

Computer Science Department  
VU Amsterdam, 1081HV Amsterdam



Bachelor Thesis

---

## Sorting Evolving Data: An Experimental Study

---

**Author:** Yasin Kerim Karyagdi (2745681)

*1st supervisor:* Ali Mehrabi  
*2nd reader:* Yasamin Nazari

*A thesis submitted in fulfillment of the requirements for  
the VU Bachelor of Science degree in Computer Science*

June 30, 2025

# Contents

<b>1</b>	<b>Abstract</b>	<b>3</b>
<b>2</b>	<b>Introduction</b>	<b>4</b>
2.1	General Framework . . . . .	4
2.2	Related Work . . . . .	4
2.3	Our contribution . . . . .	5
<b>3</b>	<b>Sorting Evolving Data</b>	<b>6</b>
3.1	Formal Setting . . . . .	6
3.2	Randomized Quicksort . . . . .	7
3.3	Blocksort . . . . .	8
3.4	Repeated Insertion Sort . . . . .	9
3.5	Improved Repeated Insertion Sort . . . . .	10
3.6	Possible Further Improved Repeated Insertion Sort . . . . .	11
3.7	Brief Overview Algorithmic Behavior . . . . .	12
<b>4</b>	<b>Experiments</b>	<b>13</b>
4.1	Experimental Setup . . . . .	13
4.2	Validating convergence time and behavior . . . . .	14
4.3	Validating Theoretical Bounds . . . . .	15
4.4	Effect of Block Size on Performance . . . . .	18
4.5	Performance of Repeated Quicksort Followed by Repeated Insertion Sort . .	23
<b>5</b>	<b>Conclusions</b>	<b>28</b>
	<b>References</b>	<b>29</b>
<b>A</b>	<b>Model Implementation</b>	<b>30</b>
<b>B</b>	<b>Algorithms Implementation</b>	<b>34</b>
<b>C</b>	<b>Running Experiments</b>	<b>39</b>

# 1 Abstract

We perform an experimental study on sorting algorithm within the evolving data model. In this model, the algorithms maintain a solution that is close to the true order at each time step. To this end, the algorithm is allowed to perform one comparison at each time step, where after each comparison a constant number  $\alpha$  of random swaps of consecutive elements occur in the true order, so the true order changes with time. Previous work studies two variants of quicksort and two variants of insertion sort. We build upon this work and see whether the theoretical claims regarding the algorithms hold in practice. Additionally, we see whether the claims regarding algorithm convergence and behavior hold with our implementation. Furthermore, we study the effect of the block size on the maintained Kendall tau distance. Lastly, we introduce a new variant of insertion sort and study how it compares with regards to quicksort and insertion sort.

We find that all algorithms converge within  $O(n^2)$ , even when starting from a maximal Kendall tau distance, and that the maintained Kendall tau and the steady behavior coincide with previous practical findings. Additionally, we find that for  $\alpha \in \{1, 2, 10\}$  the theoretical bounds clearly hold in practice for most algorithms, and that it is highly likely that the theoretical bound holds in practice for blocksort. Furthermore, we find that the block size does have an effect on the maintained Kendall tau distance and that this effect increases for higher values of the input size and the change rate. Moreover, we find that smaller block sizes perform better in general, but that it does not necessarily hold that smaller is better. There seems to be a trade-off between performance and robustness. Lastly, we find that our variant of insertion sort exhibits the expected behavior of both quicksort and insertion sort, and that it seems to be viable for cases where the change rate changes significantly throughout the run.

**Keywords:** Sorting, Evolving data, Quicksort, Insertion sort, Experimental

## 2 Introduction

In the classical paradigm, an algorithm is given an input, performs some computations, and returns an output within some time and space. We can reason about this time and space complexity in various ways, such as the best case, worst case and average case, which allows us to select the best algorithm for our use case. Though not all real world settings reflect this paradigm of input, computation, and output, which gives rise to new computational models such as the **evolving data** model introduced by Anagnostopoulos *et al.*[1].

### 2.1 General Framework

In the evolving data model, the given input changes while the algorithm performs the computations. At each time step, the algorithm accesses this changing input through **probes**, which is implemented in various ways depending on the problem setting. Furthermore, at each time step the algorithm gives some approximate output based on the information it received from probing the input and the goal of the algorithm. The output is approximately close to the "true" output, which is the output at a certain time given that the true input would have been given statically to the algorithm at that time. For example, for sorting algorithms, which we focus on in this paper, the goal is to return at each time step a sorted list that is approximately sorted according to the true ranking. Lastly, due to the nature of this model, an algorithm does not have a certain runtime, it runs indefinitely until it is terminated.

### 2.2 Related Work

This thesis builds on the previous theoretical and experimental work performed by Anagnostopoulos *et al.*[1] and by Besa Vial *et al.* [3][2] regarding sorting problems in the evolving data model.

The evolving data model was introduced by Anagnostopoulos *et al.*[1], who study sorting and selection with respect to a total ordered list that slowly changes. This means that the input does not have a static total order given to the algorithm but rather a dynamic total order that changes at each time step. This change is modeled by means of a small constant number of uniformly **random swaps** of pairwise adjacent elements, where we call the mechanism that performs the swaps the **uniform adversary**. The algorithm has access to the true total order through a constant number of **probes**, which is modeled by means of comparisons between pairwise adjacent elements in the true ordering. For this setting, they prove a lower bound of  $\Omega(n)$  on the expected Kendall tau distance between the estimated order computed by any algorithm and the actual order at any time  $t$ . Furthermore, they prove that repeatedly running quicksort maintains a Kendall tau distance of  $O(n \ln n)$ , in expectation and with high probability, with respect to the true total order. Lastly, they prove that alternating steps of quicksort with that of blocksort, a variant of quicksort, can maintain a Kendall tau distance of at most  $O(n \ln \ln n)$  in expectation and with high probability.

Besa Vial *et al.* [3][2] performed an experimental study to see how these algorithms perform in practice [3], and how they perform against sorting algorithms that have a quadratic runtime in the classical paradigm, specifically insertion sort, cocktail sort and

bubble sort. They found that the algorithms empirically reach a steady behavior within  $O(n^2)$ . Furthermore, they found that quadratic runtime algorithms maintain a lower Kendall tau distance in the steady behavior; however, they take longer to reach this state. Additionally, they introduce the hot spot adversary to see how the algorithms perform under such a setting. They provide evidence that these sorting algorithms have similar robustness against the **hot spot adversary** as they do against the uniform adversary.

Besa *et al.* [2] expanded on these results with a theoretical paper where they prove that repeatedly running insertion sort in this setting maintains a Kendall tau distance of at most  $O(n)$  in expectation and with high probability. That is, after an initialization period of  $\Theta(n^2)$  steps. Additionally, they show that running an initial run of quicksort improves this initialization time to  $\Theta(n \ln n)$  steps, with high probability, by first performing a run of quicksort and then following that with repeated insertion sort.

There are several other related works, but due to space limitations, we refer an interested reader to the related work sections found in the mentioned works of Anagnostopoulos *et al.*[1] and of Besa Vial *et al.* [3][2] for more related works.

### 2.3 Our contribution

In this thesis, we conduct an experimental investigation of sorting in the evolving data model.

We start by verifying the claim that all algorithms converge within  $O(n^2)$ [3] by investigating if this holds when we start from a maximal Kendall tau distance and run the algorithms for  $n^2$  time steps. Then we compare our findings regarding the steady behavior of quicksort, blocksort[1] and insertion sort [2] with the findings of Besa Vial *et al.* [3]. To see whether we still find similar results to those of Besa Vial *et al.* [3], even when starting from a maximal Kendall.

Afterwards, we verify whether the theoretical claims regarding the maintained Kendall tau distance hold in practice.

Then, we will experimentally investigate how the algorithm that alternates between steps of quicksort and of blocksort[1] performs for different block sizes under varying amounts for the input size and the change rate. This is done to see whether we can improve the practical results found by Besa Vial *et al.* [3] regarding blocksort.

Lastly, we investigate the new algorithm introduced in section 3.6 that expands on quicksort then repeated insertion sort [2], namely **repeated quicksort then insertion sort**. We investigate how it compares to the regular variant of insertion sort and the regular variant of quicksort for varying amounts for the input size and the change rate.

To this end, we formally introduce the setting, the relevant algorithms and their behavior in Section 3. In Section 4 we describe the experiments and present the results.

**Plagiarism declaration.** I confirm that this thesis work is my own work, is not copied from any other source (person, Internet, or machine), and has not been submitted elsewhere for assessment.

## 3 Sorting Evolving Data

### 3.1 Formal Setting

We now formally introduce the setting for sorting on a true total ordered list in the evolving data model formulated by Anagnostopoulos *et al.*[1]. From this point onward, when we mention an algorithm, we specifically refer to sorting algorithms that are run in the evolving data model.

#### Input and Output

The initial input of an algorithm is a finite domain  $U$  of  $n$  distinct elements. At each time step  $t$ , there exists a true total order  $\pi_t$  on  $U$ , and this total order changes over time. At each time step  $\pi_t$  is accessible to the algorithm through **probes**, which are comparisons between pairwise adjacent elements.

The algorithm uses probes to maintain an **approximation** of the order, denoted  $\tilde{\pi}_t$ . At each  $t$  the algorithm gives as output the current estimated ordering  $\tilde{\pi}_t$ , where  $\tilde{\pi}_t$  is a possible permutation over the  $n$  elements. The goal of the algorithm is to maintain an approximation  $\tilde{\pi}_t$  that is close to  $\pi_t$  at each  $t$ . To measure this closeness between  $\tilde{\pi}_t$  and  $\pi_t$ , we use the **Kendall tau** distance:

$$\text{KT}(\pi_1, \pi_2) = |\{(x, y) \mid x <_{\pi_1} y \wedge y <_{\pi_2} x\}|$$

This is the number of pairwise inversions between  $\tilde{\pi}_t$  and  $\pi_t$ , so in our case that is the number of pairs that are incorrectly ordered with respect to the true order. The maximal Kendall tau distance is reached when one list is reverse of the other, where we would have the Kendall tau distance equal to  $n \cdot (n - 1) / 2$ , so it is tightly bound from above by  $\Theta(n^2)$ .

#### Discrete Time Model

It is important to note that time is **discretized**, meaning we have steps  $t = 0, 1, 2, \dots$ . At each time step, an algorithm is allowed to perform a fixed number of **probes**  $p$ . We will only consider the case where  $p = 1$ , and therefore take the number of comparisons made as an indication for the time.

After each set of comparisons, the true order  $\pi_t$  changes to  $\pi_{t+1}$  through a sequence of **swaps**, depending on the adversary model. We are interested in the initial setting introduced by Anagnostopoulos *et al.*[1], therefore, we will only consider the **uniform adversary** [3]. The uniform adversary performs a constant number of **random swaps** depending on the **change rate**  $\alpha \geq 1$ . Each time, it chooses a pair of pairwise adjacent elements in  $\pi_t$  uniformly at random and swaps them.

## Metrics for Analyzing Algorithms

We do not have a definitive terminal state because an algorithm runs indefinitely in the evolving data model. Therefore, it is important to note the following two metrics to analyze the behavior of the algorithms running in the model: the **steady behavior** [3] and the **convergence time** [3]. For each algorithm the Kendall tau distance reaches a final repetitive behavior after a certain amount of time. This repetitive behavior is called the steady behavior, and the time required to reach it is called the convergence time. Every algorithm's empirical convergence time is at most  $O(n^2)$  steps[3].

Additionally, note that we can use the running time of algorithms in the classical paradigm in order to reason about its behavior and properties in the evolving data model.

### 3.2 Randomized Quicksort

In the classical paradigm, quicksort takes as an input an array of  $n$  elements and outputs the array sorted increasingly. To do so, it divides the list into two parts depending on the pivot, which is a chosen element of the array. Everything to the left of the pivot is smaller and everything to the right is bigger. Afterwards, it recursively runs quicksort on the left and right parts until the whole array is sorted. The base case for quicksort is a call in which the number of elements being sorted is less than two because there is nothing to sort in such a case. The worst-case time complexity of quicksort is  $O(n^2)$ . This occurs when the partitioning algorithm picks the smallest or largest element as the pivot in each step.

---

**Algorithm 1** Partition( $A, p, r$ )

---

**Input:** Array  $A$ , indices  $p$  and  $r$

```
1:  $x \leftarrow A[r]$ 
2:  $i \leftarrow p - 1$ 
3: for  $j \leftarrow p$  to  $r - 1$ 
4:   if  $A[j] \leq x$  then
5:      $i \leftarrow i + 1$ 
6:     Exchange  $A[i]$  with  $A[j]$ 
7:   end if
8: end for
9: Exchange  $A[i + 1]$  with  $A[r]$ 
10: return  $i + 1$ 
```

---

---

**Algorithm 2** Randomized\_Partition( $A, p, r$ )

---

**Input:** Array  $A$ , indices  $p$  and  $r$

```
1:  $i \leftarrow \text{RANDOM}(p, r)$ 
2: Exchange  $A[r]$  with  $A[i]$ 
3: return PARTITION( $A, p, r$ )
```

---

---

**Algorithm 3** Randomized\_Quicksort( $A, p, r$ )

---

**Input:** Array  $A$ , indices  $p$  and  $r$

```
1: if  $p < r$  then
2:    $q \leftarrow \text{RANDOMIZED\_PARTITION}(A, p, r)$ 
3:   RANDOMIZED_QUICKSORT( $A, p, q - 1$ )
4:   RANDOMIZED_QUICKSORT( $A, q + 1, r$ )
5: end if
```

---

In the evolving data model, we run quicksort repeatedly. Furthermore, the result of the comparisons are done using probes, whereby the result is based on the current true order. We update our current approximation when a quicksort run on the whole array terminates. We start a new run after updating the approximation. Given our setting, it is important that we use randomized quicksort, which is a variant of quicksort whereby the pivot element is randomly picked and moved into the pivot position, which is at the end of the current part of the array on which we called partition on. Randomized quicksort has an average case of  $O(n \ln n)$  and due to the random selection of the pivot we avoid the worst-case time complexity of  $O(n^2)$  with high probability. Furthermore, given our setting, it is not permissible that the recursive run of the algorithm on the left part of the array is interleaved with the recursive run on the right part. Anagnostopoulos *et al.* [1] proof that repeatedly running quicksort guarantees that for every time step  $t$  the distance between the orderings  $\pi_t$  and  $\tilde{\pi}_t$  is  $O(n \ln n)$ , with high probability.

### 3.3 Blocksor

We now introduce an algorithm that can maintain an error of  $O(n \ln \ln n)$ , with high probability [1]. This algorithm builds on the results obtained by quicksort and exploits its properties. Specifically, that after the quicksort run terminates, which due to its running time results in an error of  $\Omega(n \ln n)$ , the rank of each element in the approximation is within  $O(\ln n)$  of its actual rank. Therefore, we do not need to run an  $O(n \ln n)$  algorithm from scratch. Instead, we can run several ( $O(n/\ln n)$ ) local quicksorts on **blocks** of size  $m = \Theta(\ln n)$  to correct the ordering. We interleave steps of these local quicksorts with steps of a full quicksort in order to obtain the resulting algorithm that can maintain a distance of  $O(n \ln \ln n)$ . The blocksort algorithm is responsible for performing the local quicksorts.

---

#### Algorithm 4 BlockSort( $\rho$ )

---

**Input:** Permutation  $\rho$  = output of full Quicksort at time  $t_0$

```

1:  $B$  : array of size  $m$ 
2: for  $j = 1$  to  $m/2$ 
3:    $B(j) \leftarrow \rho^{-1}(j)$ 
4: end for
5: for  $i = 1$  to  $2n/m - 1$ 
6:   for  $j = 1$  to  $m/2$ 
7:      $B(m/2 + j) \leftarrow \rho^{-1}(i \cdot m/2 + j)$ 
8:   end for
9:   QUICKSORT( $B$ )
10:  for  $j = 1$  to  $m/2$ 
11:     $\sigma^{-1}((i - 1) \cdot m/2 + j) \leftarrow B(j)$ 
12:     $B(j) \leftarrow B(m/2 + j)$ 
13:  end for
14: end for
15: for  $j = 1$  to  $m/2$ 
16:    $\sigma^{-1}(n - m/2 + j) \leftarrow B(j)$ 
17: end for
```

---



Blocksort takes as input the last full run of quicksort, so an approximation where each element is within  $O(\ln n)$  of its actual rank. The algorithm divides the array into overlapping blocks of size  $m$ , which is  $O(\ln n)$ , and sorts them from left to right. After sorting a block, we store the smallest  $m/2$  elements to the result and create a new block with the largest  $m/2$  elements of the current block and the next  $m/2$  elements of the array. Thus, the blocks overlap with  $m/2$  elements, this is done to counteract the possible case where elements move towards neighboring blocks during execution.

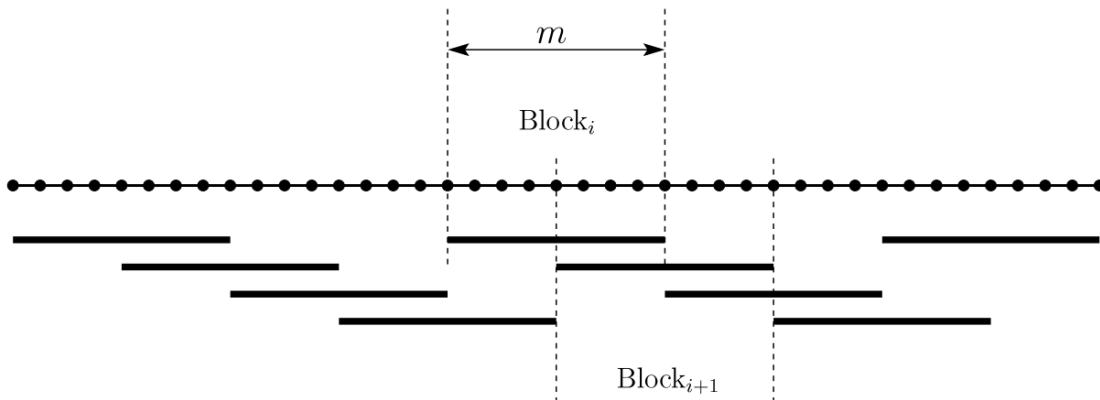


Figure 1: The overlap of blocks with their neighboring blocks.

Another possible issue is that elements eventually move far from their actual rank when we repeatedly run blocksort. To solve this, we interleave executions of blocksort with executions of regular quicksort. The two processes are executed independently with their own data structures. On even time steps we execute steps of blocksort and on odd steps those of quicksort. After each terminated run of blocksort we update our current approximation with the result obtained, and after that we start a new blocksort run with the last output of a full quicksort. Note that we start with an initial full quicksort run before being able to run blocksort.

### 3.4 Repeated Insertion Sort

In the classical paradigm insertion sort takes as an input an array of  $n$  elements and outputs the array sorted increasingly. To do so, it compares each element to the elements on its left and insert it in its correct position. We start this process from left to right and thus slowly sort the array in that order. The current element that is inserted is the element at position  $j$ . Inserting into the right position is done by means of swapping when the element to the left of it is bigger than the current element and by deducting  $j$ , until either there is nothing to the left of the element or the element to the left of it is smaller. The running time of insertion sort is  $O(n + I)$ , where  $I$  is the number of initial inversions in the list. Therefore, insertion sort runs  $O(n)$  in its best case, which is an already sorted array, and  $O(n^2)$  in its worst and average cases.

---

**Algorithm 5** Repeated Insertion Sort

---

**Input:** Array  $A$  of  $n$  elements

```
1: while true
2:   for  $i \leftarrow 1$  to  $n - 1$ 
3:      $j \leftarrow i$ 
4:     while  $j > 0$  and  $A[j] < A[j - 1]$ 
5:       Swap  $A[j]$  and  $A[j - 1]$ 
6:        $j \leftarrow j - 1$ 
7:     end while
8:   end for
9: end while
```

---

In the evolving data model we run insertion sort repeatedly. Furthermore, the result of the comparisons are done using probes, whereby the result is based on the current true order. The swaps are performed in the current approximation. So we update the current approximation directly, instead of waiting for a run of insertion sort to terminate. Note that we no longer have the property that everything to the left of an inserted element is smaller during sorting. This is because the true order changes each time step and therefore an already sorted part can become unsorted by some swaps in the true order.

Besa Vial *et al.* proof that repeatedly running insertion sort can for every time step  $t$  maintain a distance of  $O(n)$  between the orderings  $\pi_t$  and  $\tilde{\pi}_t$ , with high probability. This is the asymptotically optimal distance as proven by Anagnostopoulos *et al.*[1]. Furthermore, this linear upper bound seems to hold for higher rates of  $\alpha > 1$ [2], where Besa Vial *et al.*[3] show that it holds for  $\alpha \in \{1, 2, 10\}$ .

### 3.5 Improved Repeated Insertion Sort

Besa Vial *et al.*[2] show that the convergence time of repeated insertion sort can be improved from  $\Theta(n^2)$  to  $\Theta(n \ln n)$  by starting with a run of quicksort. This is because quicksort takes  $\Theta(n \ln n)$  comparisons to finish and because of this the number of inversions left after a run is  $O(n \ln n)$ .

---

**Algorithm 6** Quicksort Followed by Repeated Insertion Sort

---

**Input:** Array  $A$  of  $n$  elements

```
1: QUICKSORT( $A$ )
2: while true
3:   for  $i \leftarrow 1$  to  $n - 1$ 
4:      $j \leftarrow i$ 
5:     while  $j > 0$  and  $A[j] < A[j - 1]$ 
6:       Swap  $A[j]$  and  $A[j - 1]$ 
7:        $j \leftarrow j - 1$ 
8:     end while
9:   end for
10: end while
```

---

### 3.6 Possible Further Improved Repeated Insertion Sort

We now introduce our attempt at an improvement. Besa Vial *et al.*[3] found that for higher change rates  $\alpha$  quicksort performs better, this is most likely due to the insertion sort fixing at most one inversion each probe, whereby we are limited to one probe each time step, while  $\alpha$  swaps occur, each of which can either introduce or fix an inversion. Therefore, it is more dependent on the swaps that occur instead of the inversion fixed by the insertion sort. The algorithm attempts to improve the performance of Algorithm 6 for higher values of  $\alpha$  by running quicksort after a certain amount of runs of insertion sort, namely  $k$  runs. This could result in the average behavior of the two algorithms, where we can pick a higher  $k$  if we want the insertion sort behavior to dominate.

---

**Algorithm 7** Repeated Quicksort Followed by Repeated Insertion Sort

---

**Input:** Array  $A$  of  $n$  elements and an integer  $k$

```

1: while true
2:   QUICKSORT( $A$ )
3:    $m \leftarrow 0$ 
4:   while  $m < k$ 
5:     for  $i \leftarrow 1$  to  $n - 1$ 
6:        $j \leftarrow i$ 
7:       while  $j > 0$  and  $A[j] < A[j - 1]$ 
8:         Swap  $A[j]$  and  $A[j - 1]$ 
9:          $j \leftarrow j - 1$ 
10:      end while
11:    end for
12:     $m \leftarrow m + 1$ 
13:  end while
14: end while

```

---

We may have that the algorithm performs worse than repeated insertion sort for low values of  $\alpha$ , while performing worse than repeated quicksort for high values of  $\alpha$ . But it could still be a viable option if it performs better than repeated quicksort for low values of  $\alpha$  and better than repeated insertion sort for high values of  $\alpha$ . This would make it perform relatively well in systems where  $\alpha$  fluctuates constantly and where we are unsure of the current  $\alpha$  or where it costs resources to switch from one algorithm to the other. Our hypothesis is that higher values of  $k$  should maintain a lower Kendall tau distance compared to lower values of  $k$  for lower values of  $\alpha$ , while the reverse should also hold.

### 3.7 Brief Overview Algorithmic Behavior

For clarity we summarize the different possible ways of updating, the effect it has on the steady behavior, and the theoretical bounds of each algorithm.

In terms of updating, the algorithms we consider can be divided into three categories: **periodic**, **instant** and **mixed**.

Periodic algorithms update the approximation when the current run is terminated, after which they start a new run. At each time step, the algorithm gives as output the most recently updated approximation. So, while the algorithms are computing the Kendall tau distance is determined by the random swaps performed by the uniform adversary, because the output is the same while the true order changes. In the steady behavior, this leads to a slow increase in the Kendall tau distance, until the current run terminates, and the Kendall tau distance drops. This results in an oscillating behavior.

Instant algorithms can update the approximation at each time step, depending on whether an inversion has been found with the probe. The maintained Kendall tau distance at each time step is determined by the inversions fixed by the algorithm and by the random swaps performed by the uniform adversary. In the steady behavior, this leads to a somewhat straight line around a certain Kendall tau distance, which is determined by the probability that inversions are fixed or introduced.

Mixed algorithms make use of both types of updating. This is because they alternate between two different algorithms where one updates periodically and the other updates instantly. Repeated Quicksort Followed by Repeated Insertion Sort falls under this category because it repeatedly runs one quicksort run followed by a certain number of insertion sort runs. In the steady behavior, this should result in phases where it displays the behavior of quicksort and phases where it displays the behavior of insertion sort.

Algorithm	Update type	Theoretical bounds
Repeated Quicksort	Periodic	$O(n \ln n)$
Repeated Blocksrt	Periodic	$O(n \ln \ln n)$
Repeated Insertion sort	Instant	$O(n)$
Quicksort Followed by Repeated Insertion Sort	Mix	$O(n)$
Repeated Quicksort Followed by Repeated Insertion Sort	Mix	?

Table 1: Summary of algorithms, update types, and theoretical bounds.

## 4 Experiments

The main goal of our experimental framework is to address the following questions:

1. Do all the algorithms converge within  $O(n^2)$  when starting from a maximal Kendall tau distance? If so, do they converge to Kendall tau distances similar to those found by Besa Vial *et al.* [3]?
2. Do the various theoretical claims regarding the maintained Kendall tau distance hold in practice?
3. What is the effect of the block size on the maintained Kendall tau distance in the steady behavior for varying input sizes and change rates?
4. How does **Repeated Quicksort Followed by Repeated Insertion Sort** perform in the steady behavior compared to repeated quicksort and repeated insertion sort for varying input sizes and change rates?

We present results below that address each of these questions.

### 4.1 Experimental Setup

We implemented the various algorithms in Python<sup>1</sup>. Randomness for experiments was generated using randint from the Python library random. The code was run in Python 3.10.12, and the computer running it was a remote machine with the following CPU: AMD EPYC-Milan (48) @ 2.399GHz.

We take the time measurement as the number of comparisons performed because the probe rate is equal to one in all experiments. For each configuration, the different algorithms ran on the same seeds. This is so that the changes in the true ordering were the same for each algorithm given a certain configuration.

To evaluate algorithm performance, we computed the average Kendall tau distance across multiple runs. Specifically, at each time step, we calculated the average distance over all runs with different seeds. This averaging process provides a clearer view of the general performance trends of each algorithm.

Due to time constraints we used 30 seeds for the experiment performed in 4.3, while using 100 seeds in all other experiments. We terminate a run after  $n^2$  comparisons and during a run the Kendall tau distance is measured every  $n/20$  comparisons for the experiment in Section 4.2 and  $n/10$  for the rest of the experiments.

In Section 4.2 the choice of a sampling rate of  $n/20$  and a running time of  $n^2$  is the same as that of Besa et al. [3]. We chose to do so to compare our results with theirs. Note that while Besa *et al.* [3] evaluate their results on a single seed, we averaged 100 seeds to better capture the overall behavior.

---

<sup>1</sup>The code can be found on Github on the following link: <https://github.com/YasinKaryagdi/Thesis>

## 4.2 Validating convergence time and behavior

We begin by answering the first question, with the additional goal of comparing our findings with those of Besa Vial *et al.*[3]. To that end, we choose the same block size, that is, the block size is the first even number larger than  $10 \ln n$  that divides  $n$  [3]. To validate the algorithmic behavior, we chose a reverse-sorted list as input, so that the Kendall tau distance is maximal at the start. We do so because this is the worst starting point in terms of the distance between the true order and the approximation. We want to see if the algorithms still reach the steady behavior within  $O(n^2)$ [3]. Note that this differs from the choice of Besa et al.[3]. They performed this experiment with a uniformly random shuffled list, resulting in a Kendall tau distance that is on average half the maximal Kendall tau distance (in our case of  $n = 1000$ , it is equal to 250000). Additionally, we want to see how Algorithm 7 behaves compared to the rest and if it will also reach the steady behavior within  $O(n^2)$ , therefore we also run it for this setup.

We start with a reverse-sorted list and run the algorithms for the case where  $n = 1000$  and  $\alpha = 1$  for  $n^2$  time steps. Figure 2 shows the results of the whole run and Figure 3 shows the behavior at the endpoint, which is the last 5%. Lastly, Figure 4 shows the results obtained by Besa Vial *et al.*[3].

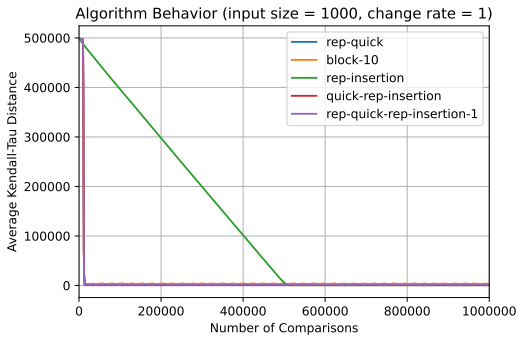


Figure 2: Algorithm behavior for  $n^2$  time steps, for  $n = 1000$  and  $\alpha = 1$

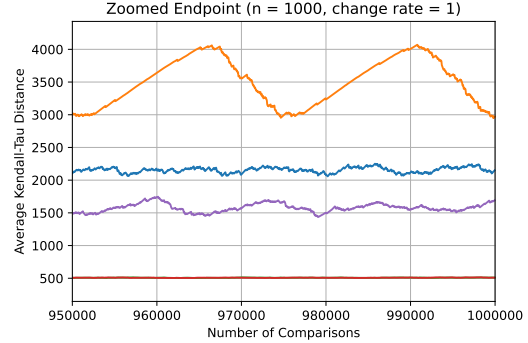


Figure 3: Algorithm behavior at the endpoint for  $n = 1000$  and  $\alpha = 1$

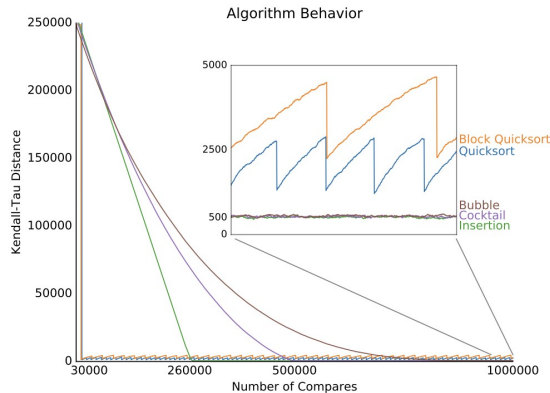


Figure 4: Algorithm behavior at the for  $n = 1000$  and  $\alpha = 1$  found by Besa Vial *et al.*[3] starting from a randomly shuffled list.

We find that all algorithms reach the steady state within  $O(n^2)$ , even when starting from the maximal Kendall tau distance. Additionally, we find that running quicksort before repeatedly running insertion sort improves the convergence time, while maintaining the same Kendall tau distance in the steady state. Furthermore, we find that the quicksort variants all oscillate, while insertion sort maintains a constant Kendall tau distance. Lastly, we find that Algorithm 7 falls between quicksort and insertion sort in terms of maintained Kendall tau distance, while displaying oscillating behavior similar to quicksort.

The steady behavior and the Kendall tau distance maintained in this state are similar to the results obtained by Besa Vial *et al.*[3]. Note that we differ in the amount of time that it takes insertion sort to reach the steady behavior; we find that it takes longer. This can be explained by the fact that we start with a different Kendall tau distance.

The second difference is in the behavior itself. Their plot shows the drop that occurs when a blocksort or quicksort run terminates, which looks like a straight line from a peak to a valley. We also see these peaks and valleys, but they are less distinct and the distance between a peak and a valley is smaller than what Besa Vial *et al.*[3] find. This is most likely because we plot the Kendall tau distance averaged over 100 runs, whereas they plot it for only one run. This explains the difference in behavior and in the exact Kendall tau distances maintained within the steady state.

From all these findings, we conclude that all algorithms converge within  $O(n^2)$  when starting from a maximal Kendall tau distance, and that the Kendall tau distances found in this case are similar to those found by Besa Vial *et al.*[3], who started from a randomly shuffled list.

### 4.3 Validating Theoretical Bounds

We are now interested in validating the theoretical bounds mentioned in Section 3.7 by looking at whether they hold in practice. To this end, we run the same algorithms from the previous experiment for input sizes between 100 and 10000. However, note that we exclude Algorithm 7 because it does not have a theoretical bound. From the previous experiment, we find that all algorithms converge within  $O(n^2)$  and therefore we only consider the Kendall tau distance maintained in the last 10% in each run, which we consider to be the endpoint. Furthermore, we only consider  $\alpha \in \{1, 2, 10\}$  because Besa Vial *et al.* [3] have shown that the theoretical bounds hold for these values of the change rates. Lastly, we start with a reverse-sorted list in order to ensure that we start from the maximal Kendall tau distance.

In the following, we plot the average Kendall tau distance over the endpoint divided by the input size against the input size. We rescale the average Kendall tau distance because the theoretical bounds are all divisible by  $n$ , and dividing by  $n$  makes it easier to identify the bounds.

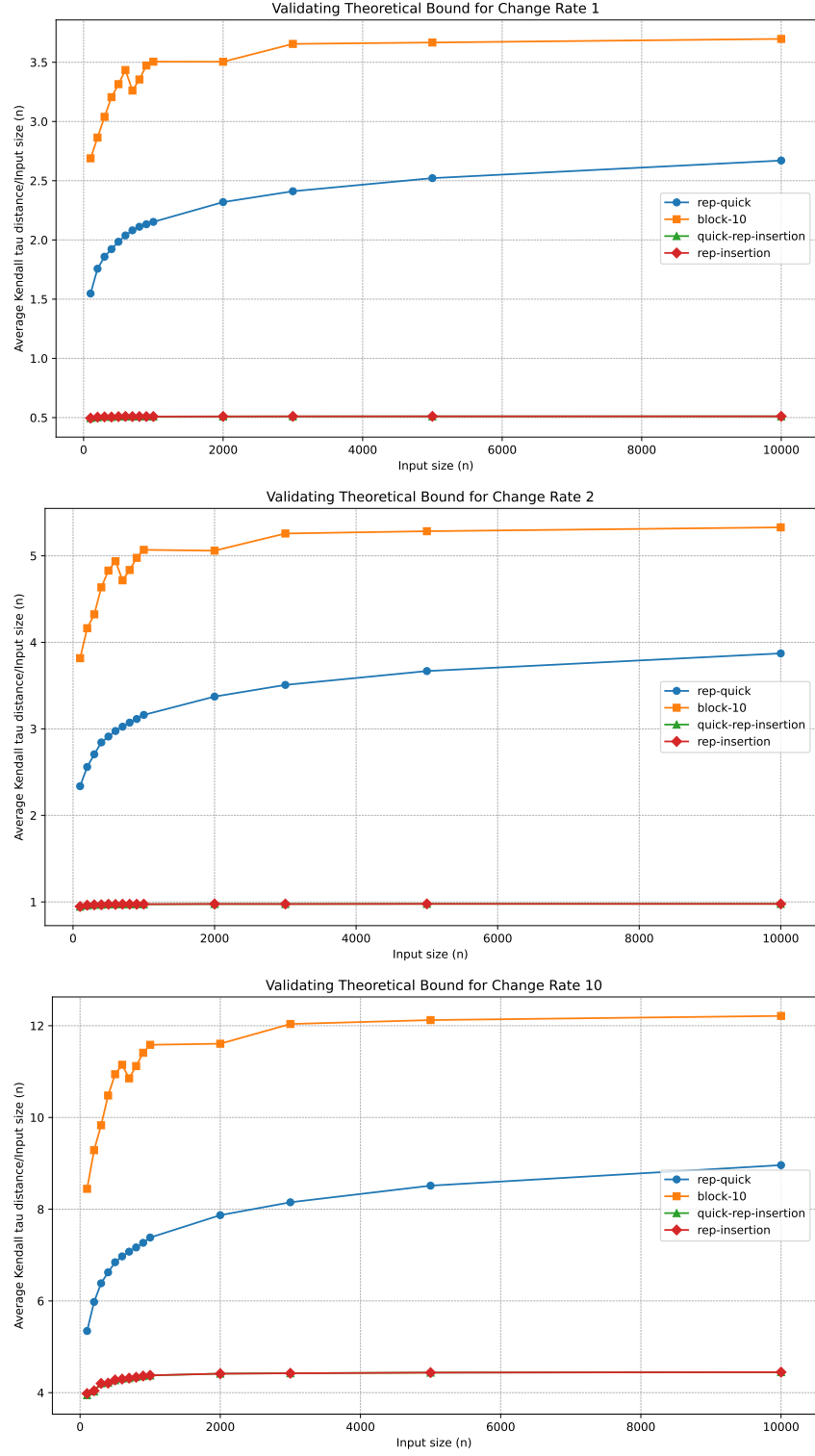


Figure 5: The attained average maintained Kendall tau distance over the endpoint divided by the input size for various input sizes and for increasing change rates.



We see that for all  $\alpha \in \{1, 2, 10\}$  we have that both repeated insertion sort, as well as quicksort followed by repeated insertion sort, maintain a rescaled average Kendall tau distance of  $O(1)$  in the endpoint. This means that they maintain an average Kendall tau distance of  $O(n)$ . So, we find that the theoretical bound holds in practice for these algorithms.

Additionally, we see that for all  $\alpha \in \{1, 2, 10\}$  we have that repeated quicksort maintains a rescaled average Kendall tau distance of  $O(\ln n)$  in the endpoint. This means that it maintains an average Kendall tau distance of  $O(n \ln n)$ , so we again find that the theoretical bound holds in practice.

We now look at blocksort. First we note that there seems to be an initial drop after  $n = 600$ , this is most likely due to the way we pick block sizes, which causes the block sizes to be inconsistent over the various input sizes. This leads to this case where the block size of  $n = 600$  is 100, while that of  $n = 700$ ,  $n = 800$  and  $n = 900$  are respectively 70, 80 and 90. This affects our ability to clearly distinguish whether the theoretical bound holds in practice. However, despite this we see that blocksort appears to be a logarithmic algorithm for all  $\alpha \in \{1, 2, 10\}$ . It is unclear from these results whether it is double logarithmic. This is because double logarithmic algorithm are supposed to have a slower increase in the average maintained Kendall tau distance compared to quicksort, which is a logarithmic algorithm. However, we find that the increase is faster compared to a logarithmic algorithm until an input size of  $n = 3000$ , after which it seems to have a slower increase. So, from  $n = 3000$  onward it appears to be double logarithmic, meaning  $O(\ln \ln n)$ , which would imply that the theoretical bound of  $O(n \ln \ln n)$  holds in practice, but this does not seem to hold for the input sizes smaller than  $n = 3000$ .

There could be multiple reasons for this; the most likely candidate is the way we pick the block sizes. Moreover, we could also have that blocksort requires the input to be at least of a certain size before it maintains a Kendall tau distance of  $O(n \ln \ln n)$ . Another reason could be the way we average, where we take the average at the endpoint of 10%. We could have that taking the endpoint to be 10% for these input sizes does not accurately depict the average. Lastly, we could have that running for 30 seeds is not sufficient to accurately depict the average behavior. In any case, from these results we can not conclude for certain whether the theoretical bound of  $O(n \ln \ln n)$  holds in practice, but it seems highly likely because we see that from  $n = 3000$  onward the increase is lower than that of quicksort, which we identified to be logarithmic.

From these results, we conclude that for  $\alpha \in \{1, 2, 10\}$  the theoretical bounds clearly hold in practice for most algorithms, and that it is highly likely that the theoretical bound holds in practice for blocksort.

#### 4.4 Effect of Block Size on Performance

As mentioned previously, Anagnostopoulos *et al.*[1] do not specify a specific constant for the block size and only mention that it should be of size  $O(\ln n)$ . We are interested in what effect the block size has on the performance of the algorithm in terms of the Kendall tau distance measured in the steady behavior.

To investigate this, we introduce the block constant  $b$ , where the block size depends on the block constant. We chose the block size for blocksort to be the first even number larger than  $b \cdot \ln(n)$  that divides  $n$ , where  $b$  is the block constant specified in the experiment.

We run the algorithm mentioned in section 3.3 for  $b \in \{1, 5, 10, 20, 40\}$  and investigate its behavior at the endpoint. Each will be run for various input sizes, for  $n \in \{100, 500, 1000\}$ , and for various change rates, for  $\alpha \in \{1, 5, 10, 20\}$ . Note that it is possible that the block constants may resolve into the same block size or into block sizes close to each other due to our choice for the block size. Furthermore, note that we are only interested in the performance of the algorithms in their steady behavior and relative to each other, and less in their convergence time, so we start with a sorted array as the initial input and only consider the endpoints. The endpoints for the various values of  $n$  are as follows: the last 40% for  $n = 100$ , the last 10% for  $n = 500$  and the last 5% for  $n = 1000$ .

In the following, we present the results of the experiment using figures to show the average Kendall tau distance at each time step, and tables to show the average Kendall tau distance over the endpoint.

Note that not all lines seem to be plotted; this is due to some block constants that result in the same exact block size and thus the same behavior as others. This happens for  $b = 5$  and  $b = 10$  given that  $n = 100$ , and for  $b = 20$  and  $b = 40$  in the case where  $n = 100$  and  $n = 500$ .

	100	500	1000
<b>block-1</b>	4	10	8
<b>block-5</b>	50	50	40
<b>block-10</b>	50	100	100
<b>block-20</b>	100	250	200
<b>block-40</b>	100	250	500

Table 2: The exact block size given a block constant for the various input sizes  $n$ .

## Results for $n = 100$

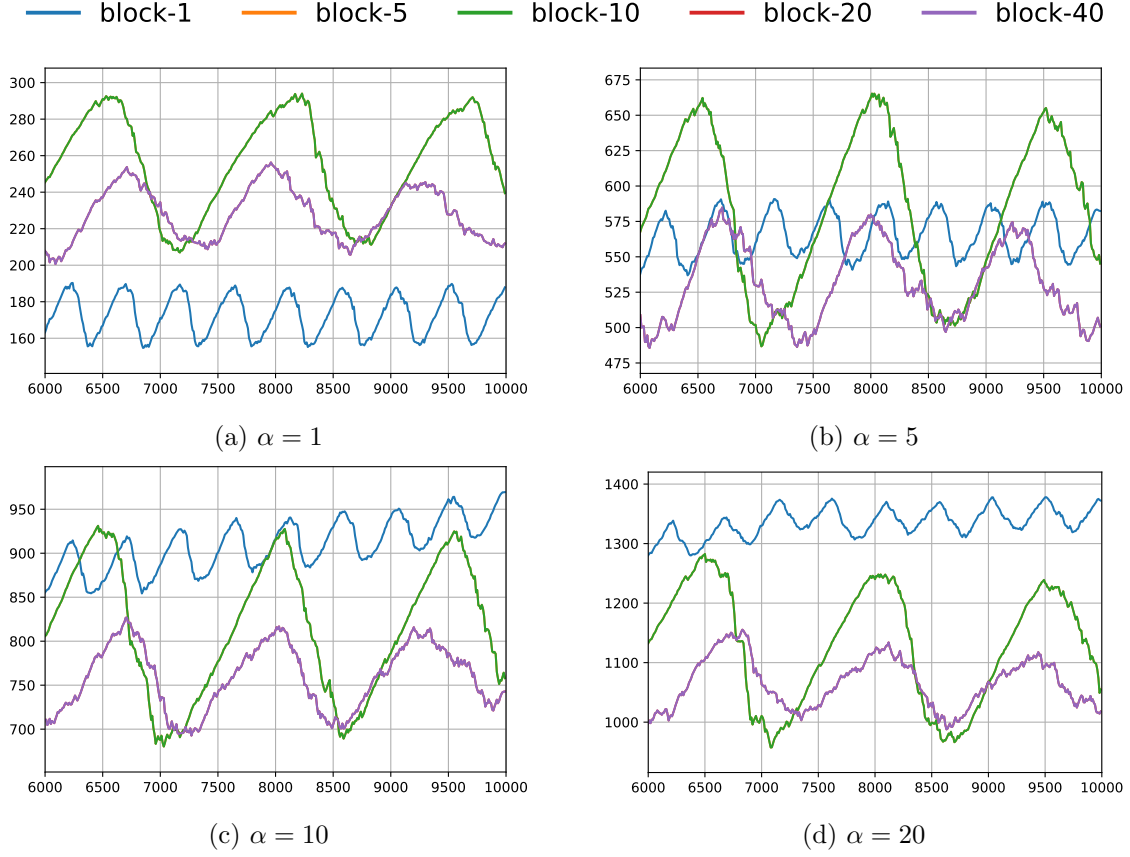


Figure 6: Comparison between blocksort performance across input sizes and alpha values for  $n = 100$ .

	1	5	10	20
<b>block-1</b>	172.79	565.78	909.19	1335.69
<b>block-5</b>	256.42	584.14	816.05	1130.99
<b>block-10</b>	256.42	584.14	816.05	1130.99
<b>block-20</b>	227.99	531.45	758.70	1066.19
<b>block-40</b>	227.99	531.45	758.70	1066.19

Table 3: The average Kendall tau distance for  $n = 100$  over the endpoint. Color coding: 1st, 2nd, 3rd.

## Results for $n = 500$

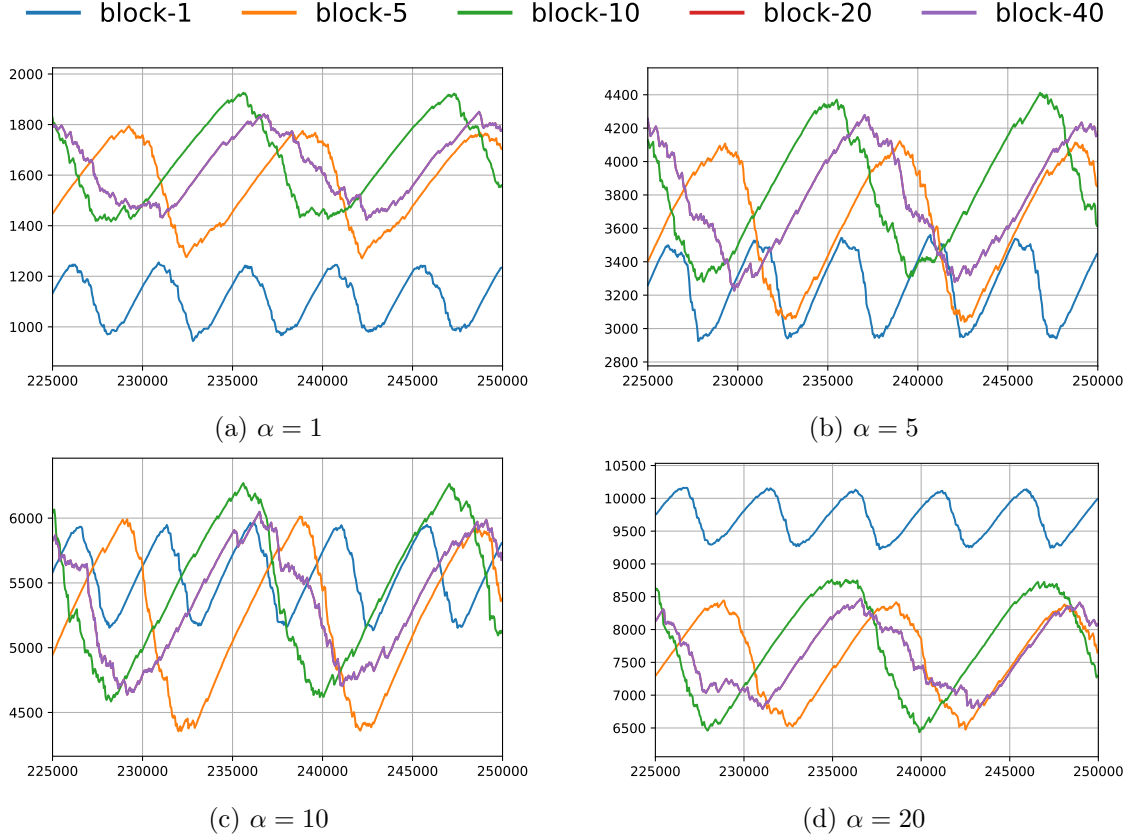


Figure 7: Comparison between blocksort performance across input sizes and alpha values for  $n = 500$ .

	1	5	10	20
<b>block-1</b>	1111.52	3254.90	5563.13	9708.67
<b>block-5</b>	1573.06	3659.18	5253.59	7599.51
<b>block-10</b>	1656.07	3841.70	5442.17	7748.60
<b>block-20</b>	1633.09	3773.40	5372.05	7614.21
<b>block-40</b>	1633.09	3773.40	5372.05	7614.21

Table 4: The average Kendall tau distance for  $n = 500$  over the endpoint. Color coding: 1st, 2nd, 3rd, 4th.

## Results for $n = 1000$

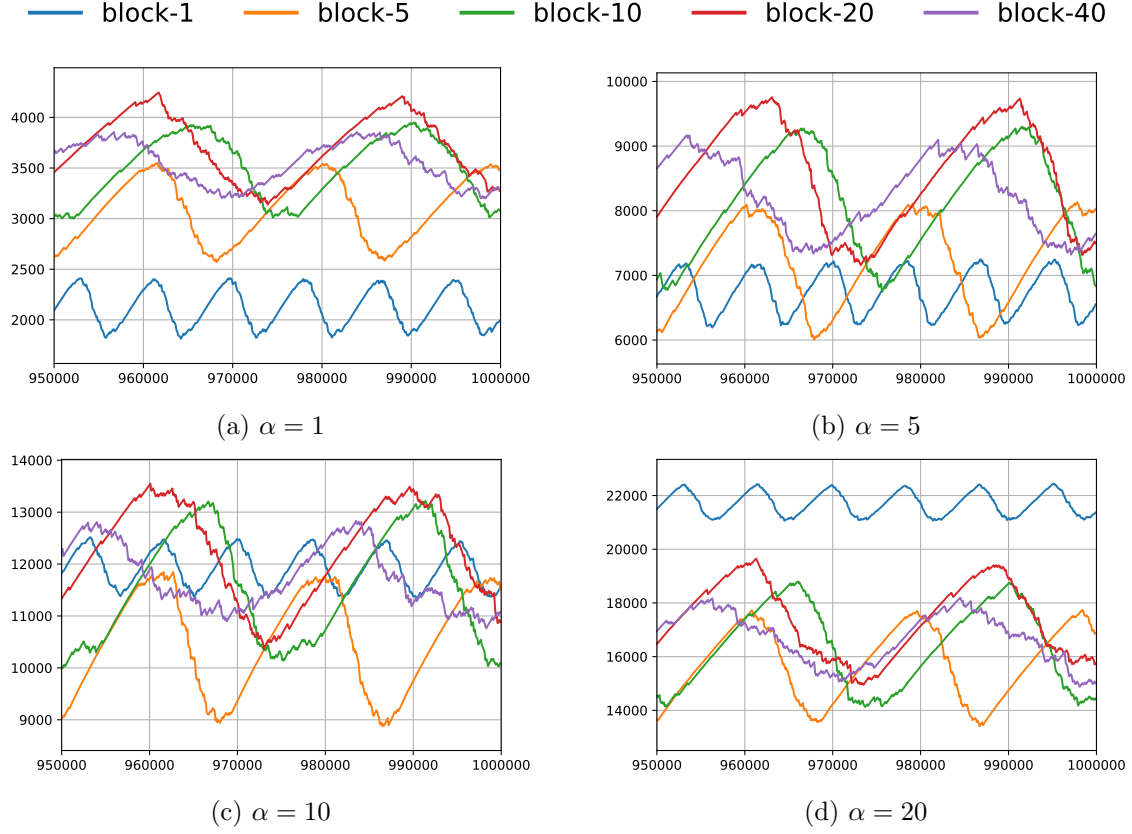


Figure 8: Comparison between blocksort performance across input sizes and alpha values for  $n = 1000$ .

	1	5	10	20
<b>block-1</b>	2123.13	6751.11	11926.26	21708.86
<b>block-5</b>	3095.18	7215.01	10531.16	15780.96
<b>block-10</b>	3499.70	8118.51	11550.61	16360.95
<b>block-20</b>	3720.99	8601.65	12252.52	17377.69
<b>block-40</b>	3547.91	8241.86	11747.07	16693.90

Table 5: The average Kendall tau distance for  $n = 1000$  over the endpoint. Color coding by rank: 1st, 2nd, 3rd, 4th, 5th. The differences between the different ranking at  $\alpha = 10$  is very small, this is the turning point for block-1. Block-20 consistently performs worse than the rest.

From Tables 3, 4 and 5 we see that the smallest block constant of 1 results in a better average Kendall tau distance maintained within the steady behavior for lower values of  $\alpha$ , namely  $\alpha = 1$  and  $\alpha = 10$ . Although it performs worse for the higher change rates, where it performs considerably worse than the other block constants. This trend is consistent over all the input sizes, and we have that the difference in performance between it and the others is more significant for higher input sizes. For example, in the case of  $\alpha = 1$  and  $n = 100$  it ranks first and there is a difference of 50 with the second ranking, while in the case of  $\alpha = 1$  and  $n = 1000$  this difference is 900, which is a significant improvement. This increase in effect also applies when it ranks last in the cases of  $\alpha = 10$  and  $\alpha = 20$ , where the differences become greater for larger input sizes. We can clearly see this difference in performance in the figures, where we see that it initially oscillates between Kendall tau distances below the other lines, but slowly starts to rise for higher  $\alpha$ , until it clearly oscillates above the other lines.

Moreover, interestingly, we find that in these cases of  $\alpha = 10$  and  $\alpha = 20$  the first ranking block constant is that of 5, our second smallest block constant, which ranked second in the other cases. From this we conclude that, in general, smaller block constants perform better given these low change rates, but that this increase in performance is less robust if the block constant is too small. This is most likely due to the elements shifting further from their true rank for higher values of  $\alpha$ . The purpose of blocksort is to counteract the  $O(\ln n)$  shift in ranking mentioned in Section 3.3. We could have that certain block sizes are better at counteracting this shift in ranking caused by the higher change rates. A small block size is most likely advantageous for a low change rate because the elements are less likely to shift too far, while it becomes a major disadvantage for higher change rates because the elements most likely constantly move to neighboring blocks, and thus are not properly moved back to their true rank.

Surprisingly, we find that a lower block constant does not always lead to better performance. This can be seen in Table 3, where the two largest block constants rank first or second, and Table 4, where they rank second or third. Furthermore, this is also the case in Table 5 where we can see that a block constant of 40 performs better than that of 20. From this we can see that a lower block size does not necessarily correspond to an improvement, though it is unclear currently why this is the case.

Lastly, we note that the differences in the average Kendall tau distance maintained are quite close in most cases. The only exception to this is a block constant of 1 in the case where  $\alpha = 1$ , where it performs significantly better than the rest, and  $\alpha = 20$ , where it performs significantly worse than the rest. Interestingly, from the figures we see that the different block constants exhibit similar but different behaviors, even though their averages are close to each other. The larger block constants oscillate less compared to the smaller ones, and they also have different maxima and minima.

From these results, we conclude that the block size has an effect on the Kendall tau distance maintained and that this possible improvement increases for higher values of  $n$ . Although there does seem to be a trade-off between performance and robustness, where a lower block constant performs better for lower change rates but is not robust for higher change rates. Lastly, a smaller block size does not always perform better compared to a bigger block size.

## 4.5 Performance of Repeated Quicksort Followed by Repeated Insertion Sort

Recall the potential improvement mentioned in Section 3.6. We are interested in how it compares to repeated quicksort and to repeated insertion sort in the steady behavior for differing values of the input size and the change rate. Furthermore, we are interested in how different amounts for  $k$  affect this attained behavior. Therefore, we consider running Algorithm 7 for  $k \in \{1, 2\}$ , where  $k$  is the number of insertion sort runs that we perform after a quicksort run. Once more, we are only interested in the performance of the algorithms in their steady behavior and relative to each other, so we start with a sorted array as the initial input and only consider the endpoints. The endpoints for the various values of  $n$  are as follows: the last 40% for  $n = 100$ , the last 10% for  $n = 500$  and the last 5% for  $n = 1000$ . The only exception is that we consider the last 40% in the cases where  $\alpha = 250$  to properly display the steady behavior. Furthermore, note that we only consider Algorithm 5, as opposed to both Algorithm 5 and Algorithm 6. This is because both maintain the same Kendall tau distance in the steady behavior, so either is a valid choice.

We run each algorithm for various input size  $n$ , namely  $n \in \{100, 500, 1000\}$ , and reasonable change rates, namely  $\alpha \in \{1, 5, 10, 20\}$ , plus an extreme case, namely  $\alpha = 250$ .

In the following, we present the results of the experiment using figures to show the average Kendall tau distance at each time step, and tables to show the average Kendall tau distance over the endpoint.

## Results for $n = 100$

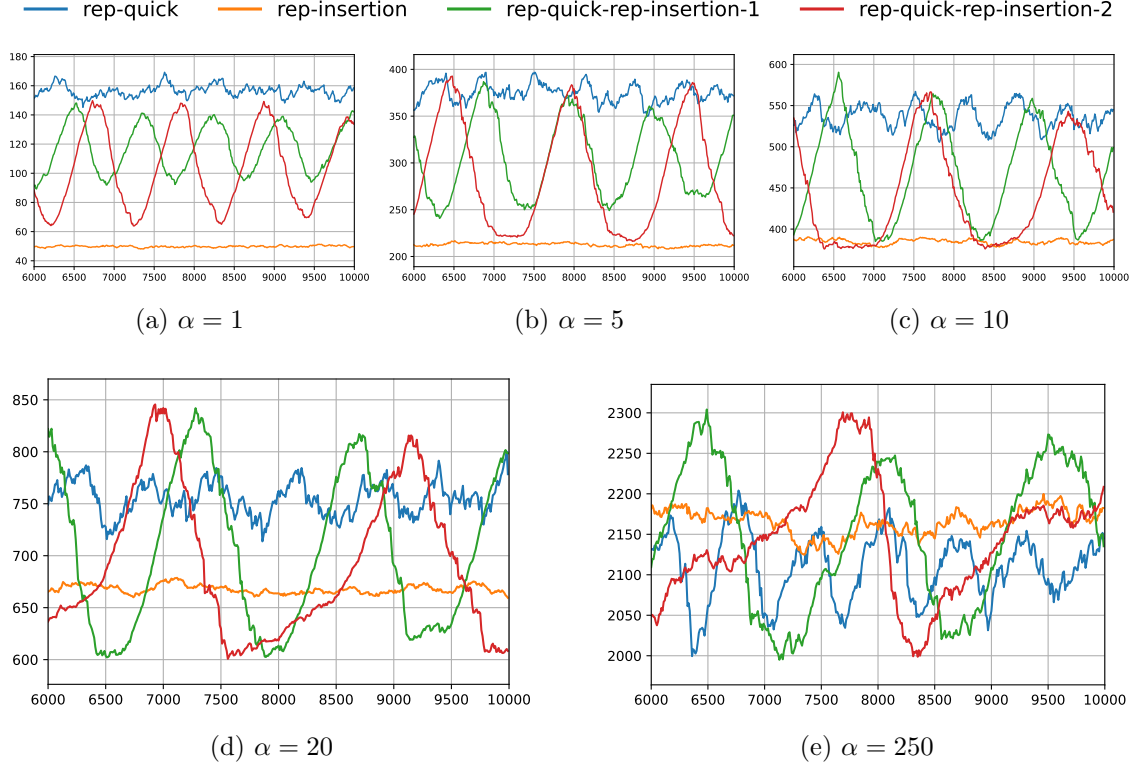


Figure 9: Comparison between blocksort performance across input sizes and  $\alpha$  values for  $n = 100$ .

	1	5	10	20	250
rep-quick	156.44	374.99	536.17	755.55	2109.81
rep-insertion	49.80	212.52	384.73	667.87	2165.24
rep-quick-rep-insertion-1	117.17	303.53	467.05	702.63	2151.45
rep-quick-rep-insertion-2	103.68	287.97	442.00	693.11	2144.20

Table 6: The average Kendall tau distance for  $n = 100$  over the endpoint. Color-coded per column: 1st, 2nd, 3rd, 4th.



## Results for $n = 500$

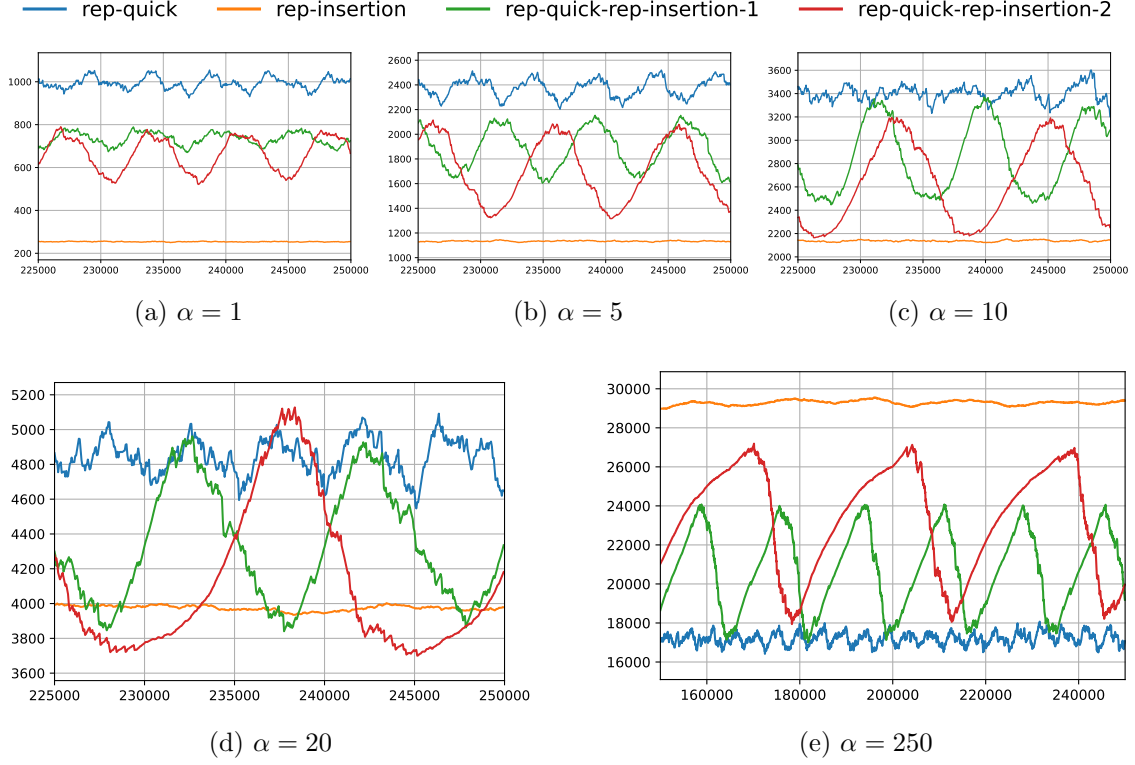


Figure 10: Comparison between blocksort performance across input sizes and  $\alpha$  values for  $n = 500$ .

	1	5	10	20	250
rep-quick	993.54	2375.23	3400.54	4846.21	17202.95
rep-insertion	254.00	1133.52	2136.98	3974.20	29285.92
rep-quick-rep-insertion-1	734.25	1881.40	2850.26	4320.68	20781.91
rep-quick-rep-insertion-2	665.83	1717.16	2631.10	4101.72	23389.68

Table 7: The average Kendall tau distance for  $n = 500$  over the endpoint. Color-coded per column: 1st, 2nd, 3rd, 4th.

## Results for $n = 1000$

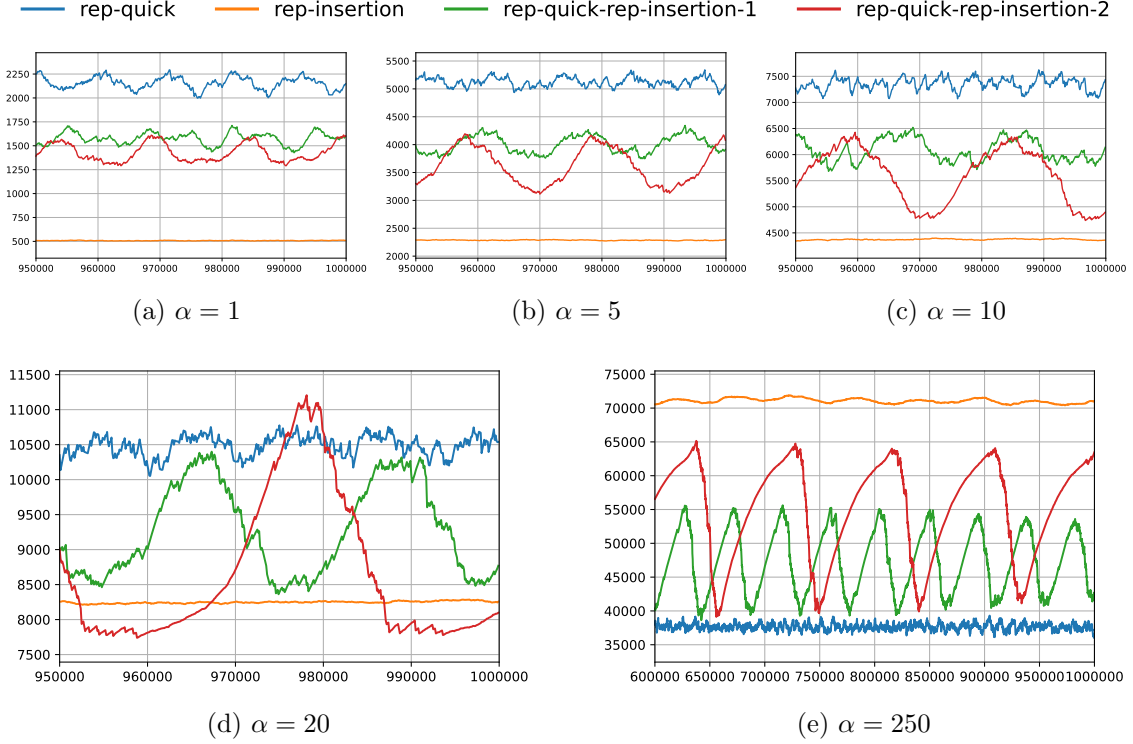


Figure 11: Comparison between blocksort performance across input sizes and  $\alpha$  values for  $n = 1000$ .

	1	5	10	20	250
<b>rep-quick</b>	2156.34	5126.42	7353.53	10486.82	37658.42
<b>rep-insertion</b>	508.41	2285.63	4374.12	8247.88	71064.09
<b>rep-quick-rep-insertion-1</b>	1576.71	4001.89	6082.36	9262.29	47305.99
<b>rep-quick-rep-insertion-2</b>	1427.56	3621.57	5596.82	8631.01	55274.90

Table 8: The average Kendall tau distance for  $n = 1000$  over the endpoint. Color-coded per column: 1st, 2nd, 3rd, 4th.

We start by looking at the performance of repeated quicksort and of insertion sort, and reason about how those affect Algorithm 7. Tables 6, 7, and 8 show that for all  $\alpha \in \{1, 5, 10, 20, 250\}$  there is a gap between the maintained Kendall tau distance of repeated quicksort and that of repeated insertion sort, and that this gap is larger for higher values of  $n$ . We find that for all  $\alpha \in \{1, 5, 10, 20\}$  repeated insertion sort performs best, while repeated quicksort performs worst. The reverse holds for  $\alpha = 250$ . Additionally, we see that the gap between the performance of these two algorithms is slowly narrowed, and we find that there is most likely a certain point between  $20 < \alpha < 250$  where repeated quicksort performs better than repeated insertion sort, and most likely continues doing so from that point onward.

We now look at how this affects the ranking and the averages of Algorithm 7, namely for  $k \in \{1, 2\}$ . We see that they consistently rank third and fourth. We have that running with  $k = 2$  outperforms  $k = 1$  in all the cases where repeated insertion sort ranks first, while  $k = 1$  outperforms  $k = 2$  in the case that repeated insertion sort ranks last. This is to be expected; as mentioned in Section 3.7 we slowly converge towards a certain Kendall tau distance when we repeatedly run insertion sort. Although in the case of Algorithm 7 we most likely do not reach that Kendall tau distance because we switch over to quicksort at a certain point, which in most cases seems to be too soon for insertion sort to have converged. However,  $k = 2$  spends more time running insertion sort and therefore gets closer to that Kendall tau distance. This is a benefit when this Kendall tau distance is lower than that maintained by quicksort, while it is a disadvantage when it is higher, which is the case when  $\alpha = 250$ . This explains how  $k = 1$  and  $k = 2$  obtain their rankings and why they switch at  $\alpha = 250$ .

From Figures 9, 10 and 11 we clearly see that for  $\alpha \in \{1, 5, 10\}$  both  $k = 1$  and  $k = 2$  slowly reach the Kendall tau distance maintained by repeated insertion sort, while we see that  $k = 2$  gets closer in all cases. From Figure 9c, we can clearly see what happens when they do reach this Kendall tau distance. Here we find that both  $k = 1$  and  $k = 2$  reach the maintained Kendall tau distance of insertion sort, namely  $\approx 400$ , and that  $k = 2$  maintains this longer than  $k = 1$ .

Furthermore, from Figures 9d, 10d and 11d we see that both  $k = 1$  and  $k = 2$  maintain a Kendall tau distance below that of insertion sort at certain time periods, but do not do so at all times and therefore do not reach a lower average maintained Kendall tau distance over the endpoint.

Lastly, we see that, in terms of the average Kendall tau distance maintained over the endpoint  $k = 1$  and  $k = 2$  seem to fall right in between repeated insertion sort and repeated quicksort, so around the mean of these two algorithms. Additionally, we see that, if we exclude  $\alpha = 250$ , this holds more closely for higher  $\alpha$ , most likely due to the narrowing of the gap between the performance of repeated insertion sort and repeated quicksort.

From these results, we conclude the following; Algorithm 7 can be viable in cases where the  $\alpha$  is unknown or changes continuously. This is because it performs better than repeated quicksort for low  $\alpha$ , while performing better than repeated insertion sort for high  $\alpha$ . Additionally, performing more insertion sort runs, so higher values for  $k$ , is beneficial in the cases when repeated insertion sort is the best performing algorithm, while it is disadvantageous when it is the worst performing algorithm. So, depending on how quickly  $\alpha$  changes or on how high  $\alpha$  is, we can select an ideal  $k$  to maintain a lower average Kendall tau distance over the endpoint.

## 5 Conclusions

In this thesis, we performed an experimental study regarding sorting in the evolving data model. To this end, we introduced the setting and algorithms before we perform the experiments.

Our initial experiment shows that all algorithms converge within  $O(n^2)$  when starting from a maximal Kendall tau distance, and that the Kendall tau distances found in this case are similar to those found by Besa Vial *et al.*[3], who started from a randomly shuffled list.

In our second experiment we validate the theoretical bounds by seeing if they hold in practice. We find that for  $\alpha \in \{1, 2, 10\}$  the theoretical bounds clearly hold in practice for most algorithms, and that it is highly likely that the theoretical bound holds in practice for blocksort.

In our third experiment, we explore the effect of the block size on the maintained Kendall tau distance. Our results show that the block size does have an effect on the maintained Kendall tau distance and that this effect increases for higher values of  $n$ . We find that smaller block constants seem to perform better in general, but that if the resulting block constant becomes too small to account for the changes it performs considerably worse. So, there seems to be a trade-off between performance and robustness, where a lower block constant performs better for lower change rates but is not robust for higher change rates. Lastly, we find that smaller block sizes do not necessarily result in a lower maintained Kendall tau distance in the steady behavior.

In our last experiment, we looked at the performance and behavior of our new introduced algorithm, Repeated Quicksort Followed by Repeated Insertion Sort, compared to repeated quicksort and repeated insertion sort. Here we find that our new variant exhibits the expected behavior of both quicksort and insertion sort, and falls between the two in terms of maintained Kendall tau distance in the steady behavior. Furthermore, we find that this algorithm can be viable in cases where the  $\alpha$  is unknown or changes continuously. This is because it performs better than repeated quicksort for low  $\alpha$ , while performing better than repeated insertion sort for high  $\alpha$ . Additionally, performing more insertion sort runs, so higher values for  $k$ , is beneficial in the cases when repeated insertion sort is the best performing algorithm, while it is disadvantageous when it is the worst performing algorithm. So, depending on how quickly  $\alpha$  changes or on how high  $\alpha$  is, we can select an ideal  $k$  to maintain a lower average Kendall tau distance over the endpoint.

Given our results, it would be interesting to explore the exact nature of this effect of the block size on the maintained Kendall tau distance. It would be interesting to answer the question on why exactly this occurs and what an ideal block size would be given the input size and the change rate. Additionally, it would be interesting to see how the current block constants perform for larger input sizes to see whether they perform better, worse, or the same as they did.

Additionally, it would be interesting to further explore our variant of insertion sort. Our results seem to imply that after a certain value for  $\alpha$  repeated quicksort performs better than repeated insertion sort, which is the turning point for our algorithm. It would be interesting to try the same experiment for higher values of  $\alpha$ , where  $20 < \alpha < 250$ , to see when this turning point is and if it is reasonable to expect it to occur. Lastly, it would be interesting to see how the algorithm performs in practice for higher rates of  $k$ .

## References

- [1] Aris Anagnostopoulos et al. “Sorting and selection on dynamic data”. In: *Theoretical Computer Science* 412.24 (2011). Selected Papers from 36th International Colloquium on Automata, Languages and Programming (ICALP 2009), pp. 2564–2576. ISSN: 0304-3975. DOI: <https://doi.org/10.1016/j.tcs.2010.10.003>. URL: <https://www.sciencedirect.com/science/article/pii/S0304397510005475>.
- [2] Juan José Besa Vial et al. “Optimally Sorting Evolving Data”. In: *CoRR* abs/1805.03350 (2018). arXiv: [1805.03350](https://arxiv.org/abs/1805.03350). URL: <http://arxiv.org/abs/1805.03350>.
- [3] Juan José Besa Vial et al. “Quadratic Time Algorithms Appear to be Optimal for Sorting Evolving Data”. In: *CoRR* abs/1805.05443 (2018). arXiv: [1805.05443](https://arxiv.org/abs/1805.05443). URL: <http://arxiv.org/abs/1805.05443>.

## A Model Implementation

We start with the implementation of our `DynamicList` class and the relevant `Stats` class. We give an instance of the `DynamicList` class to an sorting algorithm at the start of a run. The true order and our approximation are stored in the `DynamicList`. Furthermore, probing within the sorting algorithm is done through the instance with the *probe\_with\_swap* method. Within this method we implement the random swaps that occur in the true ranking by means of the *random\_swap* method. Additionally, within this method we store the current Kendall tau distance every so often depending on the time, which is the amount of probes performed, and the sample rate. To store the current distance we call the *add\_curr\_distance* with the real and approx lists as arguments.

```
1 from model.statistics import Stats
2 import random
3
4
5 class DynamicList:
6     real: list[int]
7     curr_approx: list[int]
8     change_rate: int
9     stats: Stats
10    start_time: int
11
12    # Is currently unused but it seems to be too late to remove it without
13    # potentially breaking something, can be used in the future but would
14    # most likely require
15    # you to change the way the probe method works.
16    probe_rate: int
17
18    def __init__(
19        self,
20        rand_seed: int,
21        probe_rate: int,
22        change_rate: int,
23        n: int,
24        sample_rate: int,
25        time_limit: int,
26        start_time: int
27    ):
28        random.seed(rand_seed)
29        self.stats = Stats()
30        self.probe_rate = probe_rate
31        self.change_rate = change_rate
32        self.sample_rate = sample_rate
33        self.time_limit = time_limit
34        self.start_time = start_time
35
36        # We start with a real that goes from 0, to n
37        self.real = []
38        for i in range(0, n):
39            self.real.append(i)
40
41        # We start with an approximation that goes from 0, to n
42        self.curr_approx = []
43        for i in range(0, n):
```

```

43         self.curr_approx.append(i)
44
45     # Performs change_rate amount of uniformly random swaps before probe,
46     # sorts based on the index of int i and j in the real list, so based on
47     # their positions in the real
48     def probe_with_swap(self, i: int, j: int):
49         self.stats.add_probe(i, j)
50
51         if (self.get_time() % self.sample_rate == 0) and (self.get_time()
52 >= self.start_time) and (self.get_time() <= self.time_limit):
53             self.stats.add_curr_distance(self.real, self.curr_approx)
54
55             self.random_swap(self.change_rate)
56
57             index_i = self.real.index(i)
58             index_j = self.real.index(j)
59             return index_i < index_j
60
61     # Used to perform random swaps in the real
62     def random_swap(self, change_rate):
63         n = len(self.real)
64
65         # Perform change rate amount of random swaps in the real, swapping
66         # i with i + 1
67         for x in range(0, change_rate):
68             i = random.randint(0, n - 2)
69             temp = self.real[i]
70             self.real[i] = self.real[i + 1]
71             self.real[i + 1] = temp
72
73     # Used to update the current approximation
74     def permute_answer(self, ans_perm: list[int]):
75         self.curr_approx = ans_perm.copy()
76
77     # Used to get the size of the dynamic list
78     def size(self):
79         return len(self.real)
80
81     # Used to get the current number of probes performed,
82     # we only consider the case where the probe rate = 1,
83     # so this is the same as taking the len of the list keeping track of
84     # all the probes
85     def get_time(self):
86         return len(self.stats.probes)
87
88     # Used to start with a real list that is reverse sorted,
89     # so it is sorted decreasingly
90     def reverse_order(self):
91         n = self.size()
92         for i in range(0, n):
93             self.real[i] = (n - 1) - i

```

In the Stats class we store the probes and the calculated Kendall tau distances. Note that we calculate the distance between two lists with the *merge\_sort* method, which we call from within the *add\_curr\_distance* method. The *merge\_sort* method has access to the real list by means of the *probe*, which does not apply any changes in the true order, which allows merge sort to make as many probes as it requires to count all the inversions.

The probes are stored in the probes list, while the distances get stored in the distances list. When we call *get\_time* in the DynamicList we return the length of the probes list because we only consider the case where the probe rate is equal to one. After a full execution, we store the distances list as an excel sheet as the result of that specific run.

```

1 # Used for tracking the time by means of storing the probes,
2 # and for tracking the Kendall tau distances at certain time steps,
3 # which depends on the sampling rate
4 class Stats:
5     probes: list[list[int]]
6     distances: list[int]
7
8     def __init__(self):
9         self.distances = []
10        self.probes = []
11
12    # Adds the current probe and the probes indexes to the list
13    def add_probe(self, i, j):
14        self.probes.append([i, j])
15
16    # Allows comparison based on the index of int i and j in the real list,
17    # so based on their positions in the real list,
18    # no changes are made to the real list
19    def probe(self, real: dict, i: int, j: int):
20        index_i = real.get(i)
21        index_j = real.get(j)
22        return index_i < index_j
23
24    # Merge sort that counts the number of inversions in the list,
25    # which is the number of discordant pairs in the current approximation,
26    # to this end we give access to the real list as a dictionary,
27    # where the keys are the integers and the values are their indices in
28    # the real list
29    def merge_sort(self, real: dict, temp: list[int]):
30        if len(temp) <= 1:
31            return 0
32
33        # Splitting the list into the left and right halves
34        # and recursively sorting them
35        left_size = len(temp) // 2
36        right_size = len(temp) - left_size
37
38        left = []
39        for i in range(0, left_size):
40            left.append(temp[i])
41
42        right = []
43        for i in range(left_size, len(temp)):
44            right.append(temp[i])

```



```

45
46     left_invs = self.merge_sort(real, left)
47     right_invs = self.merge_sort(real, right)
48     between_invs = 0
49     i = 0
50     j = 0
51     k = 0
52     while i < left_size and j < right_size:
53
54         if self.probe(real, left[i], right[j]):
55             temp[k] = left[i]
56             k += 1
57             i += 1
58         else:
59             temp[k] = right[j]
60
61             # Something on the right being smaller means that there are
62             # of inversions, namely (left_size - i), where i is the
63             # index of the current element on the left,
64             # so we essentially have that this right side element
65             # agrees with i elements but is discordant with the rest of the left
66             # array
67
68             between_invs += (left_size - i)
69             k += 1
70             j += 1
71
72     # Copy over the remaining elements in either left or right,
73     # the while loop will fail for at least one
74     while i < left_size:
75         temp[k] = left[i]
76         k += 1
77         i += 1
78
79     while j < right_size:
80         temp[k] = right[j]
81         k += 1
82         j += 1
83
84     return left_invs + right_invs + between_invs
85
86 # Adds the current distance between the real order and our
87 # approximation
88 def add_curr_distance(self, real: list[int], approx: list[int]):
89     temp = approx.copy()
90     real_dict = dict(zip(real, range(0, len(real))))
91
92     distance = self.merge_sort(real_dict, temp)
93     self.distances.append(distance)

```

## B Algorithms Implementation

We implemented each algorithm according to the pseudocode, where we implemented the comparisons using *probe\_with\_swap*. Note how insertion sort makes changes directly to the approximation list, while the other sorting algorithms only update the answer after "termination", so they exchange the whole current approximation for a new one, while insertion sort gradually updates it. Lastly, we want to mention that we looked at the implementation of blocksort of Besa Vial *et al.*[3] and implemented it in Python, in order to have an accurate comparison of algorithm behavior. For our stack implementation of blocksort we use the QSortState in order to fully store a sort state during blocksort.

```
1 from model.dynamic_list import DynamicList
2 from model.q_sort_state import QSortState
3 import random
4 import math
5
6
7 def partition(list: DynamicList, to_sort, low: int, high: int):
8     ranges = high - low + 1
9
10    # We randomly pick a pivot, probably should have done random.randint(
11    low, high)
12    pivotChoice = low + random.randint(0, 32767) % ranges
13
14    temp = to_sort[high]
15    to_sort[high] = to_sort[pivotChoice]
16    to_sort[pivotChoice] = temp
17
18    i = low - 1
19    for j in range(low, high):
20        # Extra check that makes sure that we don't go over time
21        if list.get_time() >= list.time_limit:
22            return -1
23
24        if list.probe_with_swap(to_sort[j], to_sort[high]):
25            i += 1
26            temp = to_sort[i]
27            to_sort[i] = to_sort[j]
28            to_sort[j] = temp
29
30    temp = to_sort[i + 1]
31    to_sort[i + 1] = to_sort[high]
32    to_sort[high] = temp
33    return i + 1
34
35 def quicksort(list: DynamicList, to_sort, low: int, high: int):
36     if low < high:
37         pivot = partition(list, to_sort, low, high)
38
39         # Pivot of -1 means that list.get_time() >= list.time_limit so we
40         # return,
41         # introduced in order to stop running the algorithm after the time
42         # limit has passed
43         if pivot == -1:
44             return
```

```

42         quicksort(list, to_sort, low, pivot - 1)
43         quicksort(list, to_sort, pivot + 1, high)
44
45
46
47 # Given a random list that goes from 0 to n-1 sort it based on the order in
48 # the real list,
49 # afterwards update the curr_approx to the result,
50 # this is the initial quicksort called when a run of quicksort is performed
51 def quicksort_base(list: DynamicList, n: int):
52     to_sort = []
53     for i in range(0, n):
54         to_sort.append(i)
55
56     quicksort(list, to_sort, 0, n - 1)
57     list.permute_answer(to_sort)
58
59 def repeated_insertion_sort(list: DynamicList, time_limit: int):
60     n = list.size()
61     while list.get_time() < time_limit:
62         for i in range(1, n):
63             j = i
64             # Probe the ints at approx[j] with approx[j - 1]
65             while j > 0 and list.probe_with_swap(
66                 list.curr_approx[j], list.curr_approx[j - 1]
67             ):
68                 temp = list.curr_approx[j - 1]
69                 list.curr_approx[j - 1] = list.curr_approx[j]
70                 list.curr_approx[j] = temp
71                 j -= 1
72
73             # Extra check that makes sure that we don't go over time
74             if list.get_time() >= time_limit:
75                 return
76
77
78 # Perform one run of quicksort, followed by repeated insertion sort
79 def quick_then_insertion_sort(list: DynamicList, time_limit: int):
80     n = list.size()
81     quicksort_base(list, n)
82     repeated_insertion_sort(list, time_limit)
83
84
85 # Perform repeated runs of one quicksort, followed by a certain amount of
86 # insertion sort runs
87 def rep_quick_then_insertion_sort(list: DynamicList, time_limit: int,
88     iterations: int):
89     n = list.size()
90
91     while list.get_time() < time_limit:
92         quicksort_base(list, n)
93
94         k = 0
95         while (list.get_time() < time_limit) and (k < iterations):
96             for i in range(1, n):

```

```

95         j = i
96         # Probe the ints at approx[j] with approx[j - 1]
97         while j > 0 and list.probe_with_swap(
98             list.curr_approx[j], list.curr_approx[j - 1]
99         ):
100             temp = list.curr_approx[j - 1]
101             list.curr_approx[j - 1] = list.curr_approx[j]
102             list.curr_approx[j] = temp
103             j -= 1
104
105         # Extra check that makes sure that we don't go over
106         time
107         if list.get_time() >= time_limit:
108             return
109
110         k += 1
111
112     # Repeatedly perform quicksort until the time limit is reached
113     def repeated_quicksort(list: DynamicList, time_limit: int):
114         n = list.size()
115         while list.get_time() < time_limit:
116             quicksort_base(list, n)
117
118     def start_new_quicksort_call(
119         to_sort: list[int], call_stack: list[QSortState], low: int, high: int
120     ):
121         if low < high:
122             range = high - low + 1
123             pivotChoice = low + random.randint(0, 32767) % range
124             temp = to_sort[high]
125             to_sort[high] = to_sort[pivotChoice]
126             to_sort[pivotChoice] = temp
127             newCall = QSortState(low, high)
128             call_stack.append(newCall)
129
130     def stack_quicksort_run_step(
131         list: DynamicList, to_sort: list[int], call_stack: list[QSortState]
132     ):
133         step_performed = False
134
135         while not step_performed:
136
137             # If there are no quicksort calls on the stack, terminate
138             if len(call_stack) == 0:
139                 return False
140
141             # Otherwise attempt to run a step of the top call
142             curr_call = call_stack[-1]
143             low = curr_call.low
144             high = curr_call.high
145             i = curr_call.i
146             j = curr_call.j
147
148             if j < high:

```

```

150         if list.probe_with_swap(to_sort[j], to_sort[high]):
151             curr_call.i += 1
152             temp = to_sort[i + 1]
153             to_sort[i + 1] = to_sort[j]
154             to_sort[j] = temp
155             curr_call.j += 1
156             step_performed = True
157         else:
158             call_stack.pop()
159             # If the quicksort call finished,
160             # then move the pivot into the correct position
161             temp = to_sort[i + 1]
162             to_sort[i + 1] = to_sort[high]
163             to_sort[high] = temp
164
165             # And start the two new recursive calls
166             low_left = low
167             high_left = i
168             low_right = i + 2
169             high_right = high
170             start_new_quicksort_call(to_sort, call_stack, low_left,
high_left)
171             start_new_quicksort_call(to_sort, call_stack, low_right,
high_right)
172         return True
173
174
175 def blocked_quicksort(list: DynamicList, time_limit: int, block_constant:
int):
176     n = list.size()
177
178     # Determine block size m close to block_constant * ln(n)
179     m = int(block_constant * math.log(n))
180     if m % 2 == 1:
181         m += 1
182     if m > n:
183         m = n
184     while n % m != 0:
185         m += 2
186
187     # Initial full quicksort
188     full_stack = []
189     quicksort_base(list, n)
190     to_full_sort = list.curr_approx.copy()
191
192     last_full_answer = list.curr_approx.copy()
193     next_full_answer = list.curr_approx.copy()
194     new_full_answer = False
195
196     start_new_quicksort_call(to_full_sort, full_stack, 0, n - 1)
197     block_stack = []
198
199     while list.get_time() < time_limit:
200         to_blocksort = [None] * m
201         block_answer = [None] * n
202

```

```

203     # Switch to new full answer if available
204     if new_full_answer:
205         new_full_answer = False
206         last_full_answer = next_full_answer.copy()
207
208     # Fill first m/2 items
209     for i in range(m // 2):
210         temp = last_full_answer[i]
211         to_blocksort[i] = temp
212
213     for i in range(1, (2 * n) // m):
214         # Fill next m/2 items
215         for j in range(m // 2):
216             temp = last_full_answer[i * (m // 2) + j]
217             to_blocksort[m // 2 + j] = temp
218
219     # Sort current block
220     start_new_quicksort_call(to_blocksort, block_stack, 0, m - 1)
221     while stack_quicksort_run_step(list, to_blocksort, block_stack)
222 :
223         full_sort_done = not stack_quicksort_run_step(
224             list, to_full_sort, full_stack
225         )
226         if full_sort_done:
227             new_full_answer = True
228             next_full_answer = to_full_sort.copy()
229             start_new_quicksort_call(to_full_sort, full_stack, 0, n
230 - 1)
231
232     # Copy smallest half to block answer, shift larger half forward
233     for j in range(m // 2):
234         temp1 = to_blocksort[j]
235         block_answer[(i - 1) * (m // 2) + j] = temp1
236         temp2 = to_blocksort[m // 2 + j]
237         to_blocksort[j] = temp2
238
239     # Handle last m/2 elements
240     for j in range(m // 2):
241         temp = to_blocksort[j]
242         block_answer[n - (m // 2) + j] = temp
243
244     list.permute_answer(block_answer)
245
246     # Update all working answers
247     last_full_answer = block_answer.copy()
248     next_full_answer = block_answer.copy()
249     to_full_sort = block_answer.copy()
250
251 # Used in order keep track of the sort states during Block Sort
252 class QSortState:
253     low: int
254     high: int
255     i: int
256     j: int
257
258     def __init__(self, l, h):
259         self.low = l

```

```

10         self.high = h
11         self.i = self.low - 1
12         self.j = self.low

```

## C Running Experiments

We use the Runner class to run an algorithm in a specific configuration for different seeds. We created instances of the Runner class within a nested for loop, where the first for loop went over the input size, the second over the algorithm, the third over the change rate and the last over the seeds. This allowed us to construct a whole experiment by specifying the lists and running the nested for loop. Most important is to note the possible starting ordering for the real, which are "reverse-sorted" and "sorted". It is also important to initialize the Runner instance with a valid algorithm name, which can be found in the *run* method.

```

1 from model.dynamic_list import DynamicList
2 from algorithms.algorithms import *
3 import pandas as pd
4
5
6 class Runner:
7     alg: str
8     curr_list: DynamicList
9     time_limit: int
10
11     def __init__(
12         self,
13         rand_seed: int,
14         probe_rate: int,
15         change_rate: int,
16         algorithm: str,
17         input_size: int,
18         time_limit: int,
19         sample_rate: int,
20         config: str,
21         start_time: int = 0
22     ):
23         self.alg = algorithm
24         self.time_limit = time_limit
25         self.curr_list = DynamicList(
26             rand_seed, probe_rate, change_rate, input_size, sample_rate,
27             time_limit, start_time
28         )
29
30         # Currently only valid options are sorted and reverse-sorted
31         if config == "reverse-sorted":
32             self.curr_list.reverse_order()
33         elif config != "sorted":
34             raise Exception("Invalid config")
35
36         # Runs the relevant algorithm with the current configuration
37         def run(self):
38             if self.alg == "rep-quick":
39                 repeated_quicksort(self.curr_list, self.time_limit)

```

```

39         elif self.alg.startswith("block"):
40             temp = self.alg.split("-")
41             block_constant = temp[len(temp) - 1]
42             blocked_quicksort(self.curr_list, self.time_limit, int(
block_constant))
43         elif self.alg == "rep-insertion":
44             repeated_insertion_sort(self.curr_list, self.time_limit)
45         elif self.alg == "quick-rep-insertion":
46             quick_then_insertion_sort(self.curr_list, self.time_limit)
47         elif self.alg.startswith("rep-quick-rep-insertion"):
48             temp = self.alg.split("-")
49             i = temp[len(temp) - 1]
50             rep_quick_then_insertion_sort(self.curr_list, self.time_limit,
int(i))
51         else:
52             raise Exception("No correct alg given as input")
53
54     # Saves the results of the current run in the for of an excel sheet
55     def store_results(self, file_name):
56         results = pd.DataFrame(self.curr_list.stats.distances)
57         datatoexcel = pd.ExcelWriter(path=file_name)
58         results.to_excel(datatoexcel)
59         datatoexcel.close()

```