# Project 1

## Integration of Merge Sort & Insertion Sort

### Team 4

1. He Qi Xin
2. Dinh Pham Minh Anh
3. Hong Jia Yang

Preamble: Overhead in Merge Sort

$$MergeSort(n) = \Theta(n \log(n))$$

$$InsertionSort(n) = \Theta(n^2)$$

$$\exists k \in \mathbb{N} : c_1 k^2 \leq c_2 k \log(k), 0 < c_1 < c_2$$

When size of sublists is small enough,
insertion sort should be performed instead

# (a) Algorithm implementation: Implement the above hybrid algorithm

## Class Structure



**My Project**

Main Page    Classes ▾    Files ▾

### MergeSort Struct Reference

#### Public Member Functions

| | |
|---:|:---|
| | **MergeSort** (vector< int > &IN) |
| void | **sort** (int p, int r) |
| void | **insertionSort** (int p, int r) |
| void | **merge** (int p, int q, int r) |
| void | **info** () |
| void | **setK** (short int k) |
| short | **getK** () |
| const auto | **getKeyCmp** () const |
| const auto | **getKeyCmp_insertionSort** () const |
| vector< int > & | **getArray** () |
| const auto | **getSize** () const |
| void | **unsort** () |

The documentation for this struct was generated from the following file:

- **project1.hpp**

Main Program

```
1       void sort(int p, int r) {
2               if (r - p + 1 <= k) {
3                       insertionSort(p, r);
4                       return;
5               }
6
7               if (p < r) {
8                       int q  = (p+r) / 2;
9                       sort(p, q);
10                      sort(q+1, r);
11                      merge(p, q, r);
12              }
13      }
```

$$A = \langle a_0, a_1, \ldots, a_{n-2}, a_{n-1} \rangle$$

$$size[A] = (n-1) - 0 + 1 = n$$

$$A' = \langle a_p, a_{p+1}, \ldots, a_{r-1}, a_r \rangle$$

$$size[A'] = r - p + 1$$

$$A' \in A$$

$$sort(n) = \begin{cases} insertionSort(n) & \text{if } n \leq k \\ mergeSort(n) & \text{otherwise} \end{cases}$$

Insertion Sort (Trivial)

```
1       void insertionSort(int p, int r) {
2             for (int i = p+1; i <= r; i++) {
3                   for (int j = i ; j > p; j--) {
4                         key_cmp_insertionSort++;
5                         key_cmp++;
6                         if (A[j] < A[j-1]) {
7                               int t = A[j];
8                               A[j] = A[j-1];
9                               A[j-1] = t;
10                        }
11                        else break;
12                  }
13            }
14      }
```

Merge (Trivial)

```cpp
1     void merge(int p, int q, int r) {
2             vector<int> L, R;
3             for (int i = p; i <= q; i++) L.push_back(A[i]);
4             for (int i = q+1; i <= r; i++) R.push_back(A[i]);
5
6             int idx_L{}, idx_R{}, idx_A{ p };
7             while (idx_L < L.size() && idx_R < R.size()) {
8                     if (L[idx_L] < R[idx_R]) {
9                             A[idx_A] = L[idx_L];
10                            idx_L++;
11                    }
12                    else {
13                            A[idx_A] = R[idx_R];
14                            idx_R++;
15                    }
16                    idx_A++;
17                    key_cmp++;
18            }
......
```

Continues Next Page

## Simple Python Script

```python
1    file_index = 1
2    i = MIN_SIZE
3    while (i <= MAX_SIZE):
4        try:
5            file = open('input_%i.txt' % file_index, 'w')
6        except OSError:
7            print("[-] Error in writing file 'input_%i.txt'" % file_index)
8            sys.exit()
9
10       A = [str(random.randint(1, MAX_INT)) for k in range(1, i+1)]
11       file.write(("\n").join(A))
12       file.write("\n")
13       file_index += 1
14       i *= SCALE_FACOTR
15       file.close()
16       msg = "[+] input_{0}.txt generated. number of elements is {1}."
17       print(msg.format(file_index-1, i))
```

Array length grow by a factor of 10, for a total of 5 input files

Visualisation of Distribution

Growth of Key Comparisons against Array Size

Rather "Linear" for
the given input sizes

Time Complexity of Insertion Sort

$$\text{length of sublists} = S = k$$

$$\text{number of sublists} = \frac{n}{k}$$

$$InsertionSort() = \Theta(k^2 \cdot \frac{n}{k}) = \Theta(nk)$$

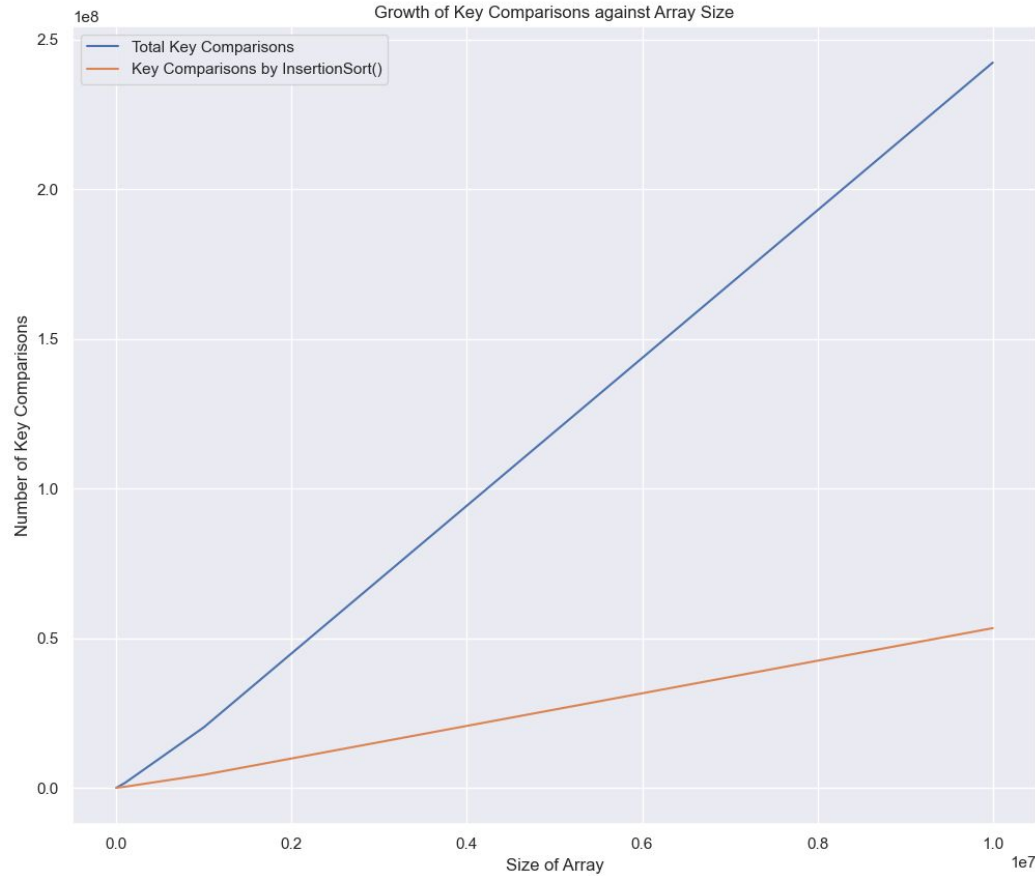## Time Complexity of Merge Sort



$cn$

$n$

$2\left(\dfrac{cn}{2}\right) = cn$

n/2          ....          n/2

$2^h = \dfrac{n}{k}$

height of tree $= \log(\dfrac{n}{k})$

$\dfrac{cn}{k}(k) = cn$

k          k          k          k          k

....

1          1          1          1          1

$$MergeSort() = \Theta\left(n\log(\dfrac{n}{k})\right)$$

Growth of Key Comparisons against Array Size

$$HybridSort(n, k) = \Theta \left( kn + n \log \frac{n}{k} \right)$$

Time Complexity of Hybrid Sort

By definition, $\quad HybridSort(n, k) = \Theta\left(kn + n\log\frac{n}{k}\right)$

if and only if

$$\exists n_0 \exists c_1 \exists c_2 \in \mathbb{N}: \quad \forall n > n_0,$$

$$0 < c_1\left(kn + n\log\frac{n}{k}\right) \leq HybridSort(n, k) \leq c_2\left(kn + n\log\frac{n}{k}\right)$$

Let c_1 = 1/2, c_2 = 2. Observe that the condition is satisfied



Growth of Key Comparisons against Array Size

Verified (RHS)

$$HybridSort(n, k) = \Theta\left(kn + n\log\frac{n}{k}\right)$$

$$k(n) \in O(\log n)$$

$$\Theta(kn + n\log\frac{n}{k}) = \Theta(kn + n\log n - n\log k)$$
$$= \Theta(n\log n + n\log n - n\log(\log n)$$
$$= O(n\log n)$$

Hybrid sort will be as slow as merge sort if the value of k = S approaches asymptotic value of O(logn)

Growth of Key Comparisons against S for Array Size of 1 million

$$HybridSort(n, k) = \Theta\left(kn + n\log\frac{n}{k}\right)$$

Key Comparisons give rough estimate of runtime

Growth of Key Comparisons against S for Array Size of 1 million

Periods of **Constant**
Periods of **Spikes**

```
1      if (r - p + 1 <= k) {
2              insertionSort(p, r);
3              return;
4      }
```

Periods of **Constant**



Growth of Key Comparisons against S for Array Size of 1 million

| | (num_sublist, ASize) | Total Key Comparisons |
|---|---|---|
| S = 600 | (1666, 600) (1, 400) | 133710976 |
| S = 800 | (1250, 800) | 133710976 |

```
1    if (r - p + 1 <= k) {
2            insertionSort(p, r);
3            return;
4    }
```

Periods of **Spikes**



Growth of Key Comparisons against S for Array Size of 1 million

| | (num_sublist, ASize) | Total Key Comparisons |
|---|---|---|
| S = 1000 | (1000, 1000) | 244639044 |
| S = 800 | (1250, 800) | 133710976 |

Growth of Key Comparisons against S for Array Size of 1 million

This shape is **non-trivial**

*Stay tuned*

$$T(n, k) = c_1 nk + c_2 \log \frac{n}{k} + c_3$$

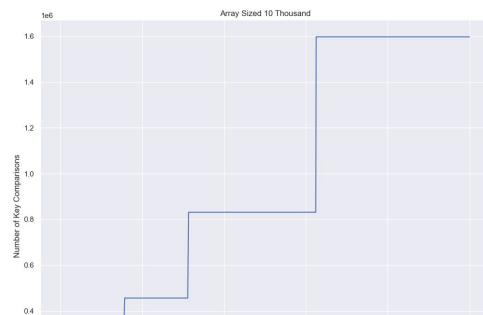Want to solve for $\quad T'(k_s) = 0$

For each n, such that $\quad T''(k_s) > 0$

In practice, it is **not feasible** to do this.
Instead, we will **run tests** and give an **estimate**

**Intuition**

Find the *largest S* for which key comparison growth is *insubstantial*

Growth of Key Comparisons against S for Various Array Sizes

Need to adjust down range of S

Growth of Key Comparisons against S for Array Size of 1 thousand

S = 61
min(S) = 61

Growth of Key Comparisons against S for Array Size of 10 thousand

S = 77

S = 77
min(S) = 61

Growth of Key Comparisons against S for Array Size of 100 thousand

S = 47
min(S) = 47

Growth of Key Comparisons against S for Array Size of 1 million

S = 60
min(S) = 47

Growth of Key Comparisons against S for Array Size of 10 million

S = 37
min(S) = 37

Growth of Key Comparisons against S for Array Size of 10 million

Why not choose **earlier values**?
Is **key comparison** the best gauge?

S = 37
min(S) = 37

Growth of Runtime(ms) against S for Various Array Sizes

**Runtime** gives the best representation

Growth of Runtime(ms) against S for Various Array Sizes

Runtime gives the best representation

Growth of Runtime(ms) against S for Various Array Sizes

S must be
sufficiently large
before time saving

Growth of Runtime (ms) against S for Array Sized 1 Thousand

Growth of Runtime (ms) against S for Array Sized 10 Thousand

Growth of Runtime (ms) against S for Array Sized 100 Thousand

Growth of Runtime (ms) against S for Array Sized 10 Million

Growth of Runtime (ms) against S for Array Sized 1 Million

Growth of Runtime (ms) against S for Array Sized 1 Thousand

Growth of Runtime (ms) against S for Array Sized 10 Thousand

Growth of Runtime (ms) against S for Array Sized 100 Thousand

Find the minimum of the set of **local minimums** of each of the graph and we get **S = 60**

| | ArraySize | S | TotalKeyCmp | KeyCmpInsertionSort | Runtime |
|---|---|---|---|---|---|
| 0 | 1000 | 37 | 13395 | 8456 | 0.3275 |
| 1 | 1000 | 60 | 13395 | 8456 | 0.3312 |
| 2 | 10000 | 37 | 143844 | 54826 | 3.9552 |
| 3 | 10000 | 60 | 183142 | 103648 | 3.7200 |
| 4 | 100000 | 37 | 1862508 | 670504 | 42.0865 |
| 5 | 100000 | 60 | 2380698 | 1284737 | 42.9996 |
| 6 | 1000000 | 37 | 23201902 | 8266274 | 486.0120 |
| 7 | 1000000 | 60 | 23201902 | 8266274 | 487.0390 |
| 8 | 10000000 | 37 | 242332523 | 53347890 | 5981.0200 |
| 9 | 10000000 | 60 | 281253476 | 101769647 | 5937.8500 |

**60** has similar runtime

Set **k = S = 0** gives merge sort

| | ArraySize | S | TotalKeyCmp | KeyCmpInsertionSort | Runtime |
|---|---|---|---|---|---|
| **0** | 10000000 | 60 | 281253476 | 101769647 | 5851.06 |
| **1** | 10000000 | 0 | 220097824 | 0 | 11886.70 |

```
1    if (r - p + 1 <= k) {
2        insertionSort(p, r);
3        return;
4    }
```

Despite having about **28%** more key comparisons,

the hybrid algorithm is about **101% faster**

Set **k = S = 0** gives merge sort

| | ArraySize | S | TotalKeyCmp | KeyCmpInsertionSort | Runtime |
|---|---|---|---|---|---|
| **0** | 10000000 | 60 | 281253476 | 101769647 | 5851.06 |
| **1** | 10000000 | 0 | 220097824 | 0 | 11886.70 |

```
1    if (r - p + 1 <= k) {
2        insertionSort(p, r);
3        return;
4    }
```

Number of key comparisons is **only an estimate**
The runtime depends on the various **overheads**

Growth of Key Comparisons against S for Array Size of 1 million

**Back to this..**

Growth of Key Comparisons against S for Array Size of 1 million

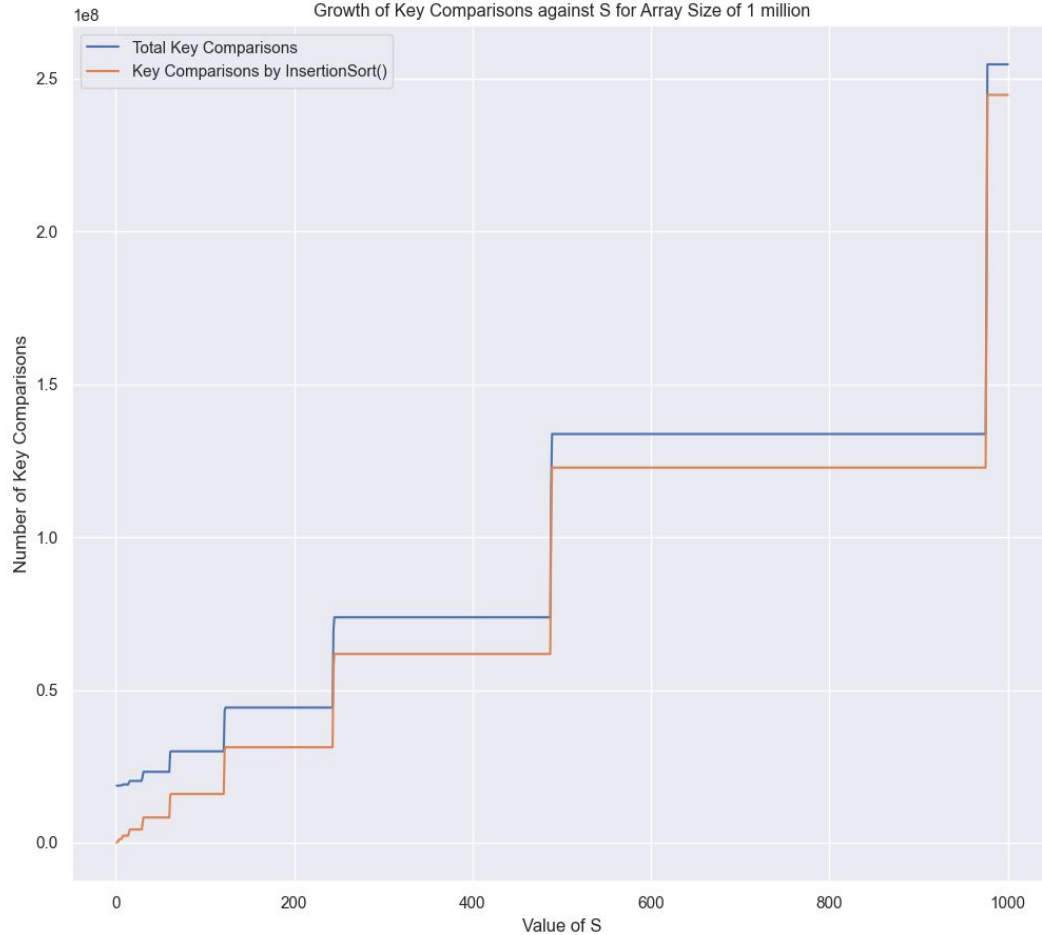**Intuition 1:**

Before a certain threshold, **S is too small** such that runtime is slower

Key comparisons also happened to be lowest for small S

Growth of Key Comparisons against S for Array Size of 1 million

Intuition 2:

Square form

Growth of Key Comparisons against S for Array Size of 1 million

What happens if we adjust the scale?

$$HybridSort(n, k) = \Theta\left(kn + n\log\frac{n}{k}\right)$$

$$k(n) \in O(\log n)$$

- When **S ≤ 1**, the key comparisons are **completely done by Merge Sort** *(slower than hybrid)*

- When **S == Size of Array**, the key comparisons are completely done by **Insertion Sort**

- The number of key comparisons made by pure merge sort will be the **minimum** number of comparisons **(local min)**

- The number of key comparisons made by pure insertion sort will be the **maximum** number of comparisons **(local max)**

$$\Theta(kn + n\log\frac{n}{k}) = \Theta(kn + n\log n - n\log k)$$
$$= \Theta(n\log n + n\log n - n\log(\log n))$$
$$= O(n\log n)$$

- When  **S ≤ threshold,** most of the sorting is still done by the slower merge sort, resulting **slower run time despite lower key comparison**

- As S increases, **more** key comparisons is done using **insertion sort**, and **less** key comparisons is needed from **merge sort** because the increasingly longer sublists are already sorted

- This leads to **increase in total key comparisons** overtime because growth in key comparisons in insertion sort is O($n^2/2$)

- The way we write the **code for key comparison** is such that comparison **terminates when either array is empty,** and this explains the period of **constants and peaks**

- Additionally, previous graphs used a **linear scale** for S, which makes it more difficult to see a clear relationship