9.08	0.349	–	–	–	–
rand-2 ⁸	10^8	2 ⁸	5.76	8.58	0.36	14.3	0.652	12.1	0.533	12.4	1.71	12.3	0.5
rand-2 ¹⁰	10^8	2 ¹⁰	6.88	11	0.456	19	0.857	15.9	0.708	15.0	1.71	15.3	0.58
rand-2 ¹²	10^8	2 ¹²	8.32	12.4	0.525	24.5	1.11	20.4	0.922	18.7	1.78	17.4	0.67
rand-2 ¹⁶	10^8	2 ¹⁶	11.2	16.4	0.655	36	1.59	29.5	1.34	34.4	3.26	32.4	1.28
rand-2 ²⁰	10^8	2 ²⁰	14	20.4	0.772	49.5	2.19	40.6	1.85	65.7	7.14	64.8	3.94

Table 1: Comparison of running times (seconds) of wavelet tree construction algorithms on a 40-core machine with hyper-threading. T_{40} is the time using 40 cores (80 hyper-threads) and T_1 is the time using a single thread.

include strings from <http://people.unipmn.it/manzini/lightweight/corpus/>, XML code from Wikipedia (*wikisamp*), protein data from [http://pizzachili.dcc.uchile.cl/texts/protein/\(proteins\)](http://pizzachili.dcc.uchile.cl/texts/protein/(proteins)), the human genome from <http://webhome.cs.uvic.ca/~thomo/HG18.fasta.tar.gz> (*HG18*), and a document array of text collections (*trec8*). The artificial inputs (*rand*), parameterized by σ , are generated by drawing each symbol uniformly at random from the range $[0, \dots, \sigma - 1]$. The lengths and alphabet sizes of the inputs are listed in Table 1.

Due to the various choices for rank/select structures, each of which has different space/time trade-offs, we do not include their construction times in the wavelet tree construction time. We modified the FEFS and FEFS2 code accordingly. The times for our code include generating the parent/child pointers for the nodes, although these could be removed using techniques from [6, 19]. FEFS and FEFS2 do not generate these pointers.

Results. Table 1 shows the single-thread (T_1) and 40-core with two-way hyper-threading (T_{40}) running times on the inputs for the various implementations. Our results show that levelWT is faster than sortWT and msortWT both sequentially and in parallel. This is because sortWT and msortWT use sorting, which has a larger overhead. msortWT is slightly faster than sortWT due to its smaller memory footprint. Compared to serialWT, levelWT is 1.2–1.8x slower on a single thread, and 13–27x faster on 40 cores with hyper-threading. The self-relative speedup of levelWT ranges from 23 to 35. On 40 cores, sortWT and msortWT are 6–19x and 7–22x faster than serialWT, respectively. sortWT and msortWT achieve self-relative speedups of 17–31 and 22–34, respectively.

FEFS2 always outperforms FEFS on 40 cores with two-way hyper-threading because FEFS splits the work among only $\log \sigma$ threads, which is less than 80 on all our inputs ($\sigma < 2^{80}$), whereas FEFS2 splits the work among all available threads in its first step. The second (merging) step of FEFS2, however, only makes use of $\log \sigma$ threads, but this is a smaller fraction of the total time. On 40 cores with hyper-threading, our best implementation (levelWT) outperforms FEFS2 by a factor of 1.3–5.6x, and FEFS by much more. Compared to msortWT, FEFS2 is faster in some cases and slower in others.

In Figure 3, we plot the speedup of the parallel implementations relative to serialWT as a function of the number of threads for HG18 and rand-2¹⁶. For HG18, our implementations and FEFS2 scale well up to 80 hyper-threads. FEFS only scales up to 2 threads due to the

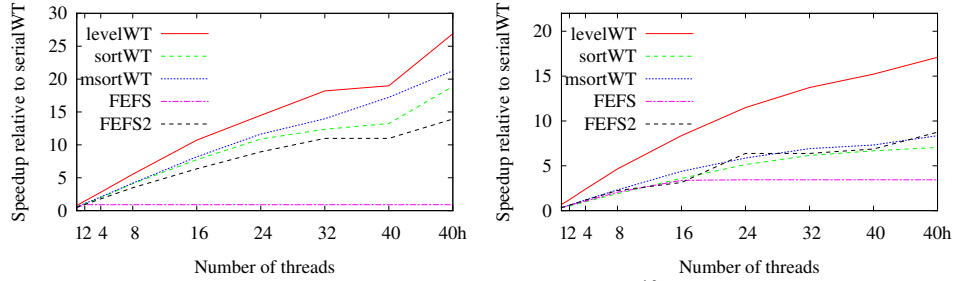


Figure 3: Speedup of implementations relative to serialWT for HG18 (left) and rand-2¹⁶ (right). (40h) is 40 cores with hyper-threading.

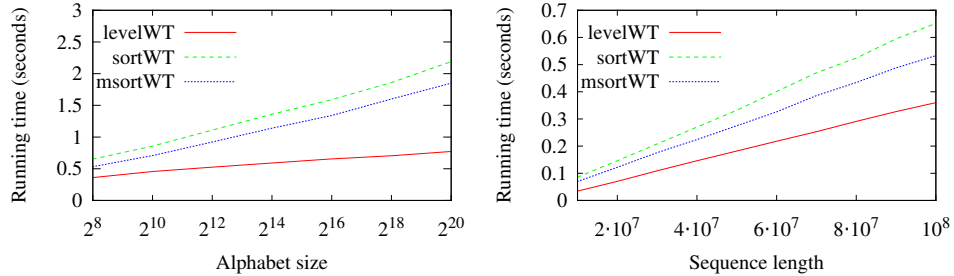


Figure 4: 40-core (with hyper-threading) running times vs. σ (left, x-axis in log-scale) and vs. n (right) on random sequences ($\sigma = 2^8$). small alphabet size. For rand-2¹⁶, our implementations again exhibit good scalability. FEFS scales up to 16 threads as there are 16 levels in the tree, while FEFS2 scales to more threads. FEFS2 is competitive with msortWT, but slower than levelWT.

In Figure 4 (left), we plot the 40-core parallel running time of our three implementations as a function of the alphabet size for random sequences of length 10^8 . We see that for fixed n , the running times increase nearly linearly with $O(\log \sigma)$, which is expected since the total work is $O(n \log \sigma)$. In Figure 4 (right), we plot the running time of the implementations as a function of n for $\sigma = 2^8$ on random sequences, and see that the times increase linearly with n as expected. Our algorithms exhibit similar speedups on 40 cores as we vary σ or n , since the core count is much lower than the available parallelism of the algorithms.

In summary, our experiments show that our three parallel algorithms for wavelet tree construction scale well with the number of threads, input length and alphabet size. levelWT outperforms sortWT and msortWT as it does not have the overheads of sorting, and achieves good speedup over serialWT. Overall, our fastest algorithm outperforms FEFS and FEFS2.

6 Parallel Construction of Rank/Select Structures

Wavelet trees make use of succinct rank/select structures which support constant work rank and select queries on binary sequences. In this section, we describe how to construct these structures in parallel using $O(n)$ work and $O(\log n)$ depth for a binary sequence of length n . We note that *sequential* construction of rank/select structures in $o(n)$ work have been described in [2, 21], however parallel construction in linear work suffices for our purposes.

The rank structure of Jacobson [17] stores the rank of every $\log^2 n$ 'th bit in a first-level directory, and the rank of every $\log n$ 'th bit in each of the ranges in a second-level directory. Rank queries in each range of size $\log n$ can be answered by at most two table look-ups, where the table stores the rank of all bit-strings of length up to $\frac{\log n}{2}$. The first- and second-level directories can be constructed by converting the bit-string to a length n array of 0's and 1's and computing a prefix sum on the array in $O(n)$ work and $O(\log n)$ depth. Entries in the second-level directory require $\log \log n$ bits each, and for space efficiency, they need to be packed into words. This can be done by processing groups of $O(\frac{\log n}{\log \log n})$ entries (the

number that fits in a word) in parallel, and packing each group into a word sequentially. There are $O(\frac{n}{\log n})$ entries, and word operations take $O(1)$ work and depth per entry, so this process takes $O(n)$ work and $O(\log n)$ depth. The look-up table can be constructed in $o(n)$ work and $O(\log n)$ depth, as we can compute the number of 1's in bit-strings of size $O(\log n)$ in $O(\log n)$ work and depth, and there are $O(2^{\frac{\log n}{2}} \log n) = O(\sqrt{n} \log n)$ such bit-strings.

Clark's select structure [5] stores the position of every $\log n \log \log n$ 'th 1 bit in a first-level directory. Then for each range r between the positions, if $r \geq \log^2 n (\log \log n)^2$ the $\log n \log \log n$ answers in the range are stored directly. Otherwise the position of every $\log r \log \log n$ 'th 1 bit is stored in a second-level directory. The sub-ranges r' in the second-level directory are again considered, and if $r' \geq \log r' \log r (\log \log n)^2$, then all answers in the range are stored directly. Otherwise, a look-up table is constructed for all bit-strings of length less than r' . To parallelize the construction, we first convert the bit-string to an array of 0's and 1's, and compute the positions of all the 1 bits using a prefix sum and filter in $O(n)$ work and $O(\log n)$ depth. This allows all of the ranges to be processed in parallel. Constructing each second-level directory again uses prefix sum and filter. Over all directories, this sums to $O(n)$ work and $O(\log n)$ depth. The look-up table can be constructed in $o(n)$ work and $O(1)$ depth, similar to the rank structure. Packing entries into words can also be done within the complexity bounds.

We note that there have been more practical variants of rank/select structures (see, e.g., [13, 27] and references therein) that have a similar high-level structure.

7 Extensions

The *Huffman-shaped wavelet tree*, where each node is placed at a level proportional to the length of its Huffman code, was introduced to improve compression and average query performance [11]. To construct it in parallel, we first compute the Huffman tree and prefix codes in $O(\sigma + n)$ work and $O(\sigma + \log n)$ depth using the algorithm of [9]. The construction then follows the strategy of levelWT. To decide how to set the bitmaps in the internal nodes and rearrange \mathbf{S} , we must know the side of the tree that each symbol is located on. We map each symbol to an integer corresponding to the location of its leaf in an in-order traversal of the Huffman tree, which can be done in parallel using an Euler tour algorithm in $O(\sigma)$ work and $O(\log \sigma)$ depth [18]. At each internal node, the symbols to the left and to the right are in consecutive ranges, and we store the highest mapped integer of nodes in its left sub-tree, which can be done during the Euler tour computation. Then the decision of how to set the bitmap and where to place the symbol in \mathbf{S}' can be made with a single comparison with the mapped integers. The rest of the computation follows the logic of levelWT. For a tree of height h , the overall work (including Huffman encoding) is $O(nh)$ as linear work is done per level. The overall depth is $O(\sigma + h \log n)$, as each level of the tree takes $O(\log n)$ depth.

Ferragina et al. [10] describe the *multiary wavelet tree* where each node has up to d children for some value d , and stores sequences of symbols in the range $[0, \dots, d - 1]$. The height of the tree is $O(\log_d \sigma)$. We describe the parallel construction for the case where $d = O(\log^\epsilon n)$ for $\epsilon < 1/3$, and d is a power of two.² We modify sortWT to process the levels one-by-one and save the sorted sequence \mathbf{S}' for the next level. On level l , \mathbf{S}' is already sorted by the top $(l - 1) \log d$ bits, so we only need to sort the next $\log d$ highest bits within each of the sub-sequences sharing the same top $(l - 1) \log d$ bits. The sub-sequence boundaries can

²The requirement $\epsilon < 1/3$ is necessary for the analysis of the rank/select structure [2] that we parallelize.

be identified with a filter, and for each sub-sequence we apply a stable integer sort using the $\log d$ appropriate bits for the level as the key. We substitute the bitmap B with a sequence of $\log d$ -bit entries, and in parallel each entry is set according to the value of the appropriate $\log d$ bits of each symbol. There will be up to d^l nodes on level l , and the offset to the Nodes array is $d^l - 1$. Since $d = O(\log^\epsilon n)$, the integer sort on each level requires $O(n)$ work and $O(\log n)$ depth, giving an overall work of $O(n \log_d \sigma)$ and depth of $O(\log n \log_d \sigma)$. We also parallelize the construction of rank/select structures on sequences with larger alphabets, used in the multiary wavelet tree. The description is in Section A.2 of the Appendix.

The *wavelet matrix* [7] is a variant of the wavelet tree where on level l , all symbols with a 0 as their l 'th highest bit are represented on the left side of the level's bitmap and all symbols with a 1 as their l 'th highest bit are represented on the right. Each level stores the number of 0's on the level. The bitmap is filled based on the $l + 1$ 'st highest bit of the symbols. To construct the wavelet matrix, we proceed level-by-level and stably reorder S based on the l 'th highest bit of the symbols using standard operations involving prefix sum (similar to levelWT) in $O(n)$ work and $O(\log n)$ depth, which also gives the number of 0's on the level. The bitmap for the level is then filled in parallel in $O(n)$ work and $O(1)$ depth. This gives an algorithm with $O(n \log \sigma)$ work and $O(\log n \log \sigma)$ depth. We can alternatively use a strategy similar to sortWT, but using the reverse of the top l bits as the key when sorting.

References

- [1] D. Arroyuelo et al. "Distributed search based on self-indexed compressed text". *Information Processing & Management* (2012).
- [2] M. A. Babenko et al. "Wavelet Trees Meet Suffix Trees". In *SODA*. 2015.
- [3] P. C. P. Bhatt et al. "Improved deterministic parallel integer sorting". *Information and Computation* (1991).
- [4] M. Burrows and D. J. Wheeler. *A Block-Sorting Lossless Data Compression Algorithm*. Tech. rep. HP Labs, 1994.
- [5] D. R. Clark. "Compact Pat Trees". PhD thesis. 1996.
- [6] F. Claude and G. Navarro. "Practical Rank/Select Queries over Arbitrary Sequences". In *SPIRE*. 2008.
- [7] F. Claude and G. Navarro. "The Wavelet Matrix". In *SPIRE*. 2012.
- [8] F. Claude et al. "Space Efficient Wavelet Tree Construction". In *SPIRE*. 2011.
- [9] J. A. Edwards and U. Vishkin. "Parallel Algorithms for Burrows-Wheeler Compression and Decompression". *Theor. Comput. Sci.* (2014).
- [10] P. Ferragina et al. "Compressed Representations of Sequences and Full-text Indexes". *ACM Trans. Algorithms* (2007).
- [11] L. Foschini et al. "When Indexing Equals Compression: Experiments with Compressing Suffix Arrays and Applications". *ACM Trans. Algorithms* (2006).
- [12] J. Fuentes-Sepulveda et al. "Efficient Wavelet Tree Construction and Querying for Multicore Architectures". In *SEA*. 2014.
- [13] S. Gog and M. Petri. "Optimized succinct data structures for massive data". *SPE* (2013).
- [14] S. Gog et al. "From Theory to Practice: Plug and Play with Succinct Data Structures". In *SEA*. 2014.
- [15] R. Grossi et al. "High-order Entropy-compressed Text Indexes". In *SODA*. 2003.
- [16] R. Grossi et al. "Wavelet Trees: From Theory to Practice". In *CCP*. 2011.
- [17] G. J. Jacobson. "Succinct Static Data Structures". PhD thesis. 1988.
- [18] J. Jaja. *Introduction to Parallel Algorithms*. Addison-Wesley Professional, 1992.
- [19] V. Makinen and G. Navarro. "Rank and select revisited and extended". *Theor. Comput. Sci.* (2007).
- [20] C. Makris. "Wavelet trees: A survey". *Comput. Sci. Inf. Syst.* (2012).
- [21] J. I. Munro et al. "Fast Construction of Wavelet Trees". In *SPIRE*. 2014.
- [22] G. Navarro. "Wavelet Trees for All". In *CPM*. 2012.
- [23] S. Rajasekaran and J. H. Reif. "Optimal and sublogarithmic time randomized parallel sorting algorithms". *SIAM J. Comput.* (1989).
- [24] R. Raman. "The Power of Collision: Randomized Parallel Algorithms for Chaining and Integer Sorting". In *Foundations of Software Technology and Theoretical Computer Science*. 1990.
- [25] J. Shun et al. "Brief announcement: the Problem Based Benchmark Suite". In *SPAA*. 2012.
- [26] G. Tischler. "On Wavelet Tree Construction". In *CPM*. 2011.

Text	n	σ	serialWT (T_1)	SDSL (T_1)	levelWT (T_1)	levelWT (T_{40})	sortWT (T_1)	sortWT (T_{40})	msortWT (T_1)	msortWT (T_{40})	FEFS (T_1)	FEFS (T_{40})	FEFS2 (T_1)	FEFS2 (T_{40})
chr22	$3.35 \cdot 10^7$	4	0.48	0.86	0.59	0.017	0.762	0.044	0.741	0.029	1.0	0.51	0.95	0.09
etext99	$1.05 \cdot 10^8$	146	2.83	2.9	5.79	0.266	9.47	0.377	9.04	0.343	–	–	–	–
HG18	$2.83 \cdot 10^9$	4	31	63.2	36.7	1.17	43.3	1.6	41.4	1.42	71.9	37.9	67.2	2.53
howto	$3.94 \cdot 10^7$	197	1.2	1.2	2.14	0.1	3.58	0.157	3.42	0.14	–	–	–	–
jdk13c	$6.97 \cdot 10^7$	113	1.51	1.13	2.88	0.152	5.23	0.223	4.9	0.207	–	–	–	–
proteins	$1.18 \cdot 10^9$	27	26.1	30.2	43.9	1.69	65.1	2.34	61	2.11	–	–	–	–
rctail96	$1.15 \cdot 10^8$	93	2.11	2.31	4.61	0.221	8.83	0.35	8.28	0.327	–	–	–	–
rfc	$1.16 \cdot 10^8$	120	2.53	2.73	5.26	0.251	8.95	0.355	8.52	0.335	–	–	–	–
sprot34	$1.1 \cdot 10^8$	66	2.55	2.6	5.19	0.24	8.53	0.343	8.14	0.31	–	–	–	–
w3c2	$1.04 \cdot 10^8$	256	2.3	1.85	5.13	0.263	9.76	0.375	8.6	0.349	9.44	1.94	8.53	0.47
wikisamp	10^8	204	2.22	1.91	5.06	0.253	8.67	0.353	8.17	0.325	–	–	–	–

Table 2: Running times (seconds) of wavelet tree construction algorithms on a 40-core machine with hyper-threading on the Burrows-Wheeler transform of the text. T_{40} is the time using 40 cores (80 hyper-threads) and T_1 is the time using a single thread.

[27] D. Zhou et al. “Space-Efficient, High-Performance Rank and Select Structures on Uncompressed Bit Sequences”. In *SEA*. 2013.

A Appendix

A.1 Performance on Burrows-Wheeler transformed inputs

For certain applications (see [20, 22]), the wavelet tree is constructed on the Burrows-Wheeler transform of the sequence. In Table 2, we report the running times of the algorithms on the Burrows-Wheeler transform of the real-world texts. For sequential times, we also report the wavelet tree implementation from SDSL [14], which performs well on sequences with many sequences of repeated characters. It uses an optimization that groups the writes into the bitmap for repeated characters into a single write. This gives an advantage for the Burrows-Wheeler transformed texts, in which there are often many sequences of repeated characters. In contrast to serialWT, the SDSL implementation generates the wavelet tree structure first before updating the bitmaps. In the column labeled *SDSL* in Table 2, we report the SDSL times using the balanced tree option and without the construction of the rank/select structures. For our implementations, we use an optimization where consecutive 1 bits in a word are written together instead of individually.³ This is similar to SDSL but it finds consecutive 1 bits on each level instead of consecutive repeated characters in the original sequence. The optimization slightly improves the running time for the Burrows-Wheeler transformed texts, but makes negligible difference on the original texts. For some inputs, SDSL performs slightly faster than serialWT, though it is slower than or performs about the same as serialWT on other inputs.

Overall, the times are faster on the Burrows-Wheeler transformed texts than on the original texts. This is because the bits of the same value are grouped into larger chunks in the bitmaps, and hence there are fewer words that actually need to be updated (words with all 0 bits do not need to be updated after initialization). Furthermore, the optimization of writing consecutive 1 bits together has more of a benefit in this setting. The relative performance among our parallel algorithms remains the same, with levelWT being the fastest, followed by msortWT and then sortWT. Compared to the faster of serialWT and SDSL for each input, levelWT is 1.7–2.8 times slower on a single thread and 7–28 times faster on 40 cores with hyper-threading. The self-relative speedups of levelWT, sortWT and msortWT are

³We thank Matthias Petri for suggesting this optimization.

18–35, 17–28 and 24–29, respectively. Our implementations are always faster than FEFS and FEFS2 in parallel, though FEFS2 gets good speedup.

A.2 Parallel Construction of Generalized Rank/Select Structures

A generalized rank/select structure supports rank and select queries on sequences with larger alphabets in $O(1)$ work. The rank query $rank_c(\mathbf{S}, i)$ returns the number of symbols less than or equal to c in \mathbf{S} from positions 0 to i (in contrast to the definition in Section 2).

We construct in parallel the rank/select structure described in the Appendix of [2], which works for an alphabet size of $\sigma = O(\log^\epsilon n)$ for $\epsilon < 1/3$. We only discuss its construction, and refer the reader to [2] for discussions on the space usage of the resulting structure. Again, the sequential construction process described in [2] takes sub-linear work, but we describe a parallel algorithm that takes linear work, which suffices for our purposes.

Rank. The rank structure stores the σ ranks (one per character) of every $\sigma \log^2 n$ 'th symbol in a first-level directory, and the σ ranks of every $\log n / (3 \log \sigma)$ 'th symbol in a second-level directory. Queries on blocks of up to length $\log n / (3 \log \sigma)$ can be answered in constant work using table lookup. The lookup table has at most $O(\sigma^{2 + \log n / (3 \log \sigma)}) = o(n)$ entries, so can be created in $o(n)$ work and $O(\log n)$ depth.

Computing ranks for the second-level directory entries again uses table lookup. Sequentially, to compute the ranks of an entry we take the ranks of the previous entry and update it with the block of $\log n / (3 \log \sigma)$ symbols between them. The results of updating the ranks of all possible second-level directory entry along with all possible blocks of $\log n / (3 \log \sigma)$ symbols can be pre-computed in a table of size $o(n)$ again in $o(n)$ work and $O(\log n)$ depth.

We now need to parallelize the computation of the ranks of the second-level directory entries. We will do this by applying a prefix sum on the entries with the combining operator defined by a lookup table. Since a prefix sum computation uses a reduction tree [18] and will combine entries separated by more than $\log n / (3 \log \sigma)$ symbols, we cannot directly apply the table lookup scheme described above. We fix this as follows. During the prefix sum, a combined entry e will represent more than one original entry, and we store the ranks of the first and last entry that it represents (call these e_F and e_L). Without loss of generality, assume e combines with an entry e' before it in the input. We then use the lookup table described previously to update e_F given e'_L and the block of $\log n / (3 \log \sigma)$ symbols between them. All ranks of e_L will then be increased by the same amount as in e_F . When combining e and its previous entry e' , we keep e'_F and e_L . Each combine operation takes constant work, so the prefix sum takes $O(n)$ work and $O(\log n)$ depth.

The ranks of a first-level directory entry can be computed by adding the ranks of the previous first-level entry and the ranks of the last second-level entry associated that first-level entry. This can be parallelized with prefix sum, where the combining operator adds all σ ranks in parallel. In the reduction tree, each combined entry e keeps track of the sum of ranks of first-level entries it represents (e_s) as well as last second-level entry e_t of its last first-level entry. When combining e with an entry e' before it in the input, we use we create a new combined entry e'' which stores the sum of e_s and e'_s as its sum, and e_t as its last second-level entry. Recall that there are only $n / (\sigma \log^2 n)$ first-level entries. Thus, this prefix sum takes $o(n)$ work and $O(\log n)$ depth.

Therefore, the construction of the rank structure takes linear work and $O(\log n)$ depth. Packing the entries into words to meet the space requirements can be done within the same

bounds.

Select. For select, a separate structure is stored for each character in the alphabet along with an array of size $O(\sigma)$ with pointers to the structures. The structure for character c stores the location of every $\sigma \log^2 n$ 'th occurrence of c in a first-level directory. Then for each range r between the locations, if $r \geq \sigma^2 \log^4 n$ then the answers in the range are directly stored. Otherwise we store the occurrence of every $\sigma(\log \log n)^2$ 'th of c in the range in a second-level directory. The sub-ranges in the second-level directory are then considered. If the length of a sub-range is at least $\sigma^3(\log \log n)^4$, then the answers are stored directly. Otherwise the range is small enough such that we can use a lookup table to find the k 'th occurrence of c in the range in $O(1)$ work. The lookup table can be constructed in $O(n^d)$ work for $d < 1$ and $O(\log n)$ depth.

We compute the entries for all σ structures together. To compute the first-level directory entries, we first split the input into chunks of size σ . Each chunk sequentially computes the number of occurrences of each character by reading characters one-by-one and updating the count associated with the character read as in [2]. This requires $O(n)$ work and $O(\sigma)$ depth. Then we perform a prefix sum on the results of the n/σ chunks, so that at the end of each chunk we know how many occurrences of each character there are up to its location. This requires a total of $O(n)$ work and $O(\log n)$ depth. With this information, we can compute which chunks the first-level directory entries lie in for each character, and store them in a packed array using prefix sum in $O(n)$ work and $O(\log n)$ depth. Note that for each character there can be at most one first-level entry per chunk since the chunk size is less than $\sigma \log^2 n$. In parallel for all chunks that contain a first-level entry, we go into the chunk and find it sequentially. Since the chunk is of length σ , this takes $O(n)$ work and $O(\sigma)$ depth. For each character c we can now compute the ranges r between first-level entries. We mark the chunks corresponding to the ranges $r \geq \sigma^2 \log^4 n$, as well as store the starting and ending point in the chunk. This takes $O(n)$ work and $O(1)$ depth. Using the initial prefix sum result, we can allocate the appropriate amount of space to store the entries that need to be stored directly inside these ranges, as well as compute appropriate offsets into the shared space among different chunks. Now we read the chunks in parallel and if the character read is inside a range (which can be checked in $O(1)$ work) we output its location to the appropriate offset in the corresponding select structure. This takes $O(n)$ work and $O(\sigma)$ depth.

We then similarly mark the chunks corresponding to the ranges $r < \sigma^2 \log^4 n$, and output every $\sigma(\log \log n)^2$ 'th occurrence of a character. This gives the sub-ranges r' for each character c and we can again mark the chunks based on whether r' is at least $\sigma^3(\log \log n)^4$, and perform the appropriate operations which can be done in $O(n)$ work and $O(\sigma)$ depth using the same ideas as above.

The entire construction of the select structure takes $O(n)$ work and $O(\sigma + \log n)$ depth. Since $\sigma = o(\log n)$, the depth simplifies to $O(\log n)$. Again, packing the entries into words to satisfy the space requirements can be done in the same bounds.