

INTRODUCTION TO ALGORITHMS

Lecture 8: Counting Sort Algorithm

Yao-Chung Fan
yfan@nchu.edu.tw

COUNTING SORT

- ▶ *key-indexed counting sort*
- ▶ *LSD radix sort*
- ▶ *MSD radix sort*

Review: summary of the performance of sorting algorithms

Frequency of operations.

| algorithm | guarantee | random | stable? | operations on keys |
|----------------|-------------------|-------------------|---------|--------------------|
| insertion sort | $\frac{1}{2} N^2$ | $\frac{1}{4} N^2$ | ✓ | compareTo() |
| mergesort | $N \lg N$ | $N \lg N$ | ✓ | compareTo() |
| quicksort | $1.39 N \lg N^*$ | $1.39 N \lg N$ | | compareTo() |
| heapsort | $2 N \lg N$ | $2 N \lg N$ | | compareTo() |

* probabilistic

Lower bound. $\sim N \lg N$ compares required by any compare-based algorithm.

Q. Can we do better (despite the lower bound)?

A. Yes, if we don't depend on key compares. ←

use array accesses
to make R-way decisions
(instead of binary decisions)

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.

| 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*

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.

$R = 6$

```
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++)
    a[i] = aux[i];
```

| 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 |

use a for 0
b for 1
c for 2
d for 3
e for 4
f for 5

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++)
    a[i] = aux[i];
```

copy
back



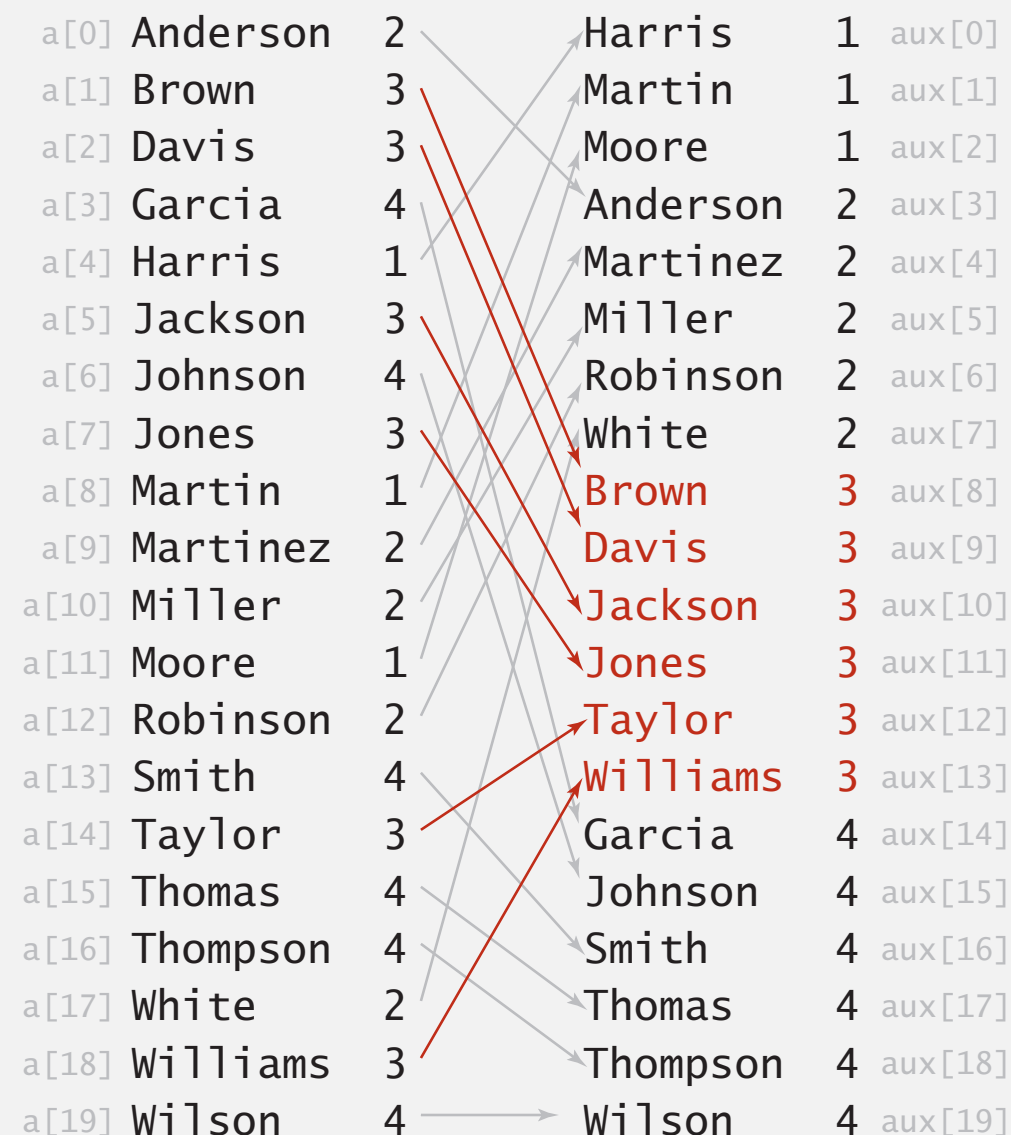
| i | a[i] | | r | count[r] | i | aux[i] |
|----|------|--|---|----------|----|--------|
| 0 | a | | | | 0 | a |
| 1 | a | | | | 1 | a |
| 2 | b | | | | 2 | b |
| 3 | b | | | | 3 | b |
| 4 | b | | a | 2 | 4 | b |
| 5 | c | | b | 5 | 5 | c |
| 6 | d | | c | 6 | 6 | d |
| 7 | d | | d | 8 | 7 | d |
| 8 | e | | e | 9 | 8 | e |
| 9 | f | | f | 12 | 9 | f |
| 10 | f | | - | 12 | 10 | f |
| 11 | f | | | | 11 | f |

Key-indexed counting: analysis

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

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

Stable? ✓



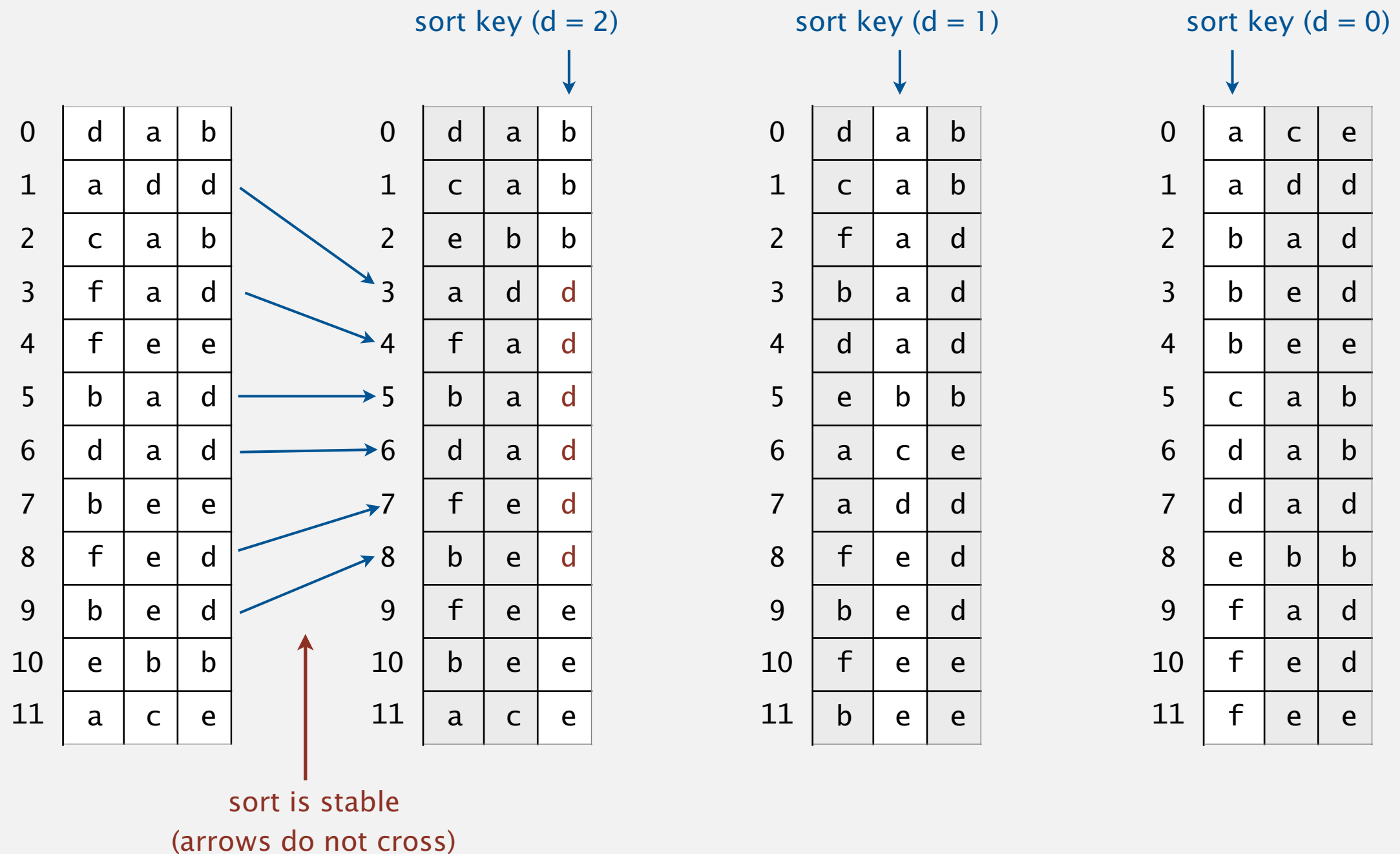
COUNTING SORT

- ▶ *key-indexed counting sort*
- ▶ *LSD radix sort*
- ▶ *MSD radix sort*

Least-significant-digit-first sort

LSD string (radix) sort.

- Consider characters from right to left.
- Stably sort using d^{th} character as the key (using key-indexed counting).



LSD string sort: correctness proof

Proposition. LSD sorts fixed-length strings in ascending order.

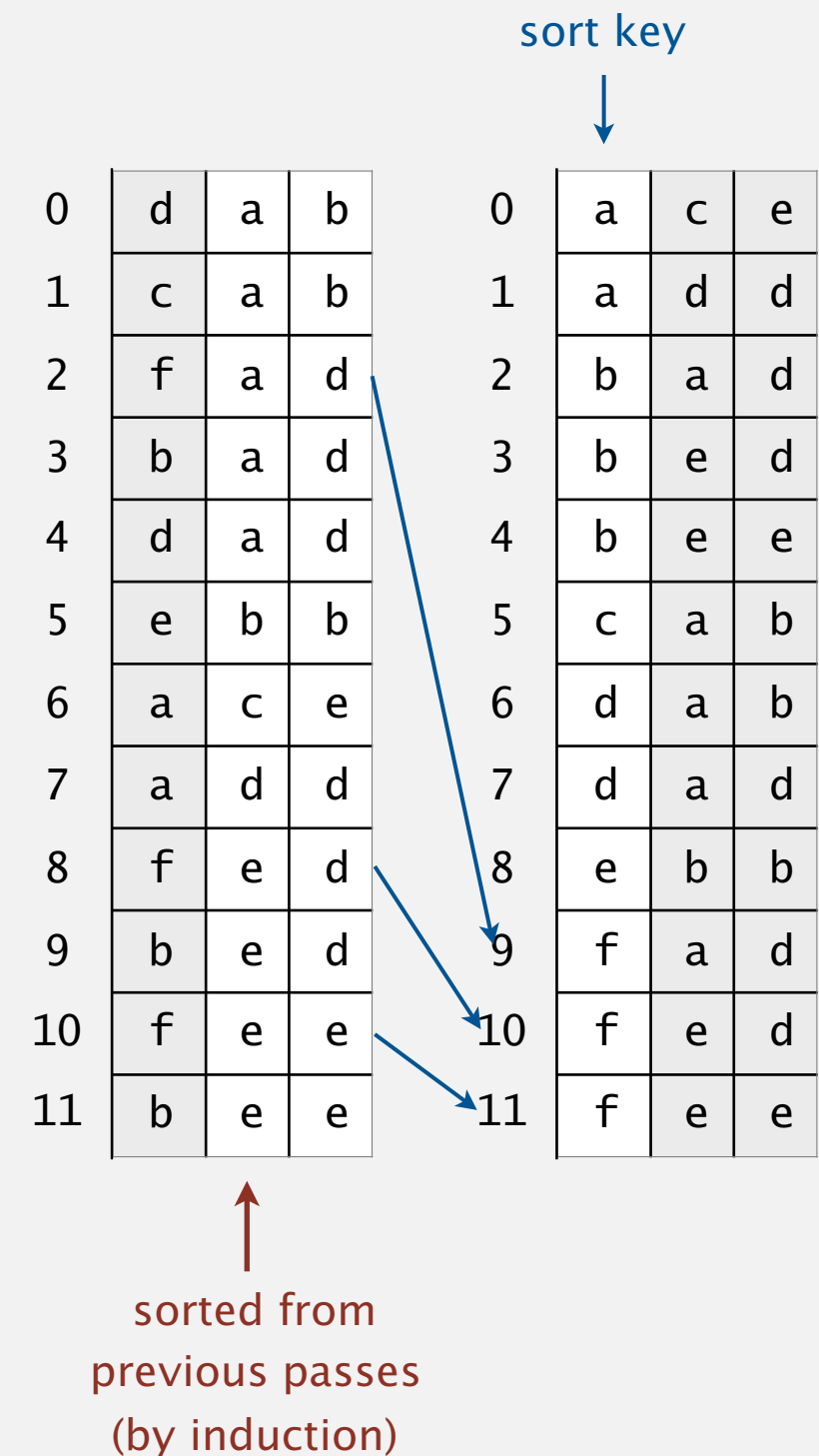
Pf. [by induction on i]

After pass i , strings are sorted by last i characters.

- If two strings differ on sort key, key-indexed sort puts them in proper relative order.
- If two strings agree on sort key, stability keeps them in proper relative order.

Proposition. LSD sort is stable.

Pf. Key-indexed counting is stable.



LSD string sort: Java implementation

```
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

Summary of the performance of sorting algorithms

Frequency of operations.

| algorithm | guarantee | random | stable? |
|----------------|-------------------|-------------------|---------|
| insertion sort | $\frac{1}{2} N^2$ | $\frac{1}{4} N^2$ | ✓ |
| mergesort | $N \lg N$ | $N \lg N$ | ✓ |
| quicksort | $1.39 N \lg N^*$ | $1.39 N \lg N$ | |
| heapsort | $2 N \lg N$ | $2 N \lg N$ | |
| LSD sort † | $2 W (N + R)$ | $2 W (N + R)$ | ✓ |

Q. What if strings are not all of same length?

String sorting challenge 1

Problem. Sort a huge commercial database on a fixed-length key.

Ex. Account number, date, Social Security number, ...

Which sorting method to use?

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



256 (or 65,536) counters;
Fixed-length strings sort in W passes.

| | | | |
|--|-------------|--|--|
| | B14-99-8765 | | |
| | 756-12-AD46 | | |
| | CX6-92-0112 | | |
| | 332-WX-9877 | | |
| | 375-99-QWAX | | |
| | CV2-59-0221 | | |
| | 287-SS-0321 | | |
| | KJ-01-12388 | | |
| | 715-YT-013C | | |
| | MJ0-PP-983F | | |
| | 908-KK-33TY | | |
| | BBN-63-23RE | | |
| | 48G-BM-912D | | |
| | 982-ER-9P1B | | |
| | WBL-37-PB81 | | |
| | 810-F4-J87Q | | |
| | LE9-N8-XX76 | | |
| | 908-KK-33TY | | |
| | B14-99-8765 | | |
| | CX6-92-0112 | | |
| | CV2-59-0221 | | |
| | 332-WX-23SQ | | |
| | 332-6A-9877 | | |

String sorting interview question

Problem. Sort one million 32-bit string.

Ex. Google (or presidential) interview.

Which sorting method to use?

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



Reverse LSD

- Consider characters from left to right.
- Stably sort using d^{th} character as the key (using key-indexed counting).

| | | sort key (d = 0) | | sort key (d = 1) | | sort key (d = 2) | | | | | |
|----|---|------------------|---|------------------|---|------------------|---|----|---|---|---|
| | | ↓ | | ↓ | | ↓ | | | | | |
| 0 | d | a | b | 0 | a | d | d | 0 | c | a | b |
| 1 | a | d | d | 1 | a | c | e | 1 | d | a | b |
| 2 | c | a | b | 2 | b | a | d | 2 | e | b | b |
| 3 | f | a | d | 3 | b | e | e | 3 | b | a | d |
| 4 | f | e | e | 4 | b | e | d | 4 | d | a | d |
| 5 | b | a | d | 5 | c | a | b | 5 | f | a | d |
| 6 | d | a | d | 6 | d | a | b | 6 | a | d | d |
| 7 | b | e | e | 7 | d | a | d | 7 | b | e | d |
| 8 | f | e | d | 8 | e | b | b | 8 | f | e | d |
| 9 | b | e | d | 9 | f | a | d | 9 | a | c | e |
| 10 | e | b | b | 10 | f | e | e | 10 | b | e | e |
| 11 | a | c | e | 11 | f | e | d | 11 | f | e | e |

not sorted!

Most-significant-digit-first string sort

MSD (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).

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

↑
sort key

count[]

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

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

sort subarrays
recursively

MSD string sort: top-level trace

| | use key-indexed counting on first character | | | recursively sort subarrays | | |
|----|---|-----------------------------|--------------------------|---|------------------|-------------|
| | count frequencies | transform counts to indices | distribute and copy back | indices at completion of distribute phase | | |
| 0 | she | 0 0 | 0 0 | 0 0 0 | sort(a, 0, 0); | 0 are |
| 1 | sells | 1 a 0 | 1 a 0 | 1 a 1 | sort(a, 1, 1); | 1 by |
| 2 | seashells | 2 b 1 | 2 b 1 | 2 b 2 | sort(a, 2, 1); | 2 sea |
| 3 | by | 3 c 1 | 3 c 2 | 3 c 2 | sort(a, 2, 1); | 3 seashells |
| 4 | the | 4 d 0 | 4 d 2 | 4 d 2 | sort(a, 2, 1); | 4 seashells |
| 5 | sea | 5 e 0 | 5 e 2 | 5 e 2 | sort(a, 2, 1); | 5 sells |
| 6 | shore | 6 f 0 | 6 f 2 | 6 f 2 | sort(a, 2, 1); | 6 sells |
| 7 | the | 7 g 0 | 7 g 2 | 7 g 2 | sort(a, 2, 1); | 7 she |
| 8 | shells | 8 h 0 | 8 h 2 | 8 h 2 | sort(a, 2, 1); | 8 she |
| 9 | she | 9 i 0 | 9 i 2 | 9 i 2 | sort(a, 2, 1); | 9 shells |
| 10 | sells | 10 j 0 | 10 j 2 | 10 j 2 | sort(a, 2, 1); | 10 shore |
| 11 | are | 11 k 0 | 11 k 2 | 11 k 2 | sort(a, 2, 1); | 11 surely |
| 12 | surely | 12 l 0 | 12 l 2 | 12 l 2 | sort(a, 2, 1); | 12 the |
| 13 | seashells | 13 m 0 | 13 m 2 | 13 m 2 | sort(a, 2, 1); | 13 the |
| | | 14 n 0 | 14 n 2 | 14 n 2 | sort(a, 2, 1); | |
| | | 15 o 0 | 15 o 2 | 15 o 2 | sort(a, 2, 1); | |
| | | 16 p 0 | 16 p 2 | 16 p 2 | sort(a, 2, 1); | |
| | | 17 q 0 | 17 q 2 | 17 q 2 | sort(a, 2, 1); | |
| | | 18 r 0 | 18 r 2 | 18 r 2 | sort(a, 2, 11); | |
| | | 19 s 0 | 19 s 2 | 19 s 12 | sort(a, 12, 13); | |
| | | 20 t 10 | 20 t 12 | 20 t 14 | sort(a, 14, 13); | |
| | | 21 u 2 | 21 u 14 | 21 u 14 | sort(a, 14, 13); | |
| | | 22 v 0 | 22 v 14 | 22 v 14 | sort(a, 14, 13); | |
| | | 23 w 0 | 23 w 14 | 23 w 14 | sort(a, 14, 13); | |
| | | 24 x 0 | 24 x 14 | 24 x 14 | sort(a, 14, 13); | |
| | | 25 y 0 | 25 y 14 | 25 y 14 | sort(a, 14, 13); | |
| | | 26 z 0 | 26 z 14 | 26 z 14 | sort(a, 14, 13); | |
| | | 27 0 | 27 14 | 27 14 | sort(a, 14, 13); | |

start of s subarray
1 + end of s subarray

MSD string sort: example

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

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

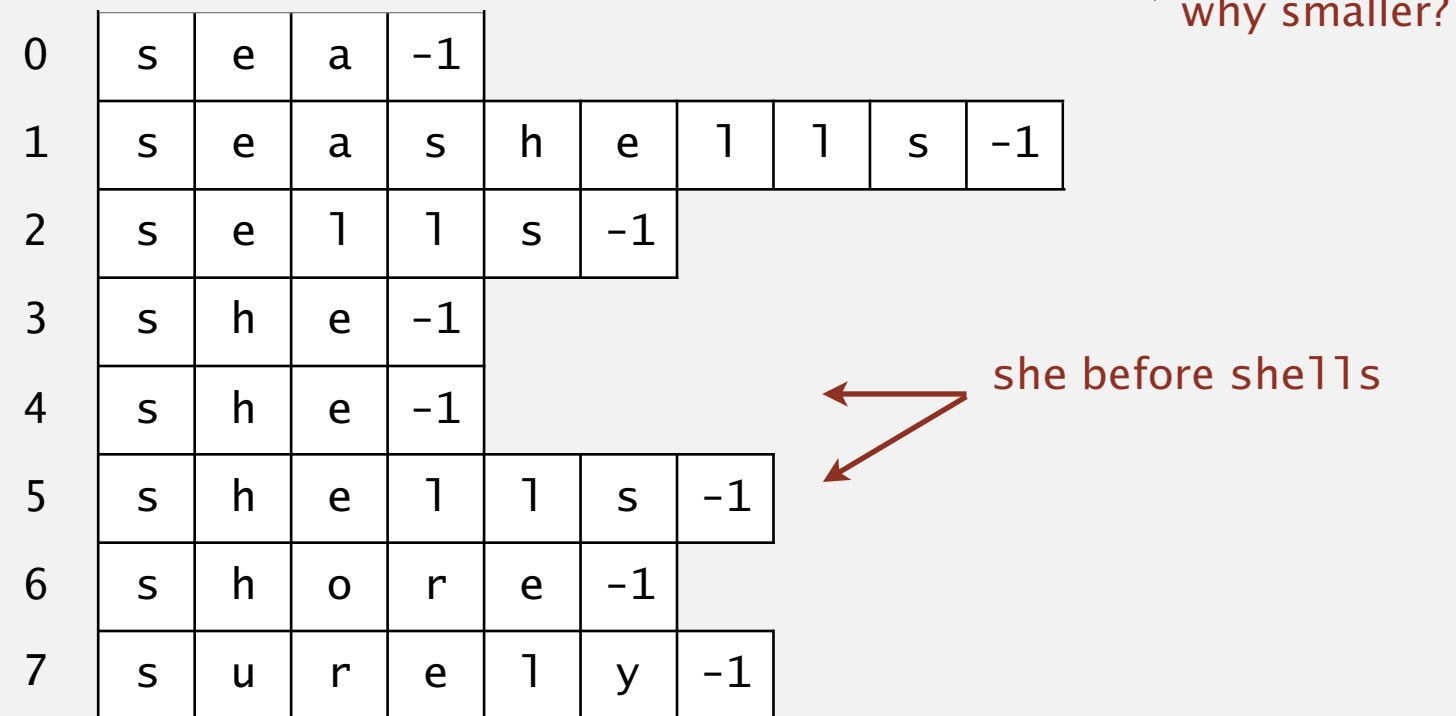
need to examine
every character
in equal keys

end of string
goes before any
char value

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

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 | | | |

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

MSD string sort: Java implementation

```
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++)
        sort(a, aux, lo + count[r], lo + count[r+1] - 1, d+1);
}
```

recycles aux[] array
but not count[] array

key-indexed counting

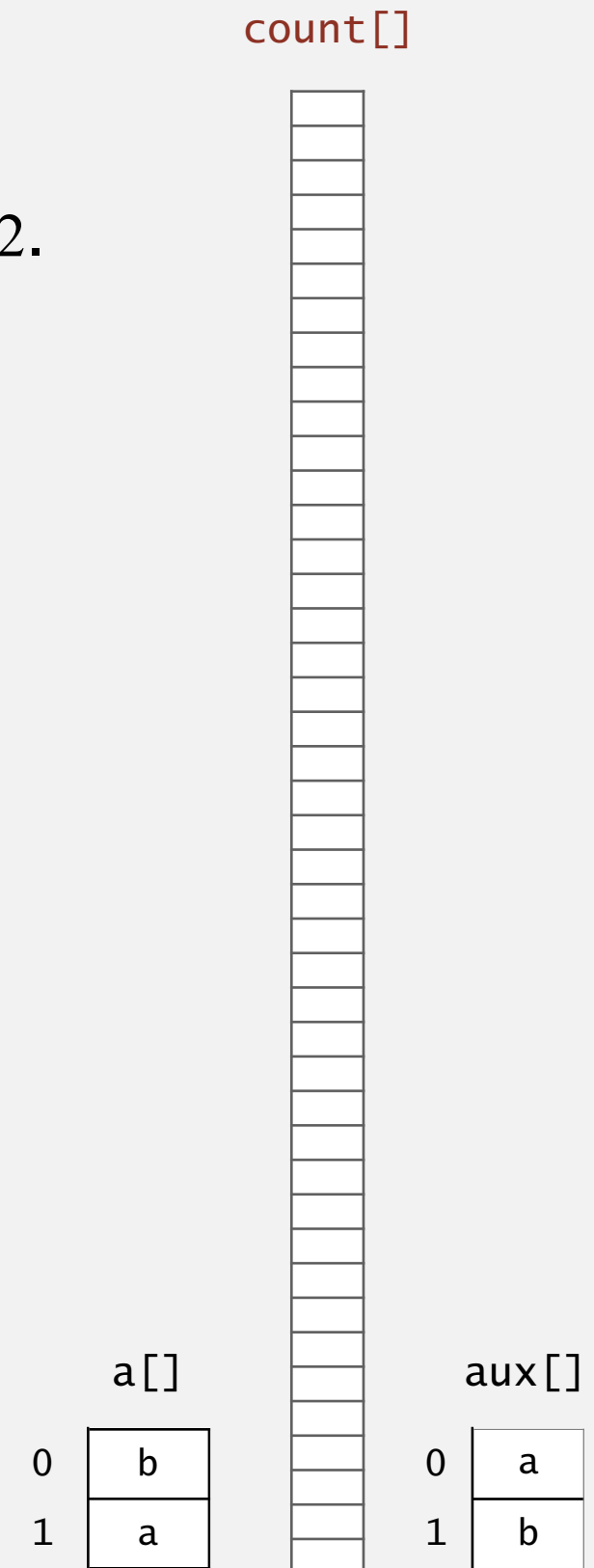
sort R subarrays recursively

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.



MSD string sort: performance

Number of characters examined.

- MSD examines just enough characters to sort the keys.
- Number of characters examined depends on keys.

| | Random (sublinear) | Non-random with duplicates (nearly linear) | Worst case (linear) |
|------------|-----------------------|--|------------------------|
| 111 | 1EIO402 | are | 1DNB377 |
| | 1HYL490 | by | 1DNB377 |
| 222 | 1ROZ572 | sea | 1DNB377 |
| | 2HXE734 | seashells | 1DNB377 |
| 333 | 2IYE230 | seashells | 1DNB377 |
| | 2XOR846 | sells | 1DNB377 |
| 444 | 3CDB573 | sells | 1DNB377 |
| | 3CVP720 | she | 1DNB377 |
| | 3IGJ319 | she | 1DNB377 |
| | 3KNA382 | shells | 1DNB377 |
| | 3TAV879 | shore | 1DNB377 |
| | 4CQP781 | surely | 1DNB377 |
| | 4QGI284 | the | 1DNB377 |
| | 4YHV229 | the | 1DNB377 |

Characters examined by MSD string sort

Summary of the performance of sorting algorithms

Frequency of operations.

| algorithm | guarantee | random | extra space | stable? |
|----------------|-------------------|-------------------|-------------|---------|
| insertion sort | $\frac{1}{2} N^2$ | $\frac{1}{4} N^2$ | 1 | ✓ |
| mergesort | $N \lg N$ | $N \lg N$ | N | ✓ |
| quicksort | $1.39 N \lg N^*$ | $1.39 N \lg N$ | $c \lg N$ | |
| heapsort | $2 N \lg N$ | $2 N \lg N$ | 1 | |
| LSD sort † | $2 W (N + R)$ | $2 W (N + R)$ | $N + R$ | ✓ |
| MSD sort ‡ | $2 W (N + R)$ | $N \log_R N$ | $N + D R$ | ✓ |

1024→512→256→128→.....

D = function-call stack depth
(length of longest prefix match)

- * probabilistic
- † fixed-length W keys
- ‡ average-length W keys