

Data Structures & Algorithms (7)

CS F211 / IS F 211

Outline of this slide

- Non-comparison-based sorting



Counting sort



The simplest setting

- Input: n **distinct** positive integers $A[1], \dots, A[n]$ with $A[i] \leq k$ for each $i = 1, 2, \dots, n$, for some known number $k = O(n)$.
- Output: a reordering of A such that the resulting list is monotonically increasing.



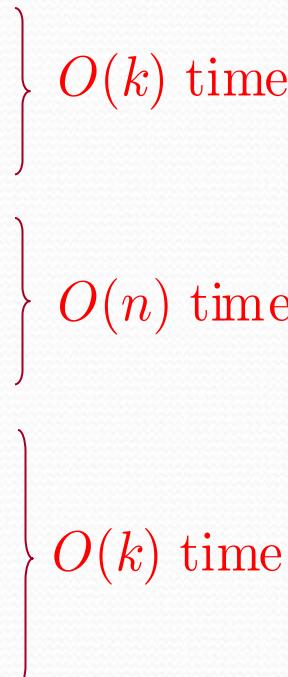
Solvable in $O(n)$ time?

Note that you are allowed to inspect the value of each $A[i]$.



An algorithm

```
allocate an array  $C[1 \dots k]$ ;            $O(1)$  time
for each  $j = 1$  to  $k$ 
    let  $C[j] = 0$ ;
for each  $i = 1$  to  $n$ 
    increment  $C[A[i]]$  by one;
for each  $j = 1$  to  $k$ 
    if  $C[j] > 0$ 
        print  $j$ ;
```



The running time is $O(n + k) = O(n)$.

Slightly harder setting

- Input: n **distinct** positive integers $A[1], \dots, A[n]$ with $A[i] \leq k$ for each $i = 1, 2, \dots, n$, for some **known** number $k = O(n)$.
- Output: a reordering of A such that the resulting list is monotonically increasing.



An algorithm

```
allocate an array  $C[1 \dots k]$ ;            $O(1)$  time  
for each  $j = 1$  to  $k$                       }  $O(k)$  time  
    let  $C[j] = 0$ ;  
for each  $i = 1$  to  $n$                       }  $O(n)$  time  
    increment  $C[A[i]]$  by one;  
for each  $j = 1$  to  $k$                       }  $O(k) + O(n)$  time  
    for each  $\ell = 1$  to  $C[j]$   
        print  $j$ ;
```

The running time is $O(n + k) = O(n)$.

A general setting

- Input: n ~~distinct~~ positive integers $A[1], \dots, A[n]$ with $A[i] \leq k$ for each $i = 1, 2, \dots, n$, for some unknown number $k = O(n)$.
- Output: a reordering of A such that the resulting list is monotonic increasing.



An algorithm

```
let  $k$  be the maximum of  $A[1], \dots, A[n]$ ;  $O(n)$  time  
allocate an array  $C[1 \dots k]$ ;  
for each  $j = 1$  to  $k$   
    let  $C[j] = 0$ ;  
for each  $i = 1$  to  $n$   
    increment  $C[A[i]]$  by one;  
for each  $j = 1$  to  $k$   
    for each  $\ell = 1$  to  $C[j]$   
        print  $j$ ;
```

The running time is $O(n + k) = O(n)$.

Counting-sort example

1	2	3	4	5	
A:	4	1	3	4	3

1	2	3	4
C:			

Loop 1

1	2	3	4	5	
$A:$	4	1	3	4	3
$C:$	0	0	0	0	0

```
for  $j \leftarrow 1$  to  $k$   
do  $C[j] \leftarrow 0$ 
```

Loop 2

1	2	3	4	5	
$A:$	4	1	3	4	3
$C:$	0	0	0	1	

```
for  $i \leftarrow 1$  to  $n$   
do  $C[A[i]] \leftarrow C[A[i]] + 1$ 
```

Loop 2

1	2	3	4	5	
$A:$	4	1	3	4	3
$C:$	1	0	0	1	

```
for  $i \leftarrow 1$  to  $n$   
do  $C[A[i]] \leftarrow C[A[i]] + 1$ 
```

Loop 2

1	2	3	4	5	
$A:$	4	1	3	4	3
$C:$	1	0	1	1	

```
for  $i \leftarrow 1$  to  $n$   
do  $C[A[i]] \leftarrow C[A[i]] + 1$ 
```

Loop 2

	1	2	3	4	5	
$A:$	4	1	3	4	3	
$C:$	1	0	1	2		

```
for  $i \leftarrow 1$  to  $n$   
do  $C[A[i]] \leftarrow C[A[i]] + 1$ 
```

Loop 2

1	2	3	4	5	
$A:$	4	1	3	4	3
$C:$	1	0	2	2	

For $i \leftarrow 1$ **to** n
do $C[A[i]] \leftarrow C[A[i]] + 1$ $\triangleright C[i] = |\{ \text{key} = i \}|$

Loop 3

1	2	3	4	5	
A:	4	1	3	4	3

1	2	3	4	
C:	1	0	2	2

```
for each  $j = 1$  to  $k$   
  for each  $\ell = 1$  to  $C[j]$   
    print  $j$ ;
```

1
1
1, 3
1, 3, 3
1, 3, 3, 4
1, 3, 3, 4, 4

Stable Sort

A sorting algorithm is **stable** if

- Numbers with the same value appear in the output array in the same order as they do in the input array
- Ties between two numbers are broken by the rule that which ever number appears first in the input array appears first in the output array

Stable Sort



Stable Sort

Counting Sort (stable) Overview

- Assumption: n input elements are integers in the range of 0 to k (integer).
 - $\Theta(n+k)$ → When $k=O(n)$, counting sort: $\Theta(n)$
- Basic Idea
 - For each input element x , determine the number of elements less than or equal to x
 - For each integer i ($0 \leq i \leq k$), count how many elements whose values are i
 - Then we know how many elements are less than or equal to i
- Algorithm storage
 - $A[1..n]$: input elements
 - $B[1..n]$: sorted elements
 - $C[1..k]$: hold the number of elements less than or equal to i

Counting Sort

```
for  $i \leftarrow 1$  to  $k$ 
    do  $C[i] \leftarrow 0$ 
```

```
for  $j \leftarrow 1$  to  $n$ 
    do  $C[A[j]] \leftarrow C[A[j]] + 1$       ▷  $C[i] = |\{\text{key} = i\}|$ 
```

```
for  $i \leftarrow 2$  to  $k$ 
    do  $C[i] \leftarrow C[i] + C[i-1]$       ▷  $C[i] = |\{\text{key} \leq i\}|$ 
```

```
for  $j \leftarrow n$  downto 1
    do  $B[C[A[j]]] \leftarrow A[j]$ 
         $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Counting-sort example

1	2	3	4	5	
A:	4	1	3	4	3

1	2	3	4
C:			

B:					
----	--	--	--	--	--

Loop 1

	1	2	3	4	5
$A:$	4	1	3	4	3

	1	2	3	4
$C:$	0	0	0	0

$B:$					
------	--	--	--	--	--

```
for  $i \leftarrow 1$  to  $k$   
do  $C[i] \leftarrow 0$ 
```

Loop 2

1	2	3	4	5
A: 4	1	3	4	3

1	2	3	4
C: 0	0	0	1

B:					
----	--	--	--	--	--

for $j \leftarrow 1$ **to** n
do $C[A[j]] \leftarrow C[A[j]] + 1$ $\triangleright C[i] = |\{ \text{key} = i \}|$

Loop 2

1	2	3	4	5
A: 4	1	3	4	3

1	2	3	4
C: 1	0	0	1

B:					
----	--	--	--	--	--

for $j \leftarrow 1$ **to** n
do $C[A[j]] \leftarrow C[A[j]] + 1$ $\triangleright C[i] = |\{ \text{key} = i \}|$

Loop 2

1	2	3	4	5
A: 4	1	3	4	3

1	2	3	4
C: 1	0	1	1

B:					
----	--	--	--	--	--

for $j \leftarrow 1$ **to** n
do $C[A[j]] \leftarrow C[A[j]] + 1$ $\triangleright C[i] = |\{ \text{key} = i \}|$

Loop 2

1	2	3	4	5
A: 4	1	3	4	3

1	2	3	4
C: 1	0	1	2

B:					
----	--	--	--	--	--

for $j \leftarrow 1$ **to** n
do $C[A[j]] \leftarrow C[A[j]] + 1$ $\triangleright C[i] = |\{ \text{key} = i \}|$

Loop 2

	1	2	3	4	5	
$A:$	4	1	3	4	3	
$C:$	1	0	2	2		
$B:$						

```
for  $j \leftarrow 1$  to  $n$ 
  do  $C[A[j]] \leftarrow C[A[j]] + 1$       ▷  $C[i] = |\{key = i\}|$ 
```

Loop 3

	1	2	3	4	5
$A:$	4	1	3	4	3

	1	2	3	4
$C:$	1	0	2	2

$B:$					
------	--	--	--	--	--

$C':$	1	1	2	2
-------	---	---	---	---

for $i \leftarrow 2$ to k

do $C[i] \leftarrow C[i] + C[i-1]$

▷ $C[i] = |\{\text{key} \leq i\}|$

Loop 3

	1	2	3	4	5
$A:$	4	1	3	4	3

	1	2	3	4
$C:$	1	0	2	2

$B:$					
------	--	--	--	--	--

$C':$	1	1	3	2
-------	---	---	---	---

for $i \leftarrow 2$ **to** k

do $C[i] \leftarrow C[i] + C[i-1]$ $\blacktriangleright C[i] = |\{\text{key} \leq i\}|$

Loop 3

	1	2	3	4	5
$A:$	4	1	3	4	3

	1	2	3	4
$C:$	1	0	2	2

$B:$					
------	--	--	--	--	--

$C':$	1	1	3	5
-------	---	---	---	---

for $i \leftarrow 2$ **to** k

do $C[i] \leftarrow C[i] + C[i-1]$ $\blacktriangleright C[i] = |\{\text{key} \leq i\}|$

Loop 4

1	2	3	4	5	
A:	4	1	3	4	3

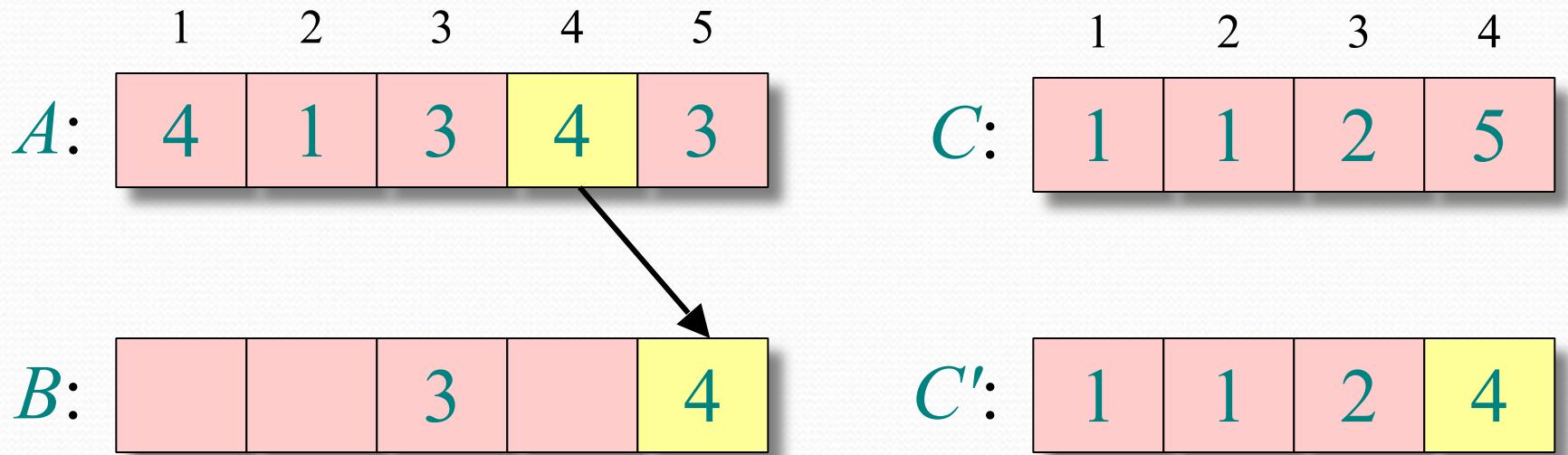
B:			3		
----	--	--	---	--	--

1	2	3	4	
C:	1	1	3	5

C':	1	1	2	5
-----	---	---	---	---

```
for  $j \leftarrow n$  downto 1  
do  $B[C[A[j]]] \leftarrow A[j]$   
 $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Loop 4



```
for  $j \leftarrow n$  downto 1  
  do  $B[C[A[j]]] \leftarrow A[j]$   
       $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Loop 4

	1	2	3	4	5
$A:$	4	1	3	4	3

$B:$		3	3		4
------	--	---	---	--	---

	1	2	3	4
$C:$	1	1	2	4

$C':$	1	1	1	4
-------	---	---	---	---

```
for  $j \leftarrow n$  downto 1  
do  $B[C[A[j]]] \leftarrow A[j]$   
 $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Loop 4

1	2	3	4	5	
A:	4	1	3	4	3

B:	1	3	3		4
----	---	---	---	--	---

1	2	3	4	
C:	1	1	1	4

C':	0	1	1	4
-----	---	---	---	---

```
for  $j \leftarrow n$  downto 1  
do  $B[C[A[j]]] \leftarrow A[j]$   
 $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Loop 4

	1	2	3	4	5
$A:$	4	1	3	4	3

$B:$	1	3	3	4	4
------	---	---	---	---	---

	1	2	3	4
$C:$	0	1	1	4

$C':$	0	1	1	3
-------	---	---	---	---

```
for  $j \leftarrow n$  downto 1  
do  $B[C[A[j]]] \leftarrow A[j]$   
 $C[A[j]] \leftarrow C[A[j]] - 1$ 
```

Analysis

$$\Theta(k) \quad \left\{ \begin{array}{l} \textbf{for } i \leftarrow 1 \text{ to } k \\ \quad \textbf{do } C[i] \leftarrow 0 \end{array} \right.$$
$$\Theta(n) \quad \left\{ \begin{array}{l} \textbf{for } j \leftarrow 1 \text{ to } n \\ \quad \textbf{do } C[A[j]] \leftarrow C[A[j]] + 1 \end{array} \right.$$
$$\Theta(k) \quad \left\{ \begin{array}{l} \textbf{for } i \leftarrow 2 \text{ to } k \\ \quad \textbf{do } C[i] \leftarrow C[i] + C[i-1] \end{array} \right.$$
$$\Theta(n) \quad \left\{ \begin{array}{l} \textbf{for } j \leftarrow n \text{ downto } 1 \\ \quad \textbf{do } B[C[A[j]]] \leftarrow A[j] \\ \quad \quad C[A[j]] \leftarrow C[A[j]] - 1 \end{array} \right.$$

$$\Theta(n + k)$$

Bucket-Sort



- Let be S be a sequence of n (key, element) entries with keys in the range $[0, N - 1]$
- Bucket-sort uses the keys as indices into an auxiliary array B of sequences (buckets)

Phase 1: Empty sequence S by moving each entry (k, o) into its bucket $B[k]$

Phase 2: For $i = 0, \dots, N - 1$, move the entries of bucket $B[i]$ to the end of sequence S

- Analysis:
 - Phase 1 takes $O(n)$ time
 - Phase 2 takes $O(n + N)$ time

Bucket-sort takes $O(n + N)$ time

Algorithm *bucketSort(S, N)*

Input sequence S of (key, element) items with keys in the range $[0, N - 1]$

Output sequence S sorted by increasing keys

$B \leftarrow$ array of N empty sequences

while $\neg S.isEmpty()$

$f \leftarrow S.first()$

$(k, o) \leftarrow S.remove(f)$

$B[k].insertLast((k, o))$

for $i \leftarrow 0$ to $N - 1$

while $\neg B[i].isEmpty()$

$f \leftarrow B[i].first()$

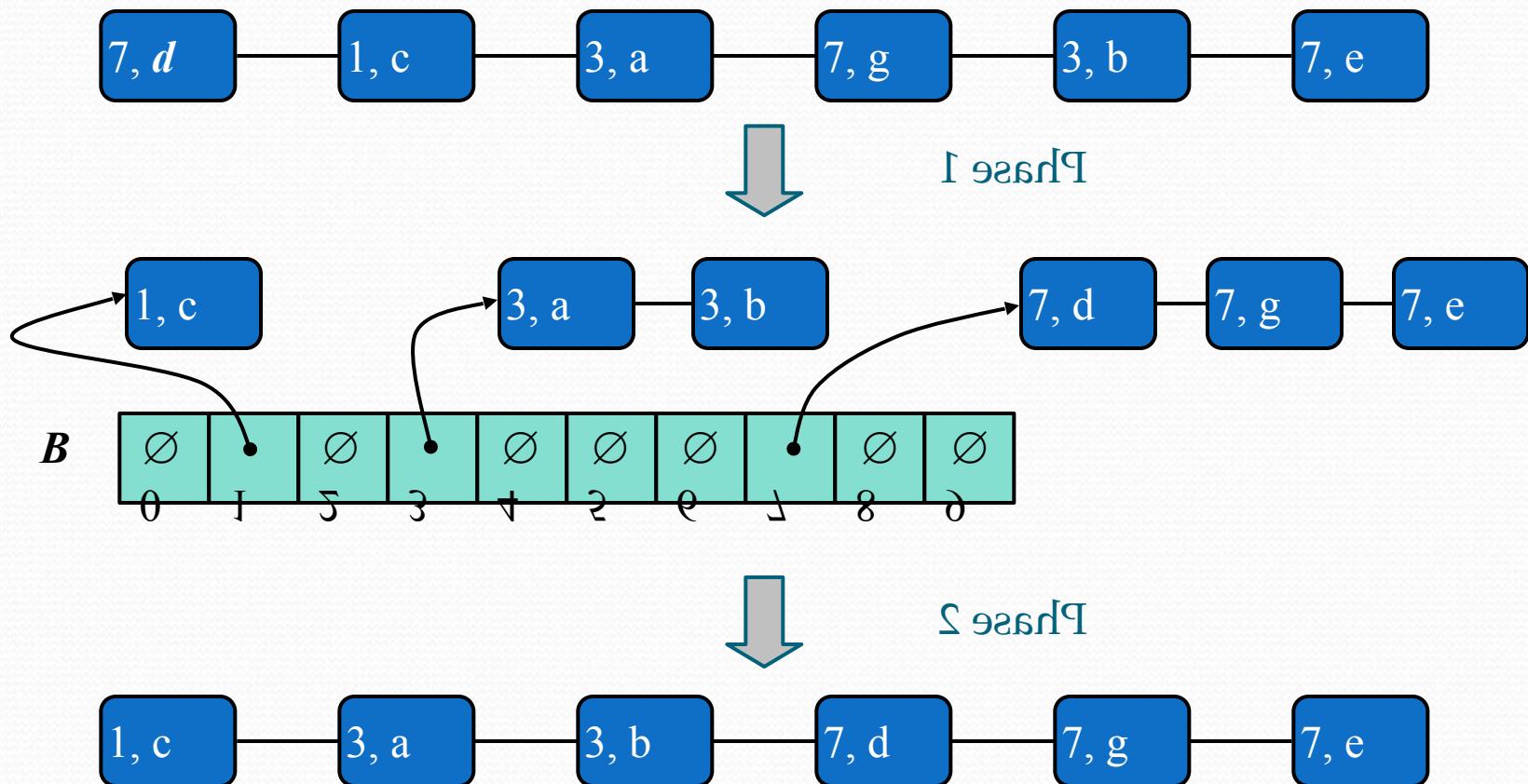
$(k, o) \leftarrow B[i].remove(f)$

$S.insertLast((k, o))$

Example



Key range [0, 9]

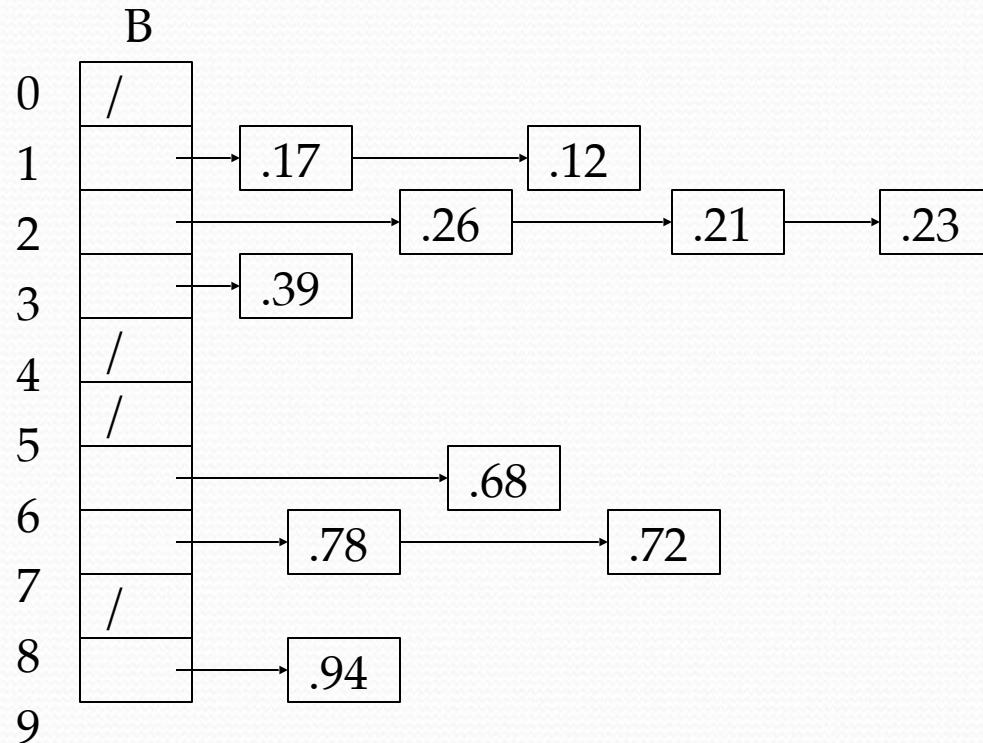


Bucket Sort Illustration

- Bucket sort runs in linear time when the input is drawn from a uniform distribution over the interval $[0, 1)$
- Basic idea
 - Divide $[0,1)$ into n equal-sized subintervals (buckets)
 - Distribute the n numbers into the buckets
 - Sort the numbers in each bucket
 - Go through the buckets in order, list the elements in each

Bucket Sort Illustration

A
1 .78
2 .17
3 .39
4 .26
5 .72
6 .94
7 .21
8 .12
9 .23
10 .68



Bucket Sort Algorithm

BUCKET-SORT(A)

- 1 $n \leftarrow \text{length}[A]$
- 2 **for** $i \leftarrow 1$ **to** n
3 **do** insert $A[i]$ into list $B[\lfloor nA[i] \rfloor]$
- 4 **for** $i \leftarrow 0$ **to** $n - 1$
5 **do** sort list $B[i]$ with insertion sort
- 6 concatenate the lists $B[0], B[1], \dots, B[n - 1]$ together in order

$$T(n) = \Theta(n) + \sum_{i=0}^{n-1} O(n_i^2).$$

Bucket Sort Algorithm – Average Case Analysis

$$T(n) = \Theta(n) + \sum_{i=0}^{n-1} O(n_i^2).$$

Taking expectations both sides and using linearity of expressions

$$\begin{aligned} E[T(n)] &= E\left[\Theta(n) + \sum_{i=0}^{n-1} O(n_i^2)\right] \\ &= \Theta(n) + \sum_{i=0}^{n-1} E[O(n_i^2)] \quad (\text{by linearity of expectation}) \\ &= \Theta(n) + \sum_{i=0}^{n-1} O(E[n_i^2]) \end{aligned}$$

Define an indicator random variable, $X_{i,j}$ as follows:

$X_{i,j} = 1$ if $A[j]$ falls in bucket i

$= 0$ otherwise

$$n_i = \sum_{j=1}^n X_{ij}.$$

Bucket Sort Algorithm – Average Case Analysis

To compute $E[n_i^2]$, we expand the square and regroup terms:

$$\begin{aligned} E[n_i^2] &= E\left[\left(\sum_{j=1}^n X_{ij}\right)^2\right] \\ &= E\left[\sum_{j=1}^n \sum_{k=1}^n X_{ij}X_{ik}\right] \\ &= E\left[\sum_{j=1}^n X_{ij}^2 + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} X_{ij}X_{ik}\right] \\ &= \sum_{j=1}^n E[X_{ij}^2] + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} E[X_{ij}X_{ik}], \end{aligned}$$

Bucket Sort Algorithm – Average Case Analysis

$$\begin{aligned}\mathbb{E}[X_{ij}^2] &= 1 \cdot \frac{1}{n} + 0 \cdot \left(1 - \frac{1}{n}\right) \\ &= \frac{1}{n}.\end{aligned}$$

When $k \neq j$, the variables X_{ij} and X_{ik} are independent, and hence

$$\begin{aligned}\mathbb{E}[X_{ij} X_{ik}] &= \mathbb{E}[X_{ij}] \mathbb{E}[X_{ik}] \\ &= \frac{1}{n} \cdot \frac{1}{n} \\ &= \frac{1}{n^2}.\end{aligned}$$

Bucket Sort Algorithm – Average Case Analysis

$$\mathbb{E}[n_i^2] = \sum_{j=1}^n \mathbb{E}[X_{ij}^2] + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} \mathbb{E}[X_{ij} X_{ik}] ,$$

$$\begin{aligned}\mathbb{E}[n_i^2] &= \sum_{j=1}^n \frac{1}{n} + \sum_{1 \leq j \leq n} \sum_{\substack{1 \leq k \leq n \\ k \neq j}} \frac{1}{n^2} \\ &= n \cdot \frac{1}{n} + n(n-1) \cdot \frac{1}{n^2} \\ &= 1 + \frac{n-1}{n} \\ &= 2 - \frac{1}{n} ,\end{aligned}$$

$$\begin{aligned}\mathbb{E}[T(n)] &= \mathbb{E}\left[\Theta(n) + \sum_{i=0}^{n-1} O(n_i^2)\right] \\ &= \Theta(n) + \sum_{i=0}^{n-1} O(\mathbb{E}[n_i^2]) \\ &= \Theta(n) + O(n) \quad (\text{since } n(2 - 1/n) = 2n - 1 = O(n))\end{aligned}$$

Herman Hollerith (1860-1929)

The 1880 U.S. Census took almost
10 years to process.

While a lecturer at MIT, Hollerith
prototyped punched-card technology.

His machines, including a “card sorter,” allowed the
1890 census total to be reported in 6 weeks.



He founded the Tabulating Machine Company in 1911,
which merged with other companies in 1924 to form
International Business Machines.

Punched cards

- Punched card = data record.
- Hole = value.
- Algorithm = machine + human operator.

1	2	3	4	W	M	0	1	5	6	Un	0	6	12	0	6	12	Ma	Mt	Vi	En	Chi	Wch	Ia	Hi	SD	Ro
5	6	7	8	9	F	10	15	18	19	S	1	7	13	1	7	13	WA	RI	CT	IN	AS	WIS	CO	MO	HO	SP
1	2	3	4	Ch	20	21	25	30	31	MO	2	8	14	2	8	N	NY	NJ	PA	ILL	NIN	ND	IR	KAN	SA	
5	6	7	8	Ap	35	40	45	50	51	MI	3	9	15	3	9	F	MD	VA	WF	WMA	DC	TEN	ALA	CLF	TY	
1	2	3	4	In	55	60	65	70	71	Wd	4	10	16	4	10	21	BL	SC	MS	LB	TEX	DRE	WGN	WI		
5	6	7	8	75	80	85	90	95+	Jn	D	5	11	17+	5	11	PE	GA	FLA	OK	IT	AKS	IDM	MEV	GT		
1	2	3	4	En	OK	0	9	4	17	II	5	Un	15	2	0	LS	Un	En	US	Un	En	UTA	AHI	AF		
5	6	7	8	Ot	NR	1	b	5	01	12	6	NG	20+	3	1	Gr	Ir	Sc	Gr	Ir	Sc	NM	COL	GP		
1	2	3	4	2	NW	4	c	6	0	13	7	I	No	4	Au	Sw	CE	Wa	Sw	CE	Wa	WYO	MNT	PP		
5	6	7	8	4	0	7	d	7	1	14	8	2	Pa	5	Sz	Nw	CF	Hu	Nw	CH	HJ	ALK	AB	UP		
1	2	3	4	6	12	10	e	8	2	15	9	3	Al	6	Po	DK	Fr	It	Dk	Fy	Ir	Au	SEA			
5	6	7	8	8+	Un	g	f	9	3	16	10	4	Un	0	01	Ru	Bu	01	Ru	Bu	Sz	Po	NS			

Replica of
punch card
from the 1900
U.S. census.
**[Howells
2000]**

Hollerith's tabulating system

- Pantograph card punch
- Hand-press reader
- Dial counters
- Sorting box

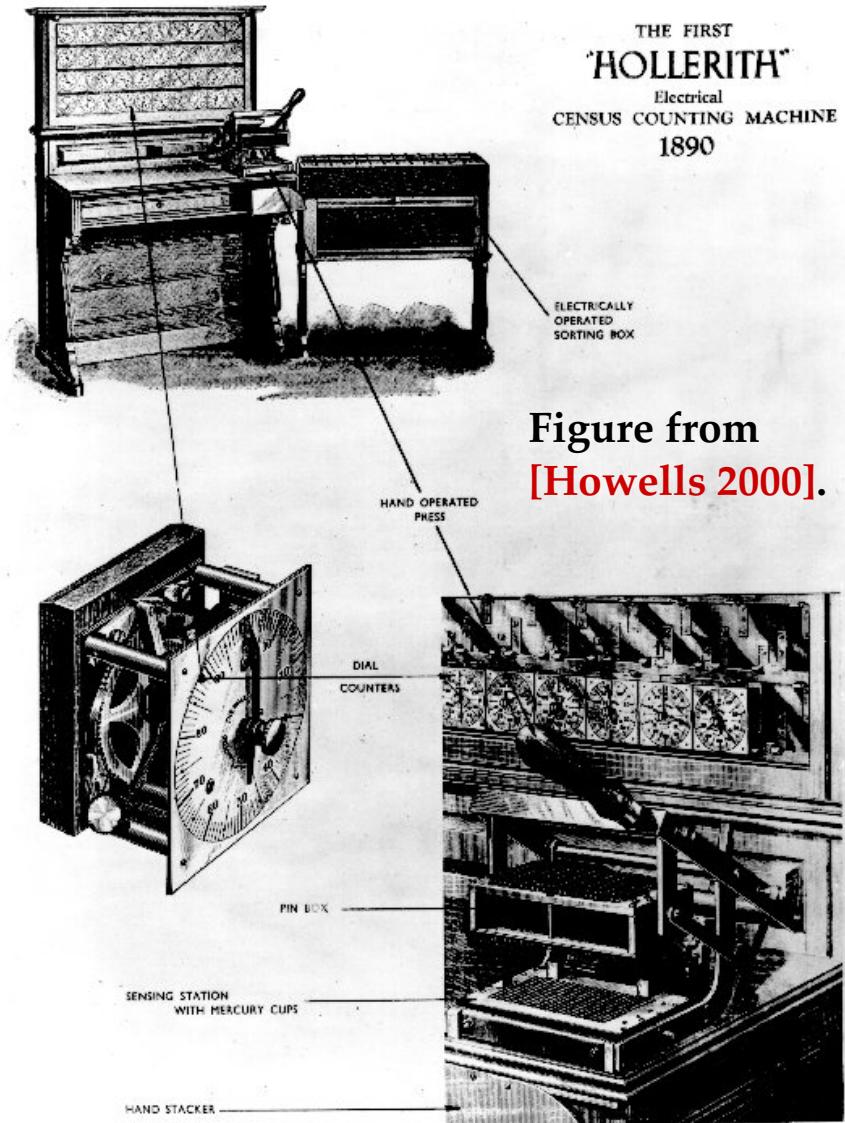


Figure from
[Howells 2000].

Origin of radix sort

Hollerith's original 1889 patent alludes to a most-significant-digit-first radix sort:

"The most complicated combinations can readily be counted with comparatively few counters or relays by first assorting the cards according to the first items entering into the combinations, then reassorting each group according to the second item entering into the combination, and so on, and finally counting on a few counters the last item of the combination for each group of cards."

Least-significant-digit-first radix sort seems to be a folk invention originated by machine operators.

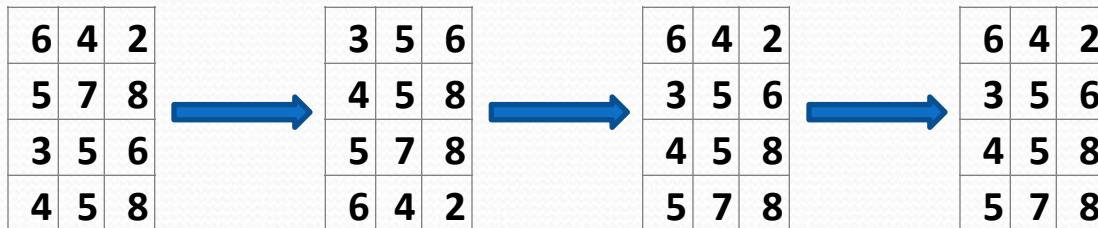
Most-significant-digit-first radix sort will not work.

Example: 642

Origin of radix sort

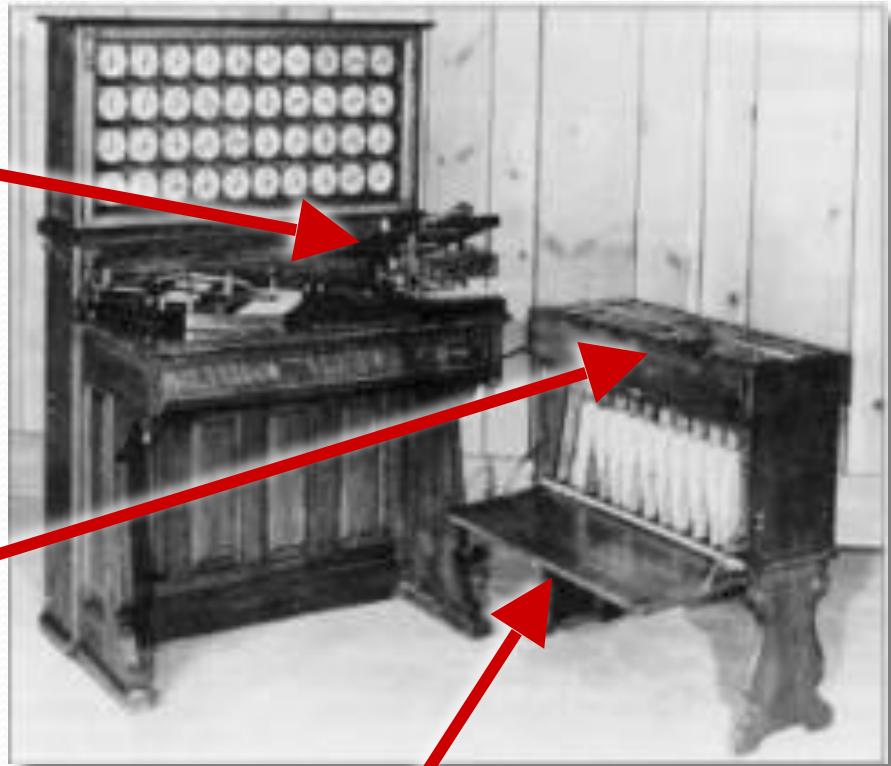
Most-significant-digit-first radix sort will not work.

Example: 642, 578, 356, 458



Operation of the sorter

- An operator inserts a card into the press.
- Pins on the press reach through the punched holes to make electrical contact with mercury-filled cups beneath the card.
- Whenever a particular digit value is punched, the lid of the corresponding sorting bin lifts.
- The operator deposits the card into the bin and closes the lid.
- When all cards have been processed, the front panel is opened, and the cards are collected in order, yielding one pass of a stable sort.



Hollerith Tabulator, Pantograph, Press, and Sorter

Radix Sort

- Intuitively, you might sort on the most significant digit, then the second most significant, etc.
- Problem: lots of intermediate piles of cards to keep track of
- Key idea: sort the *least* significant digit first

```
RadixSort(A, d)
    for i=1 to d
        StableSort(A) on digit i
```

Operation of radix sort

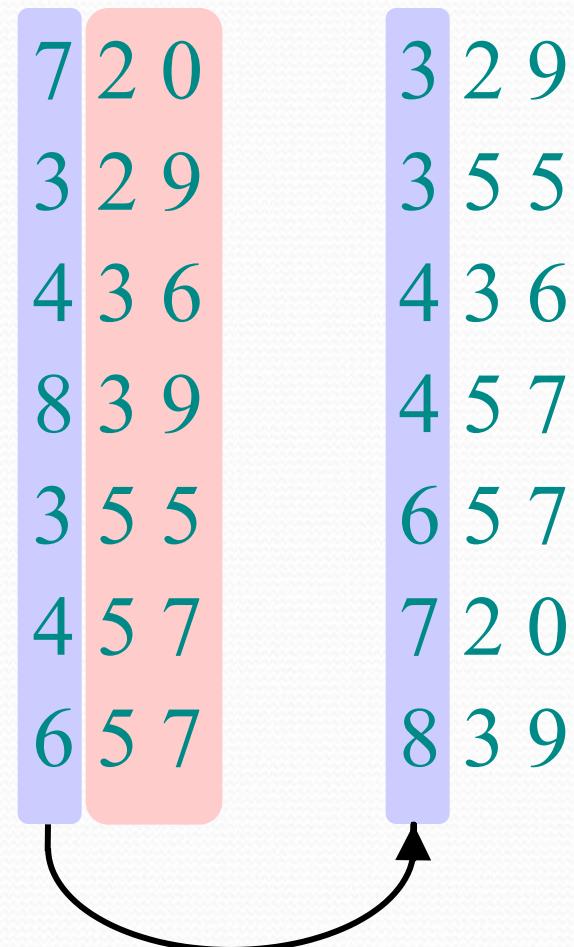
3 2 9	7 2 0	7 2 0	3 2 9
4 5 7	3 5 5	3 2 9	3 5 5
6 5 7	4 3 6	4 3 6	4 3 6
8 3 9	4 5 7	8 3 9	4 5 7
4 3 6	6 5 7	3 5 5	6 5 7
7 2 0	3 2 9	4 5 7	7 2 0
3 5 5	8 3 9	6 5 7	8 3 9



Correctness of radix sort

Induction on digit position

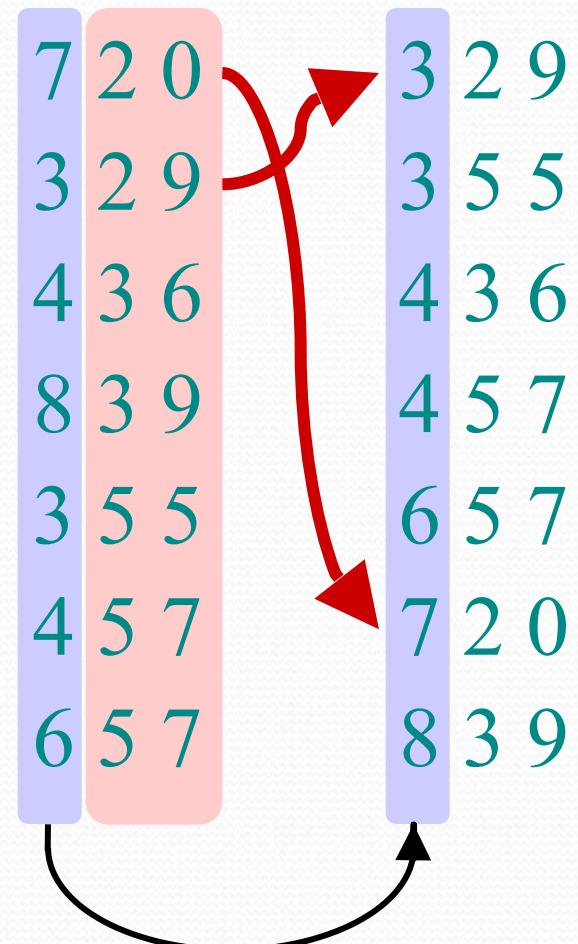
- Assume that the numbers are sorted by their low-order $t - 1$ digits.
- Sort on digit t



Correctness of radix sort

Induction on digit position

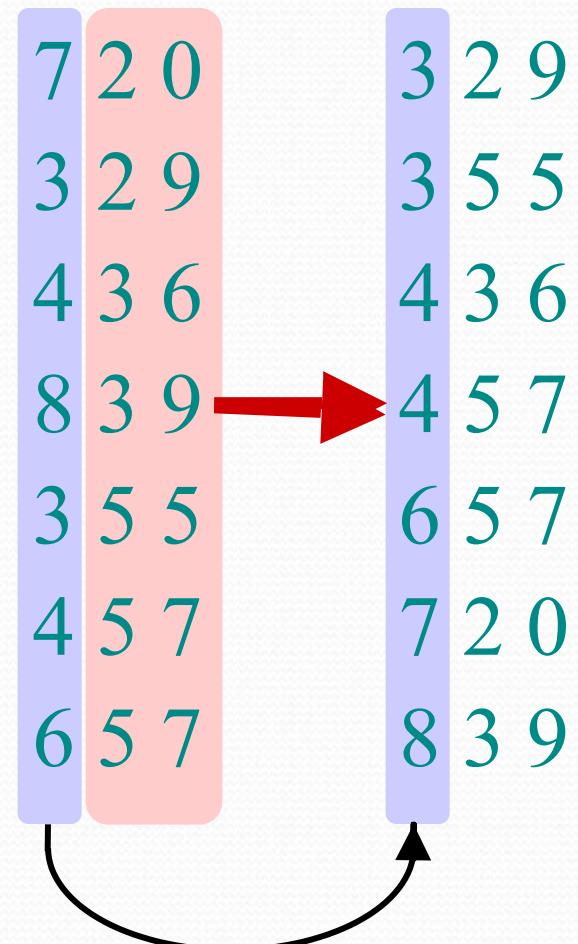
- Assume that the numbers are sorted by their low-order $t - 1$ digits.
- Sort on digit t
 - Two numbers that differ in digit t are correctly sorted.



Correctness of radix sort

Induction on digit position

- Assume that the numbers are sorted by their low-order $t - 1$ digits.
- Sort on digit t
 - Two numbers that differ in digit t are correctly sorted.
 - Two numbers equal in digit t are put in the same order as the input \Rightarrow correct order.



Radix Sort

- *What sort will we use to sort on digits?*
- Bucket sort is a good choice:
 - Sort n numbers on digits that range from 1.. N
 - Time: $O(n + N)$
- Each pass over n numbers with d digits takes time $O(n+k)$, so total time $O(dn+dk)$
 - When d is constant and $k=O(n)$, takes $O(n)$ time

Summary: Radix Sort

- Radix sort:
 - Assumption: input has d digits ranging from 0 to k
 - Basic idea:
 - Sort elements by digit starting with *least* significant
 - Use a stable sort (like bucket sort) for each stage
 - Each pass over n numbers with 1 digit takes time $O(n+k)$, so total time $O(dn+dk)$
 - When d is constant and $k=O(n)$, takes $O(n)$ time
 - Fast, Stable, Simple



Thank You!!