# String Sorts

Mark Matthews PhD

# Comparison Sorts

We have seen from before that comparison sorting algorithms requires N Log N  compares in the worst case.

There's no way for us to beat this performance taking a comparison approach no matter what clever tricks we employ.

But what if we don't compare?

# Counting sort

**Assumptions:**

- $n$ records
- Each record contains keys and data
- All keys are in the range of 1 to k

**Space**

- The unsorted list is stored in A, the sorted list will be stored in an additional array B
- Uses an additional array C of size k

# Counting sort

**Main idea:**

1. For each key value $i$, $i = 1,…,k$, count the number of times the keys occurs in the unsorted input array $A$.

   Store results in an auxiliary array, C

2. Use these counts to compute the offset. Offset is used to calculate the location where the record with key value $i$ will be stored in the sorted output list $B$.

# Key-indexed counting: assumptions about keys

Assumption. Keys are integers between $0$ and $R - 1$.

Implication. Can use key as an array index.

Applications.

- Sort string by first letter.
- Sort class roster by section.
- Sort phone numbers by area code.
- Subroutine in a sorting algorithm. [we'll see later]

Remark. Keys may have associated data $\Rightarrow$

can't just count up number of keys of each value.

| input | | sorted result | |
|-------|---------|---------------|---|
| *name* | *section* | *(by section)* | |
| Anderson | 2 | Harris | 1 |
| Brown | 3 | Martin | 1 |
| Davis | 3 | Moore | 1 |
| Garcia | 4 | Anderson | 2 |
| Harris | 1 | Martinez | 2 |
| Jackson | 3 | Miller | 2 |
| Johnson | 4 | Robinson | 2 |
| Jones | 3 | White | 2 |
| Martin | 1 | Brown | 3 |
| Martinez | 2 | Davis | 3 |
| Miller | 2 | Jackson | 3 |
| Moore | 1 | Jones | 3 |
| Robinson | 2 | Taylor | 3 |
| Smith | 4 | Williams | 3 |
| Taylor | 3 | Garcia | 4 |
| Thomas | 4 | Johnson | 4 |
| Thompson | 4 | Smith | 4 |
| White | 2 | Thomas | 4 |
| Williams | 3 | Thompson | 4 |
| Wilson | 4 | Wilson | 4 |

*keys are*
*small integers*

5

Goal. Sort an array a[] of $N$ integers between $0$ and $R - 1$.

- Count frequencies of each letter using key as index.
- Compute frequency cumulates which specify destinations.
- Access cumulates using key as index to move items.
- Copy back into original array.

offset by 1

[stay tuned]

i  a[i]

```
int N = a.length;
int[] count = new int[R+1];


for (int i = 0; i < N; i++)
    count[a[i]+1]++;



for (int r = 0; r < R; r++)
    count[r+1] += count[r];



for (int i = 0; i < N; i++)
    aux[count[a[i]]++] = a[i];



for (int i = 0; i < N; i++)
```

count frequencies

| i | a[i] |
|---|------|
| 0 | d |
| 1 | a |
| 2 | c |
| 3 | f |
| 4 | f |
| 5 | b |
| 6 | d |
| 7 | b |
| 8 | f |
| 9 | b |
| 10 | e |
| 11 | a |

r count[r]

| r | count[r] |
|---|----------|
| a | 0 |
| b | 2 |
| c | 3 |
| d | 1 |
| e | 2 |
| f | 1 |
| - | 3 |

Goal.  Sort an array a[] of $N$ integers between $0$ and $R - 1$.

- Count frequencies of each letter using key as index.

- **Compute frequency cumulates which specify destinations.**

- Access cumulates using key as index to move items.

- Copy back into original array.

i  a[i]

```
int N = a.length;
int[] count = new int[R+1];


for (int i = 0; i < N; i++)
   count[a[i]+1]++;


for (int r = 0; r < R; r++)
   count[r+1] += count[r];


for (int i = 0; i < N; i++)
   aux[count[a[i]]++] = a[i];


for (int i = 0; i < N; i++)
```

compute
cumulates

| i | a[i] |
|---|------|
| 0 | d |
| 1 | a |
| 2 | c |
| 3 | f |
| 4 | f |
| 5 | b |
| 6 | d |
| 7 | b |
| 8 | f |
| 9 | b |
| 10 | e |
| 11 | a |

r count[r]

| r | count[r] |
|---|----------|
| a | 0 |
| b | 2 |
| c | 5 |
| d | 6 |
| e | 8 |
| f | 9 |
|  | 12 |

6 keys < d, 8 keys < e
so d's go in a[6] and a[7]

# Key-indexed counting demo

Goal. Sort an array a[] of $N$ integers between $0$ and $R - 1$.

- Count frequencies of each letter using key as index.
- Compute frequency cumulates which specify destinations.
- **Access cumulates using key as index to move items.**
- Copy back into original array.

```
int N = a.length;
int[] count = new int[R+1];


for (int i = 0; i < N; i++)

    count[a[i]+1]++;


for (int r = 0; r < R; r++)

    count[r+1] += count[r];


for (int i = 0; i < N; i++)

    aux[count[a[i]]++] = a[i];



for (int i = 0; i < N; i++)
```

move items →

| i | a[i] |
|---|------|
| 0 | d |
| 1 | a |
| 2 | c |
| 3 | f |
| 4 | f |
| 5 | b |
| 6 | d |
| 7 | b |
| 8 | f |
| 9 | b |
| 10 | e |
| 11 | a |

| r | count[r] |
|---|------|
| a | 2 |
| b | 5 |
| c | 6 |
| d | 8 |
| e | 9 |
| f | 12 |
| - | 12 |

| i | aux[i] |
|---|------|
| 0 | a |
| 1 | a |
| 2 | b |
| 3 | b |
| 4 | b |
| 5 | c |
| 6 | d |
| 7 | d |
| 8 | e |
| 9 | f |
| 10 | f |
| 11 | f |

Goal. Sort an array a[] of $N$ integers between $0$ and $R - 1$.

- Count frequencies of each letter using key as index.
- Compute frequency cumulates which specify destinations.
- Access cumulates using key as index to move items.
- Copy back into original array.

```
int N = a.length;
int[] count = new int[R+1];


for (int i = 0; i < N; i++)
    count[a[i]+1]++;


for (int r = 0; r < R; r++)
    count[r+1] += count[r];


for (int i = 0; i < N; i++)
    aux[count[a[i]]++] = a[i];
```

copy
back

```
for (int i = 0; i < N; i++)
```
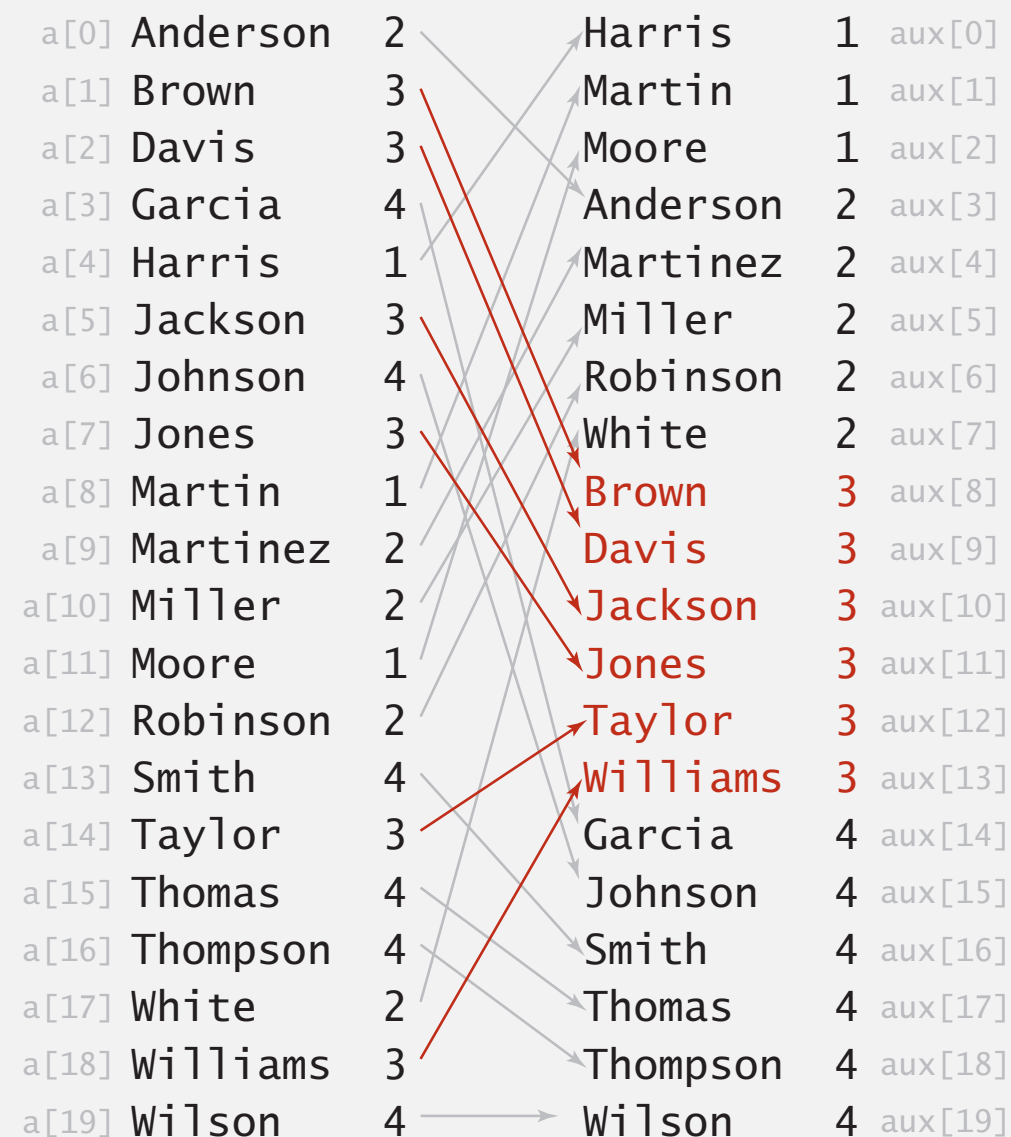
| i | a[i] |
|---|------|
| 0 | a |
| 1 | a |
| 2 | b |
| 3 | b |
| 4 | b |
| 5 | c |
| 6 | d |
| 7 | d |
| 8 | e |
| 9 | f |
| 10 | f |
| 11 | f |

| r | count[r] |
|---|----------|
| a | 2 |
| b | 5 |
| c | 6 |
| d | 8 |
| e | 9 |
| f | 12 |
| - | 12 |

| i | aux[i] |
|---|--------|
| 0 | a |
| 1 | a |
| 2 | b |
| 3 | b |
| 4 | b |
| 5 | c |
| 6 | d |
| 7 | d |
| 8 | e |
| 9 | f |
| 10 | f |
| 11 | f |

Proposition. Key-indexed takes time proportional to $N + R$.

Proposition. Key-indexed counting uses extra space proportional to $N + R$.

Stable?  ✔

# Radix Sort

- Radix: the base of a number system or logarithm.
- Uses a stable sorting algorithm (key indexing counting) so is also stable

- Radix sort is a multiple pass distribution sort.
- It distributes each item to a bucket according to part of the item's key.
- After each pass, items are collected from the buckets, keeping the items in order, then redistributed according to the next most significant part of the key.

- For numbers, it sorts keys digit-by-digit, for strings it sorts character-by-character.

# Example of Radix Sort

| 12 | 34 | 42 | 32 | 44 | 41 | 34 | 11 | 32 | 23 | 87 | 50 | 77 | 58 | 08 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

Input is an array of 15 integers. For integers, the number of buckets is 10, from 0 to 9. The first pass distributes the keys into buckets by the least significant digit (LSD). When the first pass is done, we have the following.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 50 | 41 | 12 | 23 | 34 | | | 87 | 58 | |
| | 11 | 42 | | 44 | | | 77 | 08 | |
| | | 32 | | 34 | | | | | |
| | | 32 | | | | | | | |

# Example of Radix Sort

We collect these, keeping their relative order:

| 50 | 41 | 11 | 12 | 42 | 32 | 32 | 23 | 34 | 44 | 34 | 87 | 77 | 58 | 08 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

Now we distribute by the next most significant digit, which is the highest digit in our example, and we get the following.

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|---|---|---|---|---|---|---|---|---|---|
| 08 | 11 | 23 | 32 | 41 | 50 |  | 77 | 87 |  |
|  | 12 |  | 32 | 42 | 58 |  |  |  |  |
|  |  |  | 34 | 44 |  |  |  |  |  |
|  |  |  | 34 |  |  |  |  |  |  |

When we collect them, they are in order.

| 08 | 11 | 12 | 23 | 32 | 32 | 34 | 34 | 41 | 42 | 44 | 50 | 58 | 77 | 87 |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|

# Least-significant-digit-first string sort

LSD string (radix) sort.

- Consider characters from right to left.

- Stably sort using $d^{th}$ character as the key (using key-indexed counting).

sort key (d = 2)   sort key (d = 1)   sort key (d = 0)

| | | | |
|---|---|---|---|
| 0 | d | a | b |
| 1 | a | d | d |
| 2 | c | a | b |
| 3 | f | a | d |
| 4 | f | e | e |
| 5 | b | a | d |
| 6 | d | a | d |
| 7 | b | e | e |
| 8 | f | e | d |
| 9 | b | e | d |
| 10 | e | b | b |
| 11 | a | c | e |

| | | | |
|---|---|---|---|
| 0 | d | a | b |
| 1 | c | a | b |
| 2 | e | b | b |
| 3 | a | d | d |
| 4 | f | a | d |
| 5 | b | a | d |
| 6 | d | a | d |
| 7 | f | e | d |
| 8 | b | e | d |
| 9 | f | e | e |
| 10 | b | e | e |
| 11 | a | c | e |

| | | | |
|---|---|---|---|
| 0 | d | a | b |
| 1 | c | a | b |
| 2 | f | a | d |
| 3 | b | a | d |
| 4 | d | a | d |
| 5 | e | b | b |
| 6 | a | c | e |
| 7 | a | d | d |
| 8 | f | e | d |
| 9 | b | e | d |
| 10 | f | e | e |
| 11 | b | e | e |

| | | | |
|---|---|---|---|
| 0 | a | c | e |
| 1 | a | d | d |
| 2 | b | a | d |
| 3 | b | e | d |
| 4 | b | e | e |
| 5 | c | a | b |
| 6 | d | a | b |
| 7 | d | a | d |
| 8 | e | b | b |
| 9 | f | a | d |
| 10 | f | e | d |
| 11 | f | e | e |

sort is stable

(arrows do not cross)

14

# LSD string sort:  Java implementation

```java
public class LSD
{
   public static void sort(String[] a, int W)
   {
      int R = 256;
      int N = a.length;
      String[] aux = new String[N];

      for (int d = W-1; d >= 0; d--)
      {
         int[] count = new int[R+1];
         for (int i = 0; i < N; i++)
            count[a[i].charAt(d) + 1]++;
         for (int r = 0; r < R; r++)
            count[r+1] += count[r];
         for (int i = 0; i < N; i++)
            aux[count[a[i].charAt(d)]++] = a[i];
         for (int i = 0; i < N; i++)
            a[i] = aux[i];
      }
```

fixed-length W strings

radix R

do key-indexed counting
for each digit from right to left

key-indexed counting

MSD string (radix) sort.

- Partition array into $R$ pieces according to first character (use key-indexed counting).
- Recursively sort all strings that start with each character (key-indexed counts delineate subarrays to sort).



count[]

| | |
|---|---|
| a | 0 |
| b | 2 |
| c | 5 |
| d | 6 |
| e | 8 |
| f | 9 |
| - | 12 |

sort key

sort subarrays recursively

# MSD string sort: example

**input**

| she | are | are | are | are | are | are | are | are |
|---|---|---|---|---|---|---|---|---|
| sells | by | by | by | by | by | by | by | by |
| seashells | she | sells | seashells | sea | sea | sea | sea | sea |
| by | sells | seashells | sea | seashells | seashells | seashells | seashells | seashells |
| the | seashells | sea | seashells | seashells | seashells | seashells | seashells | seashells |
| sea | sea | sells | sells | sells | sells | sells | sells | sells |
| shore | shore | seashells | sells | sells | sells | sells | sells | sells |
| the | shells | she | she | she | she | she | she | she |
| shells | she | shore | shore | shore | shore | shore | shore | shore |
| she | sells | shells | shells | shells | shells | shells | shells | shells |
| sells | surely | she | she | she | she | she | she | she |
| are | seashells | surely | surely | surely | surely | surely | surely | surely |
| surely | the | the | the | the | the | the | the | the |
| seashells | the | the | the | the | the | the | the | the |

*need to examine
every character
in equal keys*

*end of string
goes before any
char value*

**output**

| are | are | are | are | are | are | are | are |
|---|---|---|---|---|---|---|---|
| by | by | by | by | by | by | by | by |
| sea | sea | sea | sea | sea | sea | sea | sea |
| seashells | seashells | seashells | seashells | seashells | seashells | seashells | seashells |
| seashells | seashells | seashells | seashells | seashells | seashells | seashells | seashells |
| sells | sells | sells | sells | sells | sells | sells | sells |
| sells | sells | sells | sells | sells | sells | sells | sells |
| she | she | she | she | she | she | she | she |
| shore | sshore | shore | shells | she | she | she | she |
| shells | hells | shells | she | shells | shells | shells | shells |
| she | she | she | shore | shore | shore | shore | shore |
| surely | surely | surely | surely | surely | surely | surely | surely |
| the | the | the | the | the | the | the | the |
| the | the | the | the | the | the | the | the |

**Trace of recursive calls for MSD string sort (no cutoff for small subarrays, subarrays of size 0 and 1 omitted)**

17

# Variable-length strings

Treat strings as if they had an extra char at end (smaller than any char).

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 0 | s | e | a | -1 | | | | | |
| 1 | s | e | a | s | h | e | l | l | s | -1 |
| 2 | s | e | l | l | s | -1 | | | |
| 3 | s | h | e | -1 | | | | | |
| 4 | s | h | e | -1 | | | | | |
| 5 | s | h | e | l | l | s | -1 | | |
| 6 | s | h | o | r | e | -1 | | | |
| 7 | s | u | r | e | l | y | -1 | | |

she before shells

```
private static int charAt(String s, int d)
{
    if (d < s.length()) return s.charAt(d);
    else return -1;
}
```

C strings. Have extra char '\0' at end ⟹ no extra work needed.

# MSD string sort: Java implementation

```java
public static void sort(String[] a)
{

   aux = new String[a.length];

   sort(a, aux, 0, a.length - 1, 0);

}


private static void sort(String[] a, String[] aux, int lo, int hi, int d)
{

   if (hi <= lo) return;

   int[] count = new int[R+2];
   for (int i = lo; i <= hi; i++)
      count[charAt(a[i], d) + 2]++;
   for (int r = 0; r < R+1; r++)
      count[r+1] += count[r];
   for (int i = lo; i <= hi; i++)
      aux[count[charAt(a[i], d) + 1]++] = a[i];
   for (int i = lo; i <= hi; i++)
      a[i] = aux[i - lo];

   for (int r = 0; r < R; r++)
```
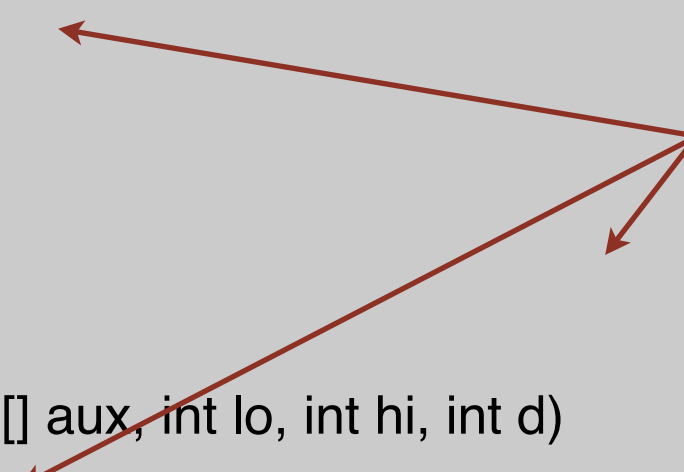
recycles aux[] array
but not count[] array

key-indexed counting

sort R subarrays recursively

# MSD string sort:  potential for disastrous performance

Observation 1.  Much too slow for small subarrays.

- Each function call needs its own count[] array.
- ASCII (256 counts):  100x slower than copy pass for $N = 2$.
- Unicode (65,536 counts):  32,000x slower for $N = 2$.

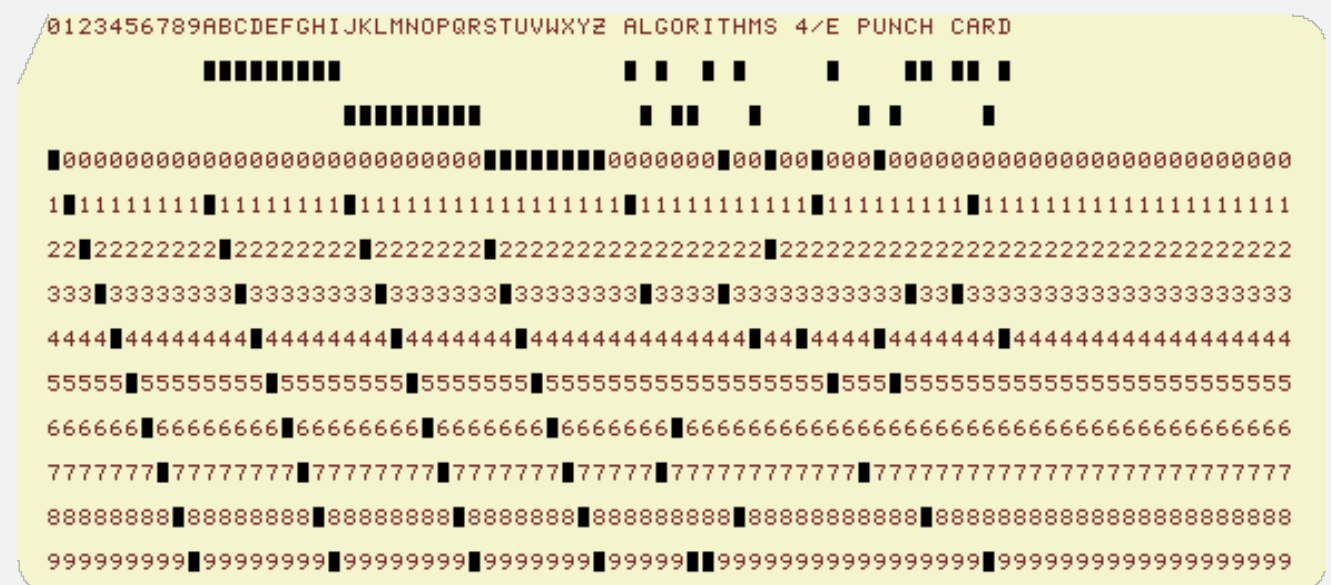Observation 2.  Huge number of small subarrays
because of recursion.

count[]

a[]

| 0 | b |
|---|---|
| 1 | a |

aux[]

| 0 | a |
|---|---|
| 1 | b |

# How to take a census in 1900s?

**1880 Census.** Took 1500 people 7 years to manually process data.

**Herman Hollerith.** Developed counting and sorting machine to automate.

- Use punch cards to record data (e.g., gender, age).

- Machine sorts one column at a time (into one of 12 bins).

- Typical question: how many women of age 20 to 30?



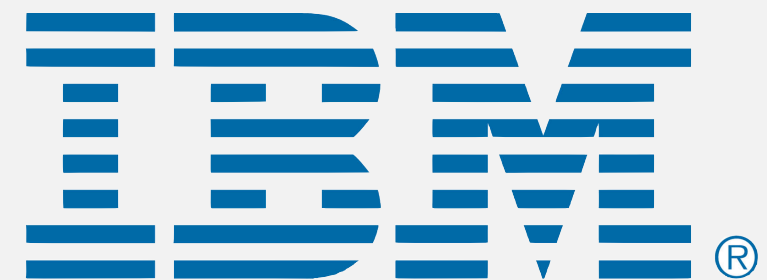**Hollerith tabulating machine and sorter**
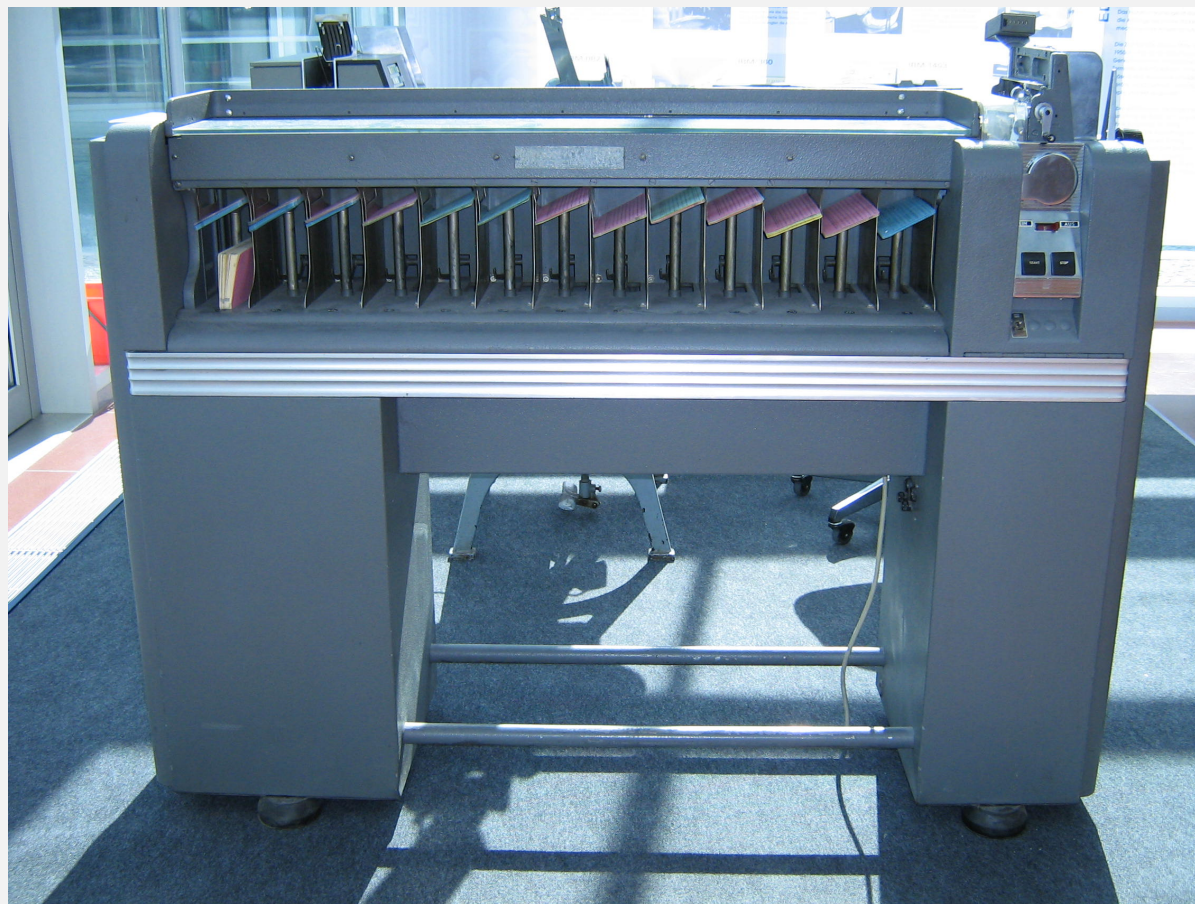


**punch card (12 holes per column)**

**1890 Census.** Finished in 1 year (and under budget)!

# How to get rich sorting in 1900s?

Punch cards. [1900s to 1950s]

- Also useful for accounting, inventory, and business processes.
- Primary medium for data entry, storage, and processing.

Hollerith's company later merged with 3 others to form Computing Tabulating Recording Corporation (CTRC); company renamed in 1924.



**IBM 80 Series Card Sorter (650 cards per minute)**

# String sorting interview question

Problem.  Sort one million 32-bit integers.

Ex.  Google (or presidential) interview.

Which sorting method to use?

- Insertion sort.
- Mergesort.
- Quicksort.
- Heapsort.
- LSD string sort.

# Radix sorts summary

Advantages

- Radix sorts are stable, preserving existing order of equal keys.
- They work in linear time, unlike most other sorts. In other words, they do not slow down when large numbers of items need to be sorted. Most sorts run in O(n log n) or O(n2) time.
- The time to sort per item is constant, as no comparisons among items are made. With other sorts, the time to sort per time increases with the number of items.
- Radix sort is particularly efficient when you have large numbers of records to sort with short keys.

Drawbacks

- Radix sorts do not work well when keys are very long, as the total sorting time is proportional to key length and to the number of items to sort.
- They are not "in-place", using more working memory than a traditional sort.