



# The Case for a Learned Sorting Algorithm

Authors: Ani Kristo, Kapil Vaidya, Ugur Cetintemel, Sanchit Misra, Tim Kraska

Presenter: Terryn Brunelle



# Motivation

- Fundamental CS problem
- Database operations
  - Sort query results
  - Perform joins



## Existing Work

- Comparison sort
- Distribution sort
  - Counting sort
  - Radix sort
- ML-enhanced algorithms

# Learned Sort

- Train CDF model
- Use predicted prob for each key to predict final position for every key in sorted output

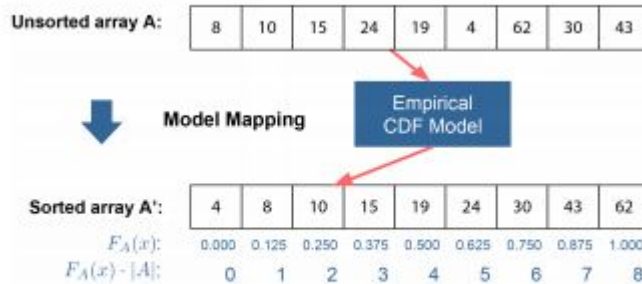


Figure 1: Sorting with the perfect CDF model

Linear time possible!  
(if have perfect model)



# Problems

- Perfect model = expensive to train
- Random-access memory problem

# Algorithm 1

---

## Algorithm 1 A first Learned Sort

---

**Input**  $A$  - the array to be sorted

**Input**  $F_A$  - the CDF model for the distribution of  $A$

**Input**  $o$  - the over-allocation rate. Default=1

**Output**  $A'$  - the sorted version of array  $A$

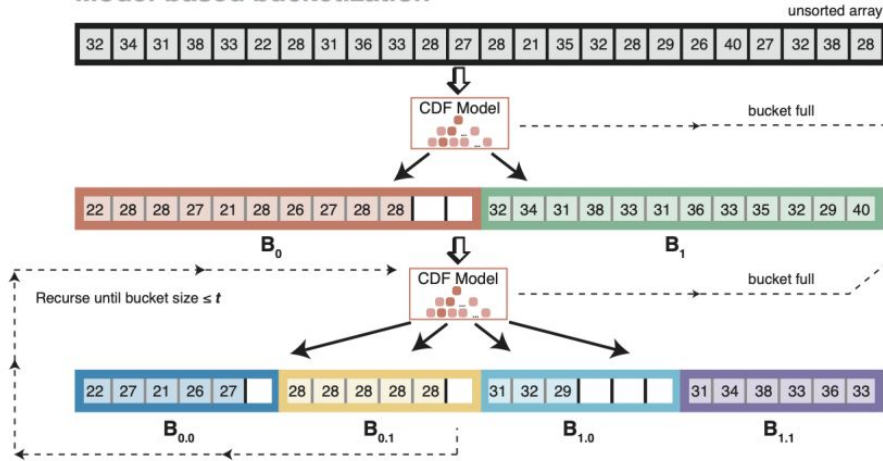
```
1: procedure LEARNED-SORT( $A, F_A, o$ )
2:    $N \leftarrow A.length$ 
3:    $A' \leftarrow$  empty array of size  $(N \cdot o)$ 
4:   for  $x$  in  $A$  do
5:      $pos \leftarrow \lfloor F_A(x) \cdot N \cdot o \rfloor$ 
6:     if EMPTY( $A'[pos]$ ) then  $A'[pos] \leftarrow x$ 
7:     else COLLISION-HANDLER( $x$ )
8:   if  $o > 1$  then COMPACT( $A'$ )
9:   if NON-MONOTONIC then INSERTION-SORT( $A'$ )
10:  return  $A'$ 
```

---

# Cache-Efficient Learned Sort

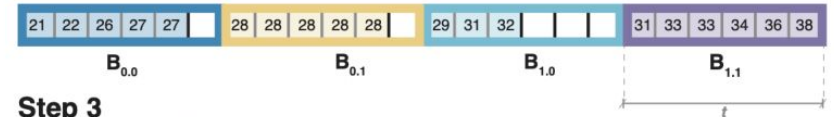
## Step 1

### Model-based bucketization



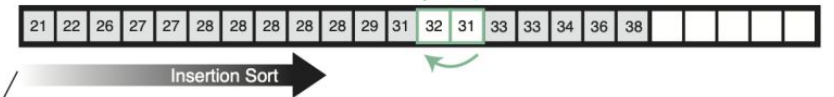
## Step 2

### In-bucket reordering



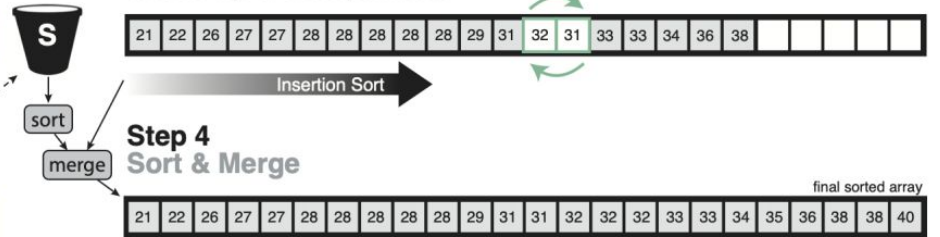
## Step 3

### Touch-up & Compaction



## Step 4

### Sort & Merge



# Pseudo-Code: Step 1

**Input**  $A$  - the array to be sorted

**Input**  $F_A$  - the CDF model for the distribution of  $A$

**Input**  $f$  - fan-out of the algorithm

**Input**  $t$  - threshold for bucket size

**Output**  $A'$  - the sorted version of array  $A$

```
1: procedure LEARNED-SORT( $A, F_A, f, t$ )
2:    $N \leftarrow |A|$                                 ▶ Size of the input array
3:    $n \leftarrow f$                                 ▶  $n$  represents the number of buckets
4:    $b \leftarrow \lfloor N/f \rfloor$                        ▶  $b$  represents the bucket capacity
5:    $B \leftarrow [] \times N$                          ▶ Empty array of size  $N$ 
6:    $I \leftarrow [0] \times n$                        ▶ Records bucket sizes
7:    $S \leftarrow []$                                 ▶ Spill bucket
8:   read_arr  $\leftarrow$  pointer to  $A$ 
9:   write_arr  $\leftarrow$  pointer to  $B$ 

10:  // Stage 1: Model-based bucketization
11:  while  $b \geq t$  do                               ▶ Until bucket capacity reaches the threshold  $t$ 
12:     $I \leftarrow [0] \times n$                        ▶ Reset array I
13:    for  $x \in$  read_arr do
14:       $pos \leftarrow \lfloor \text{INFER}(F_A, x) \cdot n \rfloor$ 
15:      if  $I[pos] \geq b$  then                         ▶ Bucket is full
16:        S.append( $x$ )                                ▶ Add to spill bucket
17:      else                                           ▶ Write into the predicted bucket
18:        write_arr[ $pos \cdot b + I[pos]$ ]  $\leftarrow x$ 
19:        INCREMENT  $I[pos]$ 
20:       $b \leftarrow \lfloor b/f \rfloor$                    ▶ Update bucket capacity
21:       $n \leftarrow \lfloor N/b \rfloor$                    ▶ Update the number of buckets
22:      PTRSWP(read_arr, write_arr)                   ▶ Pointer swap to reuse memory
```



# Pseudo-Code: Steps 2-4

```
23: // Stage 2: In-bucket reordering
24: offset  $\leftarrow$  0
25: for  $i \leftarrow 0$  up to  $n$  do                                 $\triangleright$  Process each bucket
26:      $K \leftarrow [0] \times b$                                  $\triangleright$  Array of counts

27:     for  $j \leftarrow 0$  up to  $I[i]$  do  $\triangleright$  Record the counts of the predicted positions
28:          $pos \leftarrow \lfloor \text{INFER}(F_A, \text{read\_arr}[\text{offset} + j]) \cdot N \rfloor$ 
29:         INCREMENT  $K[pos - \text{offset}]$ 

30:     for  $j \leftarrow 1$  up to  $|K|$  do                         $\triangleright$  Calculate the running total
31:          $K[j] \leftarrow K[j] + K[j - 1]$ 

32:      $T \leftarrow []$                                          $\triangleright$  Temporary auxiliary memory
33:     for  $j \leftarrow 0$  up to  $I[i]$  do                         $\triangleright$  Order keys w.r.t. the cumulative counts
34:          $pos \leftarrow \lfloor \text{INFER}(F_A, \text{read\_arr}[\text{offset} + j]) \cdot N \rfloor$ 
35:          $T[j] \leftarrow \text{read\_arr}[\text{offset} + K[pos - \text{offset}]]$ 
36:         DECREMENT  $K[pos - \text{offset}]$ 
37:     Copy  $T$  back to  $\text{read\_arr}[\text{offset}]$ 
38:     offset  $\leftarrow$  offset +  $b$ 

39: // Stage 3: Touch-up
40: INSERTION-SORT-AND-COMPACT(read_arr)

41: // Stage 4: Sort & Merge
42: SORT( $S$ )
43:  $A' \leftarrow \text{MERGE}(\text{read\_arr}, S)$ 

44: return  $A'$ 
```



# Optimizations

- Process elements in batches (cache locality)
- One bucket at a time (temporal locality)
- Bucket buffer space (reduce overflows)

# CDF Model

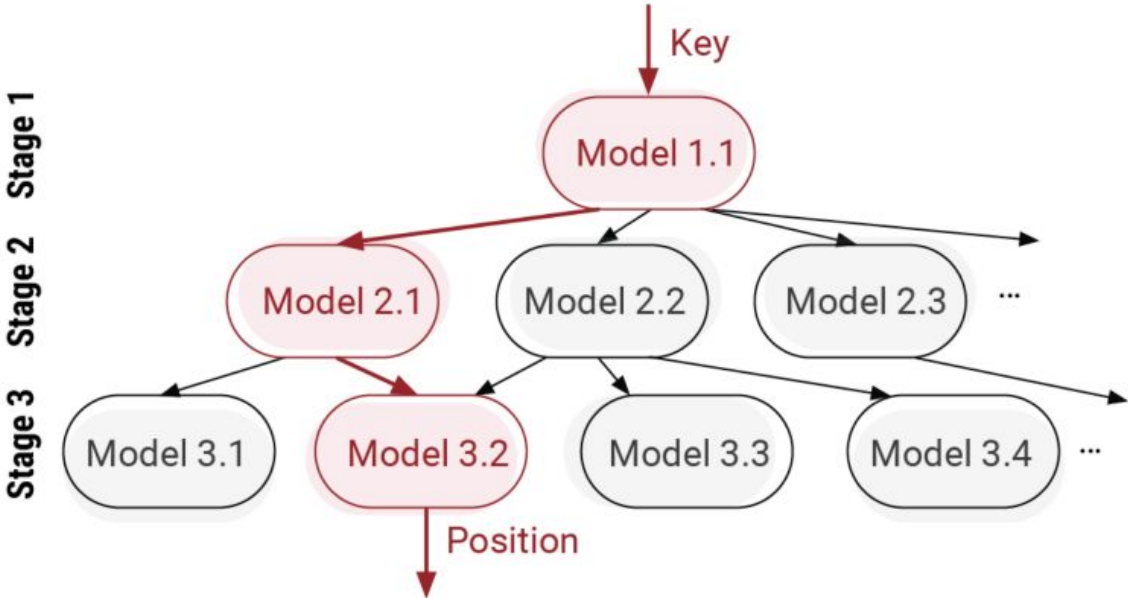


Figure 5: A typical RMI architecture containing three layers

# CDF Model

---

## Algorithm 3 The inference procedure for the CDF model

---

**Input**  $F_A$  - the trained model ( $F_A[l][r]$  refers to the  $r^{th}$  model in the  $l^{th}$  layer)

**Input**  $x$  - the key

**Output**  $r$  - the predicted rank (between 0-1)

```
1: procedure INFER( $F_A, x$ )
2:    $L \leftarrow$  the number of layers of the CDF model  $F_A$ 
3:    $M^l \leftarrow$  the number of models in the  $l^{th}$  layer of the RMI  $F_A$ 
4:    $r \leftarrow 0$ 
5:   for  $l \leftarrow 0$  up to  $L$  do
6:      $r = x \cdot F_A[l][r].\text{slope} + F_A[l][r].\text{intercept}$ 
7:   return  $r$ 
```

---

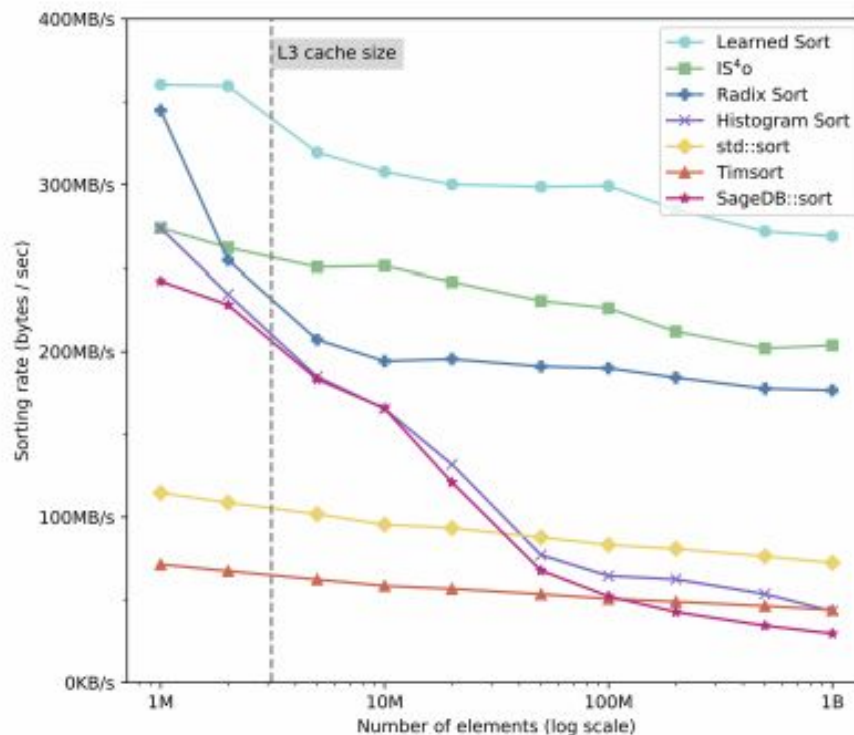


## Theoretical Results

- Step 1:  $O(N \cdot L)$
- Step 2:  $O(N)$
- Step 3:  $O(Nt)$  (non-dominant)
- Step 4:  $O(s \log s) + O(N)$

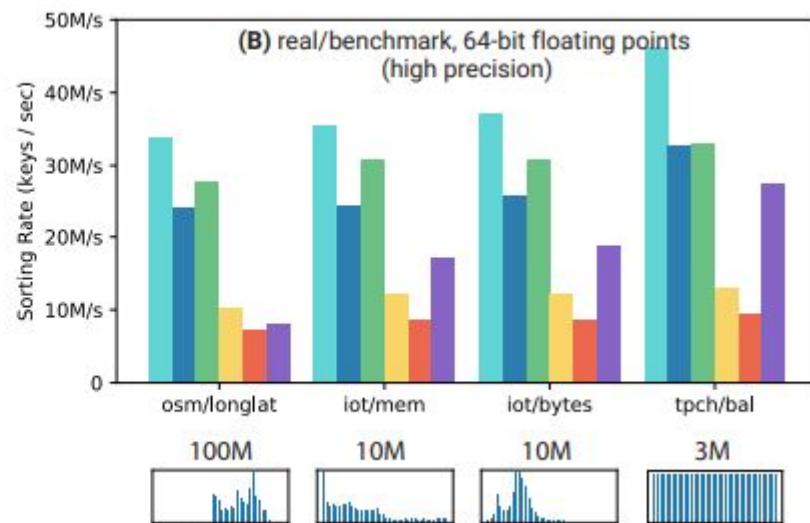
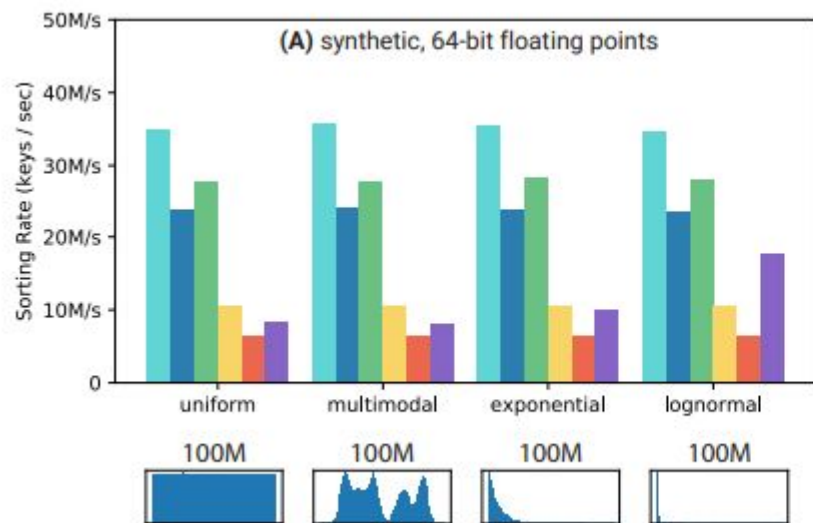
Space complexity: order of  $O(N)$

# Experimental Results



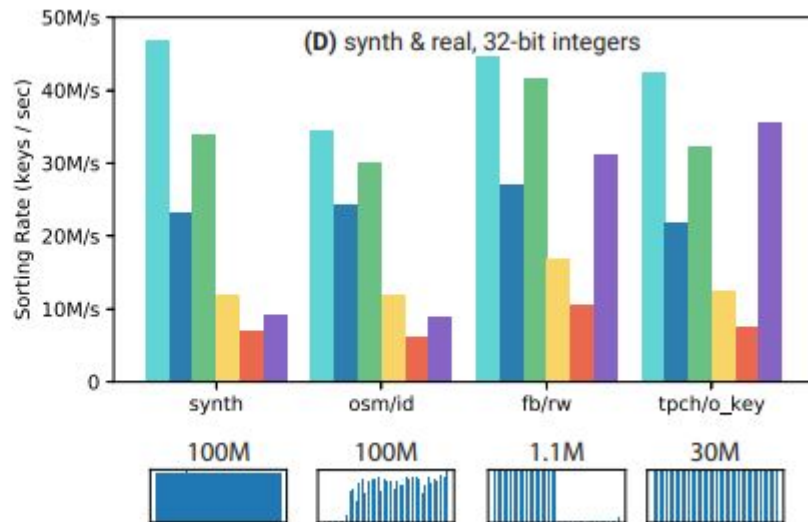
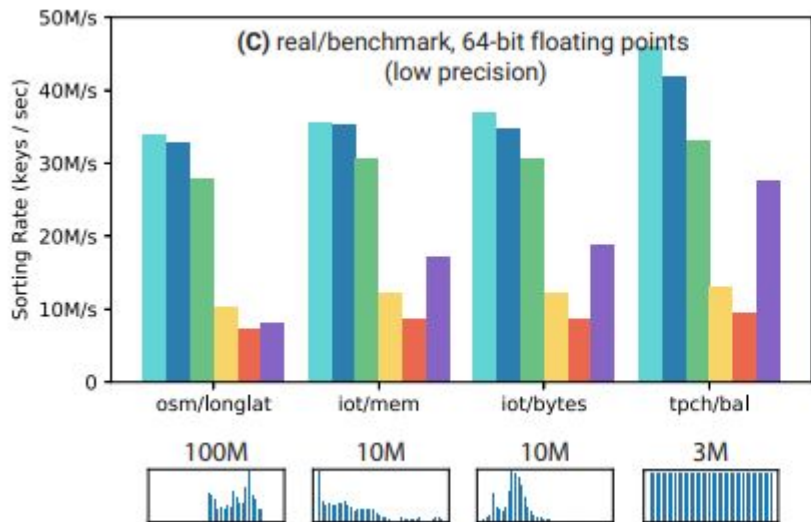
**Figure 8: The sorting throughput for normally distributed double-precision keys (higher is better).**

# Experimental Results



Legend: Learned Sort (cyan), Radix Sort (blue), IS<sup>4</sup>o (green), std::sort (yellow), Timsort (red), Histogram Sort (purple)

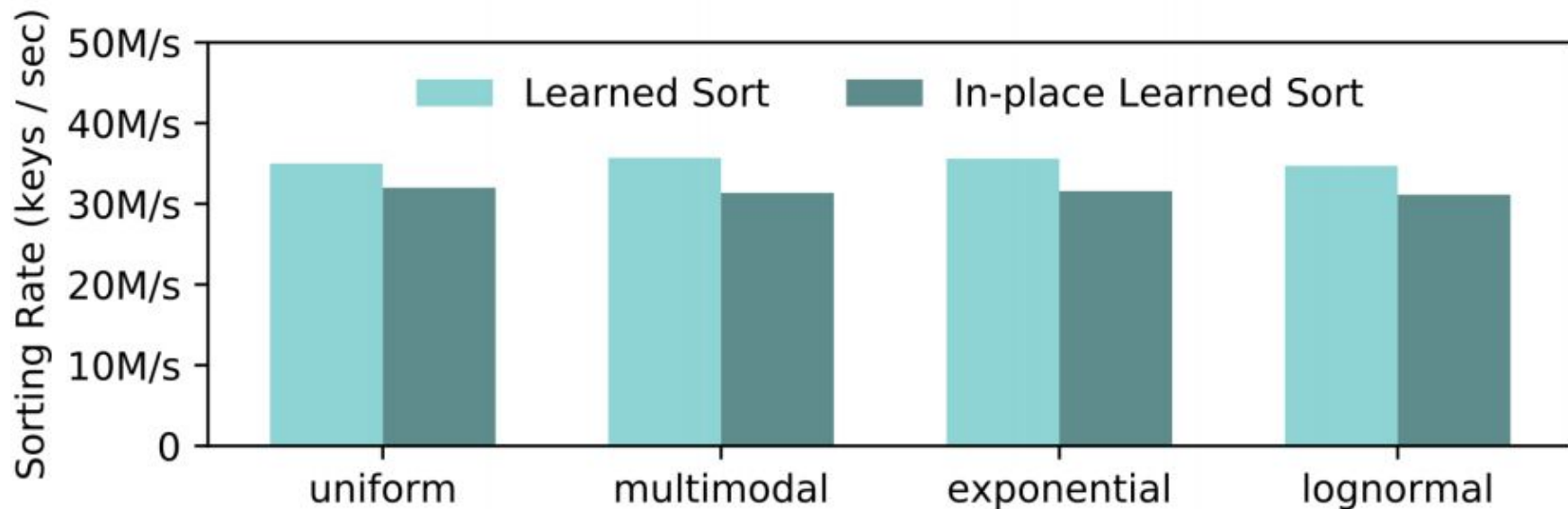
# Experimental Results



Legend: Learned Sort (cyan), Radix Sort (blue), IS<sup>4</sup>o (green), std::sort (yellow), Timsort (orange), Histogram Sort (purple)

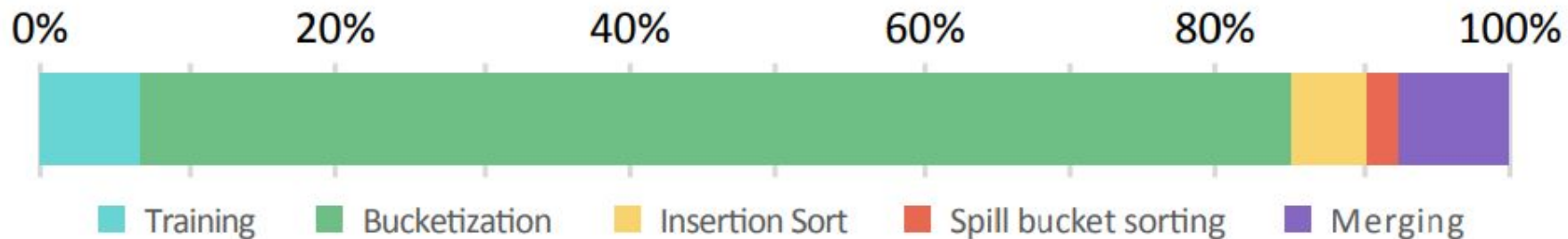


## In-Place Sorting



**Figure 12: The sorting rate of Learned Sort and its in-place version for all of our synthetic datasets.**

# Performance Decomposition



**Figure 13: Performance of each of the stages of Learned Sort.**



# Strengths/Weaknesses

## Strengths

- Performance on real-world data
- Improvement over default Java/Python sorting function
- Cache-efficient
- Model training time accounted for

## Weaknesses

- Other CDF implementations?
- Duplicate keys



## Directions for Future Work

- Sorting complex objects
- Parallel Sorting
- Using in DB systems



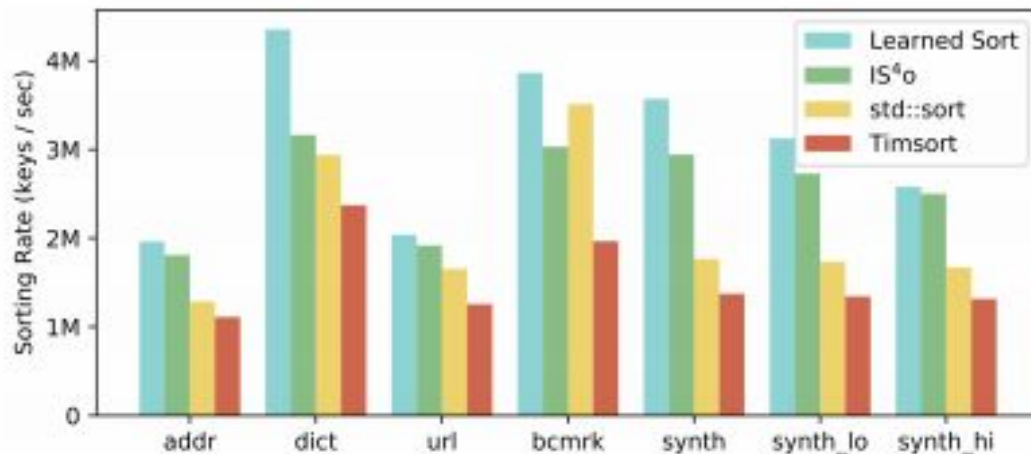
## Discussion Questions

- Can you think of adversarial inputs that may be good to evaluate this specific approach on?
- What parallelization techniques may apply to this algorithm/sorting algorithms in general?
- What are other ways through which collisions might be handled? What is attractive about the spill bucket method?

---

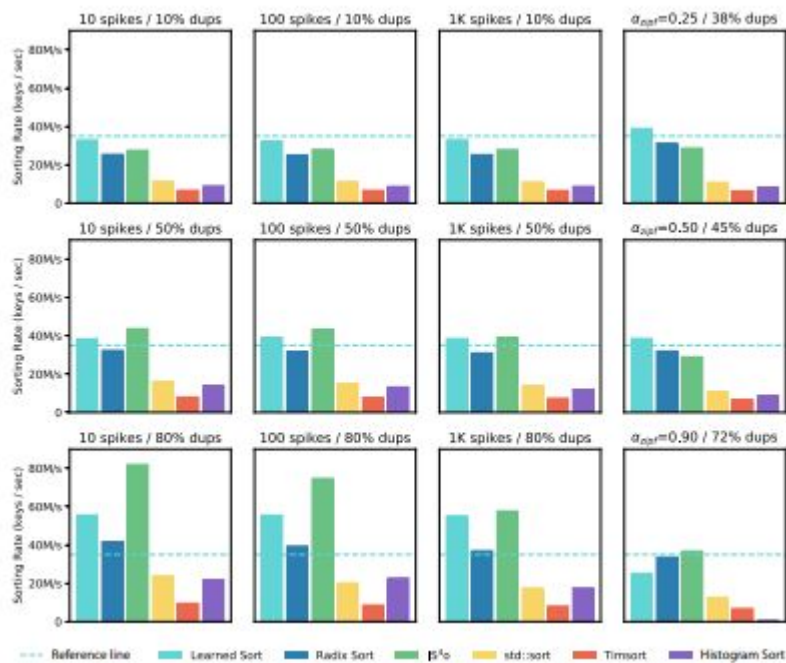
# Additional Materials

## String Sorting



**Figure 10: The sorting rate for various strings datasets.**

# Duplicates







# CDF Model Training

---

## Algorithm 4 The training procedure for the CDF model

---

**Input**  $A$  - the input array  
**Input**  $L$  - the number of layers of the CDF model  
**Input**  $M^l$  - the number of linear models in the  $l^{th}$  layer of the CDF model  
**Output**  $F_A$  - the trained CDF model with RMI architecture

```
1: procedure TRAIN( $A, L, M$ )
2:    $S \leftarrow \text{SAMPLE}(A)$ 
3:   SORT( $S$ )
4:    $T \leftarrow [][][]$  ▷ Training sets implemented as a 3D array
5:   for  $i \leftarrow 0$  up to  $|S|$  do
6:      $T[0][0].\text{add}((S[i], i/|S|))$ 
7:   for  $l \leftarrow 0$  up to  $L$  do
8:     for  $m \leftarrow 0$  up to  $M^l$  do
9:        $F_A[l][m] \leftarrow$  linear model trained on the set  $\{t \mid t \in T[l][m]\}$ 
10:      if  $l + 1 < L$  then
11:        for  $t \in T[l][m]$  do
12:           $F_A[l][m].\text{slope} \leftarrow F_A[l][m].\text{slope} \cdot M^{l+1}$ 
13:           $F_A[l][m].\text{intercept} \leftarrow F_A[l][m].\text{intercept} \cdot M^{l+1}$ 
14:           $i \leftarrow F_A[l][m].\text{slope} \cdot t + F_A[l][m].\text{intercept}$ 
15:           $T[l + 1][i].\text{add}(t)$ 
16:   return  $F_A$ 
```

---