# Strings

- Strings Sorts
- Tries
- Substring Search
- Regular Expressions
- Data Compression

# Strings

- Strings Sorts
- Tries
- Substring Search
- Regular Expressions
- Data Compression

# String processing

String. Sequence of characters.

Important fundamental abstraction.

- Information processing.
- Genomic sequences.
- Communication systems (e.g., email).
- Programming systems (e.g., Java programs).
- …

> " *The digital information that underlies biochemistry, cell biology, and development can be represented by a simple string of G's, A's, T's and C's. This string is the root data structure of an organism's biology.* " *— M. V. Olson*

# The char data type

**C char data type.** Typically an 8-bit integer.

- Supports 7-bit ASCII.
- Need more bits to represent certain characters.

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | A | B | C | D | E | F |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | NUL | SOH | STX | ETX | EOT | ENQ | ACK | BEL | BS | HT | LF | VT | FF | CR | SO | SI |
| 1 | DLE | DC1 | DC2 | DC3 | DC4 | NAK | SYN | ETB | CAN | EM | SUB | ESC | FS | GS | RS | US |
| 2 | SP | ! | " | # | $ | % | & | ' | ( | ) | * | + | , | - | . | / |
| 3 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | : | ; | < | = | > | ? |
| 4 | @ | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O |
| 5 | P | Q | R | S | T | U | V | W | X | Y | Z | [ | \ | ] | ^ | _ |
| 6 | ` | a | b | c | d | e | f | g | h | i | j | k | l | m | n | o |
| 7 | p | q | r | s | t | u | v | w | x | y | z | { | | | } | ~ | DEL |

**Hexadecimal to ASCII conversion table**

**Java char data type.** A 16-bit unsigned integer.

- Supports original 16-bit Unicode.
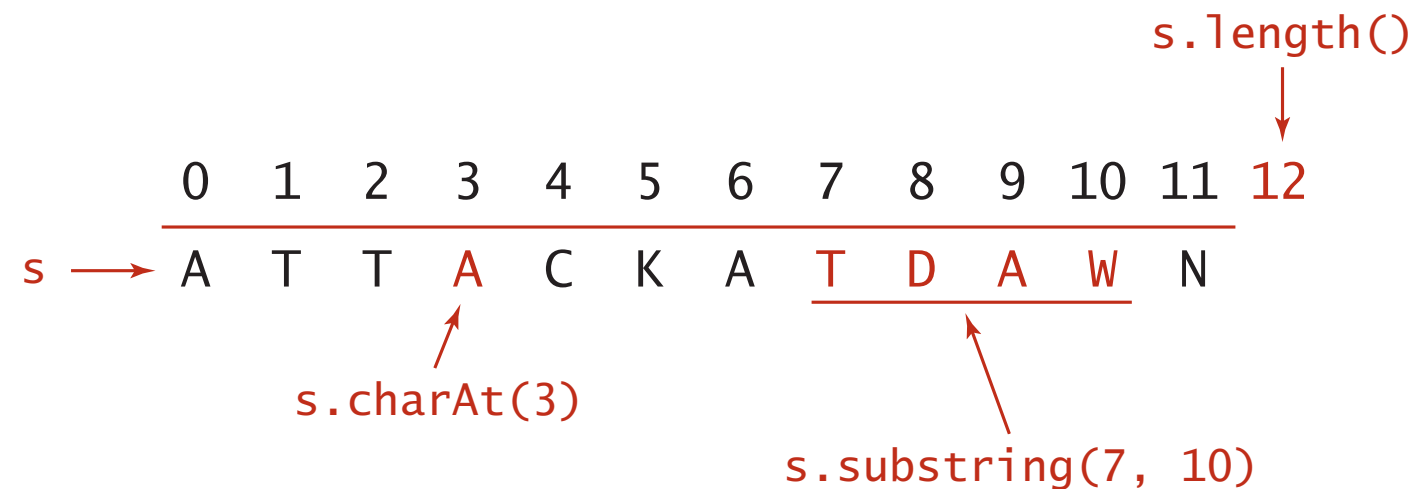- Supports 21-bit Unicode 3.0 (awkwardly).

# The String data type

String data type.  Sequence of characters (immutable).

Indexing.  Get the $i^{th}$ character.

Substring extraction.  Get a contiguous sequence of characters from a string.

String concatenation.  Append one character to end of another string.

```
                                            s.length()
                                                 ↓
      0   1   2   3   4   5   6   7   8   9  10  11  12
     ────────────────────────────────────────────────
s →   A   T   T   A   C   K   A   T   D   A   W   N
                      ↑                  ─────────
              s.charAt(3)
                                    s.substring(7, 10)
```

# The String data type: Java implementation

```java
public final class String implements Comparable<String>
{
   private char[] value;   // characters
   private int offset;     // index of first char in array
   private int count;      // length of string
   private int hash;       // cache of hashCode()

   private String(int offset, int count, char[] value)
   {
      this.offset = offset;
      this.count  = count;
      this.value  = value;
   }

   public String substring(int from, int to)
   {  return new String(offset + from, to - from, value);  }

   public char charAt(int index)
   {  return value[index + offset];  }

   public String concat(String that)
   {
      char[] val = new char[this.length() + that.length());
      ...
      return new String(0, this.length() + that.length(), val);
   }
}
```
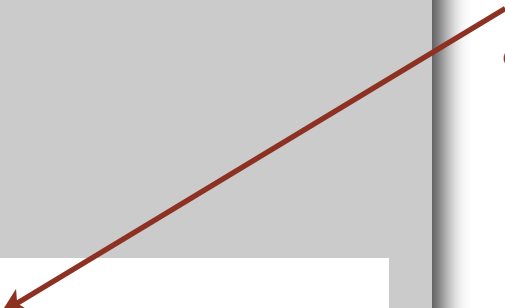
strings share underlying `char[]` array

# The String data type: performance

String data type. Sequence of characters (immutable).

Underlying implementation. Immutable `char[]` array, offset, and length.

| | String | |
| operation | guarantee | extra space |
| --- | --- | --- |
| `charAt()` | 1 | 1 |
| `substring()` | 1 | 1 |
| `concat()` | N | N |

Memory. $40 + 2N$ bytes for a virgin `String` of length $N$.

use `byte[]` or `char[]` instead of `String` to save space

# Alphabets

Digital key. Sequence of digits over fixed alphabet.

Radix. Number of digits $R$ in alphabet.

| name | R() | lgR() | characters |
|---|---|---|---|
| BINARY | 2 | 1 | 01 |
| OCTAL | 8 | 3 | 01234567 |
| DECIMAL | 10 | 4 | 0123456789 |
| HEXADECIMAL | 16 | 4 | 0123456789ABCDEF |
| DNA | 4 | 2 | ACTG |
| LOWERCASE | 26 | 5 | abcdefghijklmnopqrstuvwxyz |
| UPPERCASE | 26 | 5 | ABCDEFGHIJKLMNOPQRSTUVWXYZ |
| PROTEIN | 20 | 5 | ACDEFGHIKLMNPQRSTVWY |
| BASE64 | 64 | 6 | ABCDEFGHIJKLMNOPQRSTUVWXYZabcdefghijklmnopqrstuvwxyz0123456789+/ |
| ASCII | 128 | 7 | *ASCII characters* |
| EXTENDED_ASCII | 256 | 8 | *extended ASCII characters* |
| UNICODE16 | 65536 | 16 | *Unicode characters* |

**Standard alphabets**

# String Sorts

- key-indexed counting
- LSD string sort
- MSD string sort
- 3-way string quicksort
- suffix arrays

# Review: summary of the performance of sorting algorithms

Frequency of operations = key compares.

| algorithm | guarantee | random | extra space | stable? | operations on keys |
|-----------|-----------|--------|-------------|---------|--------------------|
| insertion sort | $N^2/2$ | $N^2/4$ | no | yes | `compareTo()` |
| mergesort | N lg N | N lg N | N | yes | `compareTo()` |
| quicksort | 1.39 N lg N * | 1.39 N lg N | c lg N | no | `compareTo()` |
| heapsort | 2 N lg N | 2 N lg N | no | no | `compareTo()` |

\* probabilistic

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

Q. Can we do better (despite the lower bound)?
A. Yes, if we don't depend on compares.

# String Sorts

- **key-indexed counting**
- LSD string sort
- MSD string sort
- 3-way string quicksort
- suffix arrays

# 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.

**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*

# Key-indexed counting

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

• Count frequencies of each letter using key as index.

•

•

•

```
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];
```

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 |

offset by 1
[stay tuned]

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

13

# Key-indexed counting

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.
- 
- 

```
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];
```

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] |
|---|---|
| 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

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 records.
- 

```
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];
```

move records

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

| i | aux[i] |
|---|--------|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | |
| 7 | |
| 8 | |
| 9 | |
| 10 | |
| 11 | |

# Key-indexed counting

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 records.
- 

```
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];
```

move records

| 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 | 0 |
| b | 2 |
| c | 5 |
| d | 7 |
| e | 8 |
| f | 9 |
| – | 12 |

| i | aux[i] |
|---|--------|
| 0 | |
| 1 | |
| 2 | |
| 3 | |
| 4 | |
| 5 | |
| 6 | d |
| 7 | |
| 8 | |
| 9 | |
| 10 | |
| 11 | |

# Key-indexed counting

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 records.
- 

```
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];
```

move records

| 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 | 1 |
| b | 2 |
| c | 5 |
| d | 7 |
| e | 8 |
| f | 9 |
| – | 12 |

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

# Key-indexed counting

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 records.
-

```
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];
```

move records

| 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 | 1 |
| b | 2 |
| c | 6 |
| d | 7 |
| e | 8 |
| f | 9 |
| – | 12 |

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

18

# Key-indexed counting

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 records.
- 

```
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];
```

move records

| 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 | 1 |
| b | 2 |
| c | 6 |
| d | 7 |
| e | 8 |
| f | 10 |
| – | 12 |

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

# Key-indexed counting

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 records.
-

```
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];
```

move records →

| 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 | 1 |
| b | 2 |
| c | 6 |
| d | 7 |
| e | 8 |
| f | 11 |
| - | 12 |

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

# Key-indexed counting

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 records.
- 

```
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];
```

move
records

| 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 | 1        |
| b | 3        |
| c | 6        |
| d | 7        |
| e | 8        |
| f | 11       |
| – | 12       |

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

21

# Key-indexed counting

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 records.
-

```
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];
```

move
records →

| 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 | 1 |
| b | 3 |
| c | 6 |
| d | 8 |
| e | 8 |
| f | 11 |
| – | 12 |

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

# Key-indexed counting

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 records.
- 

```
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];
```

move records →

| 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 | 1 |
| b | 4 |
| c | 6 |
| d | 8 |
| e | 8 |
| f | 11 |
| - | 12 |

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

# Key-indexed counting

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 records.
- 

```
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];
```

move records

| 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 | 1 |
| b | 4 |
| c | 6 |
| d | 8 |
| e | 8 |
| f | 12 |
| – | 12 |

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

# Key-indexed counting

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 records.
•

```
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];
```

move records

| 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 | 1 |
| b | 5 |
| c | 6 |
| d | 8 |
| e | 8 |
| f | 12 |
| – | 12 |

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

# Key-indexed counting

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 records.
- 

```
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];
```

move records →

| 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 | 1 |
| b | 5 |
| c | 6 |
| d | 8 |
| e | 9 |
| f | 12 |
| – | 12 |

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

# Key-indexed counting

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 records.
-

```
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];
```

move records

| 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 |

# Key-indexed counting

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 records.
- 

```
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];
```

move records

| 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 |

# Key-indexed counting

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 records.
- 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] |
|---|------|
| 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 |

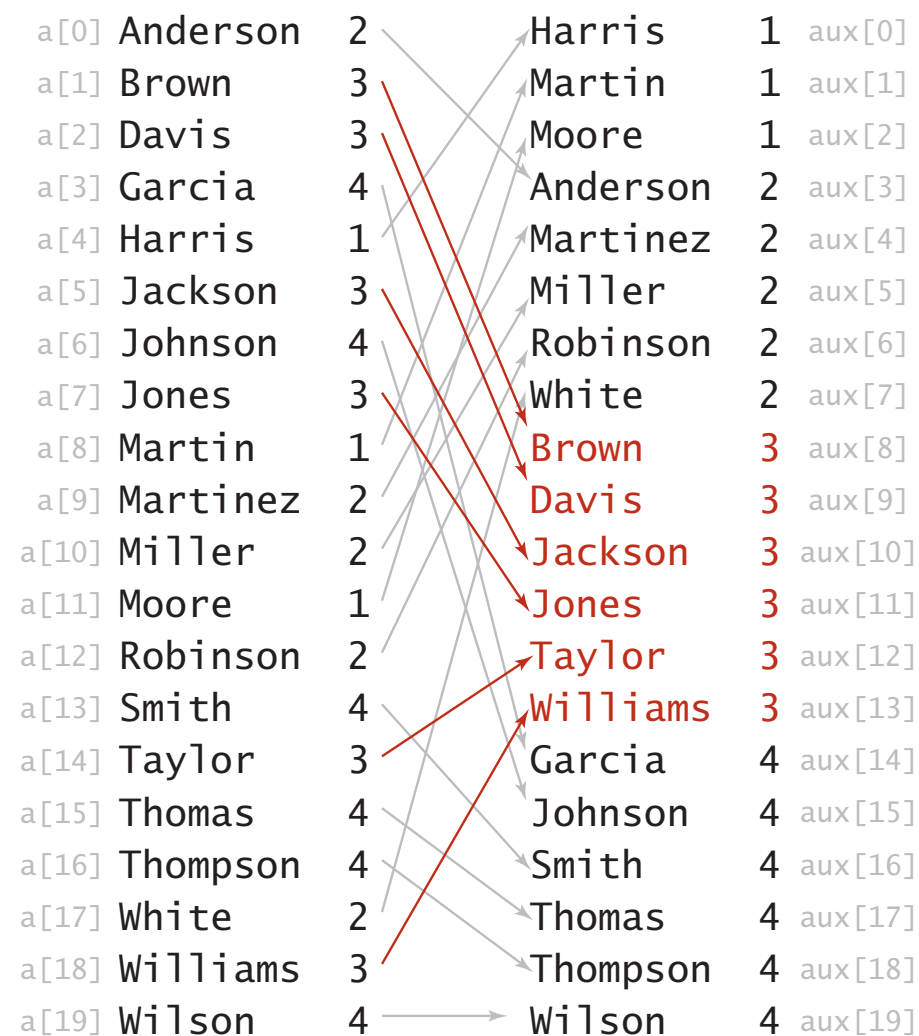| i | aux[i] |
|---|--------|
| 0 | a |
| 1 | a |
| 2 | b |
| 3 | b |
| 4 | b |
| 5 | c |
| 6 | d |
| 7 | d |
| 8 | e |
| 9 | f |
| 10 | f |
| 11 | f |

# Key-indexed counting:  analysis

Proposition.  Key-indexed counting uses $8N + 3R$ array accesses to sort
$N$ records whose keys are integers between $0$ and $R-1$.

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

Stable? Yes!

In-place? No.

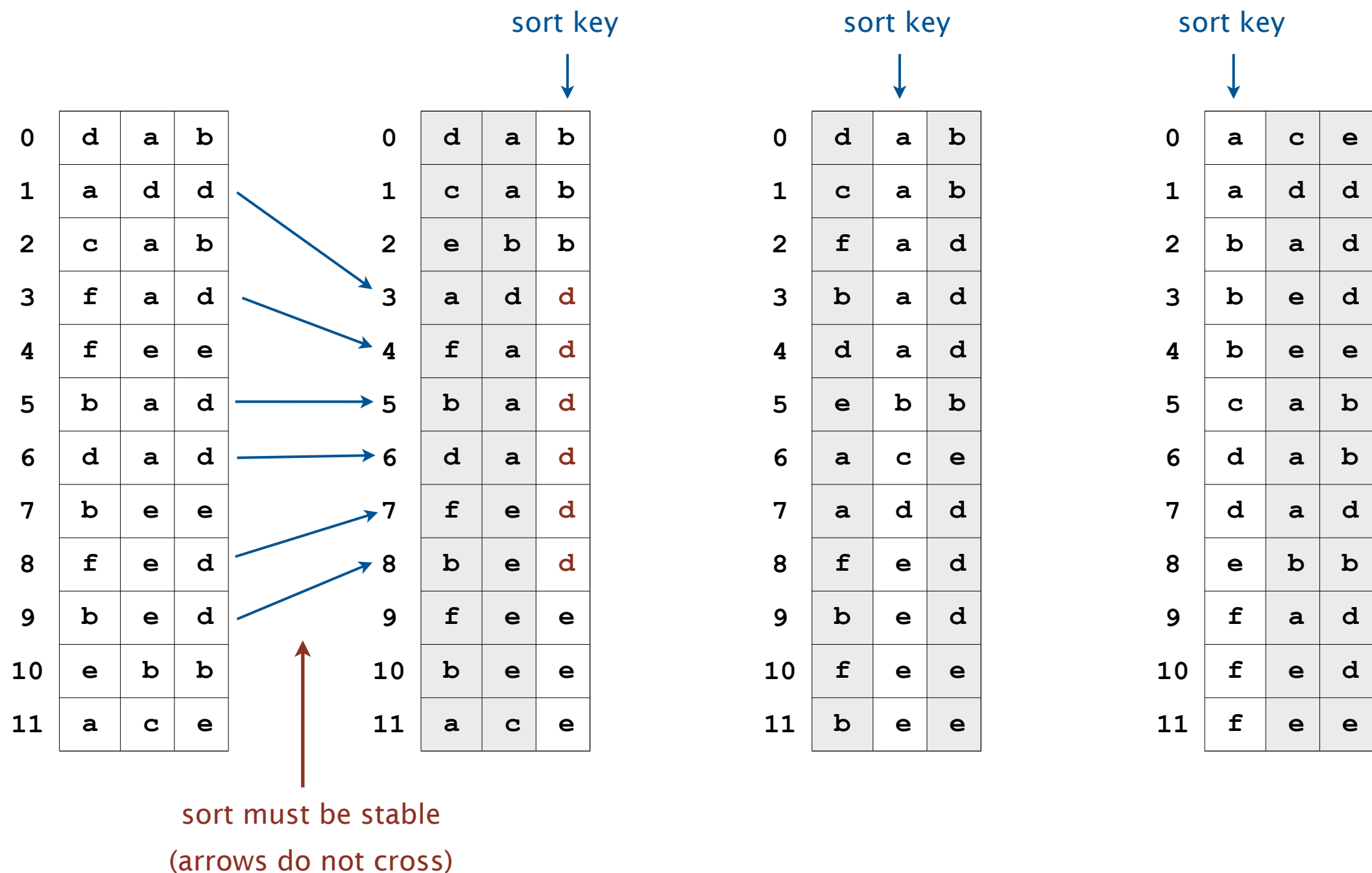| | | | | | |
|---|---|---|---|---|---|
| a[0] | Anderson | 2 | Harris | 1 | aux[0] |
| a[1] | Brown | 3 | Martin | 1 | aux[1] |
| a[2] | Davis | 3 | Moore | 1 | aux[2] |
| a[3] | Garcia | 4 | Anderson | 2 | aux[3] |
| a[4] | Harris | 1 | Martinez | 2 | aux[4] |
| a[5] | Jackson | 3 | Miller | 2 | aux[5] |
| a[6] | Johnson | 4 | Robinson | 2 | aux[6] |
| a[7] | Jones | 3 | White | 2 | aux[7] |
| a[8] | Martin | 1 | Brown | 3 | aux[8] |
| a[9] | Martinez | 2 | Davis | 3 | aux[9] |
| a[10] | Miller | 2 | Jackson | 3 | aux[10] |
| a[11] | Moore | 1 | Jones | 3 | aux[11] |
| a[12] | Robinson | 2 | Taylor | 3 | aux[12] |
| a[13] | Smith | 4 | Williams | 3 | aux[13] |
| a[14] | Taylor | 3 | Garcia | 4 | aux[14] |
| a[15] | Thomas | 4 | Johnson | 4 | aux[15] |
| a[16] | Thompson | 4 | Smith | 4 | aux[16] |
| a[17] | White | 2 | Thomas | 4 | aux[17] |
| a[18] | Williams | 3 | Thompson | 4 | aux[18] |
| a[19] | Wilson | 4 | Wilson | 4 | aux[19] |

# String Sorts

- key-indexed counting
- **LSD string sort**
- MSD string sort
- 3-way string quicksort
- suffix arrays

# Least-significant-digit-first string sort

LSD string sort.

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

|     | sort key |   |   |
| --- | --- | --- | --- |
| 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 |

|     |   |   | sort key |
| --- | --- | --- | --- |
| 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 |

sort must be stable
(arrows do not cross)

|     |   | sort key |   |
| --- | --- | --- | --- |
| 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 |

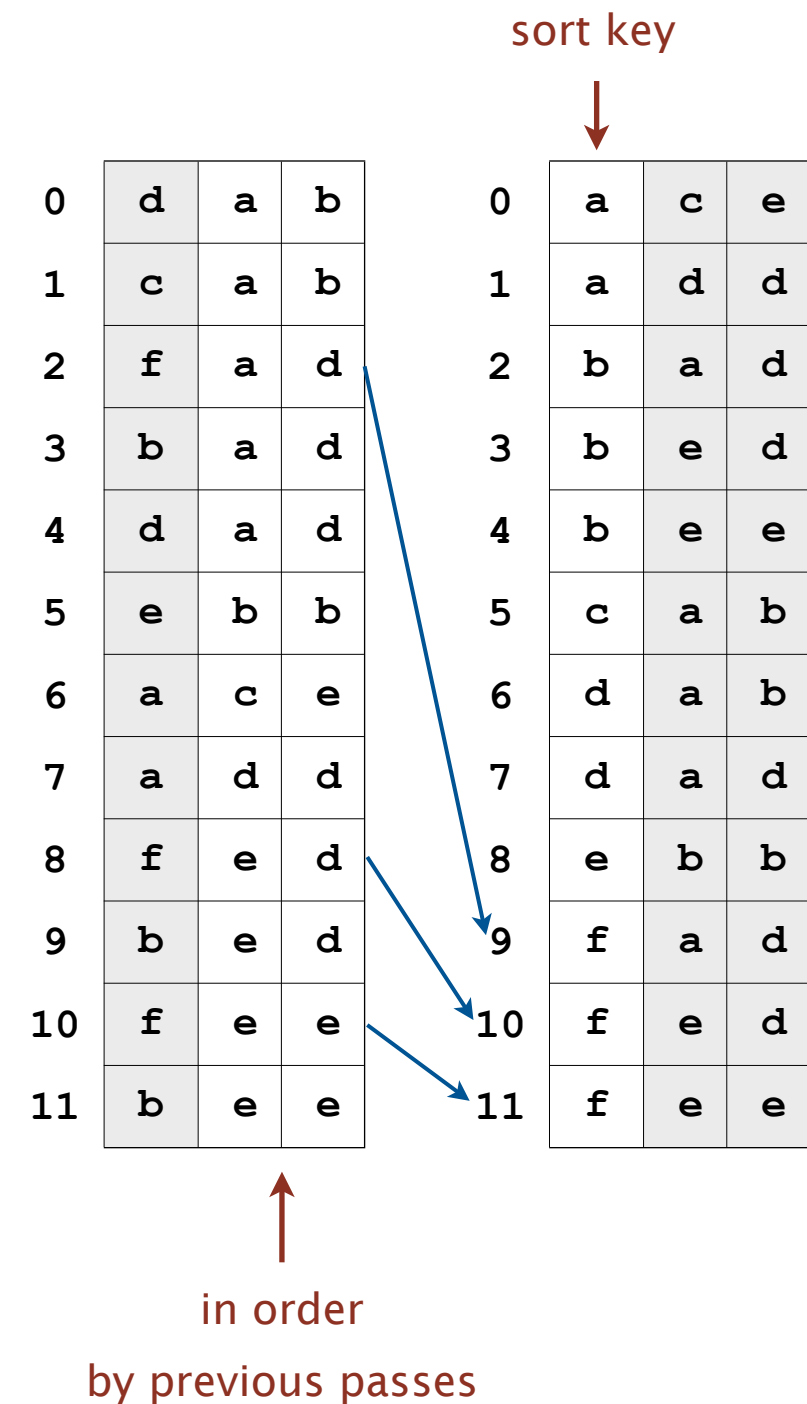|     | sort key |   |   |
| --- | --- | --- | --- |
| 0   | a | c | e |
| 1   | a | d | d |
| 2   | b | a | d |
| 3   | b | e | d |
| 4   | b | e | e |
| 5   | c | a | b |
| 6   | d | a | b |
| 7   | d | a | d |
| 8   | e | b | b |
| 9   | f | a | d |
| 10  | f | e | d |
| 11  | f | e | e |

# LSD string sort: correctness proof

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

**Pf.** [thinking about the future]

- If the characters not yet examined differ, it doesn't matter what we do now.
- If the characters not yet examined agree, stability ensures later pass won't affect order.

sort key

|    |   |   |   |
|----|---|---|---|
| 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 |

in order
by previous passes

# 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

# LSD string sort:  example

| Input | d = 6 | d = 5 | d = 4 | d = 3 | d= 2 | d= 1 | d = 0 | Output |
|---|---|---|---|---|---|---|---|---|
| 4PGC938 | 2IYE230 | 3CIO720 | 2IYE230 | 2RLA629 | 1ICK750 | 3ATW723 | 1ICK750 | 1ICK750 |
| 2IYE230 | 3CIO720 | 3CIO720 | 4JZY524 | 2RLA629 | 1ICK750 | 3CIO720 | 1ICK750 | 1ICK750 |
| 3CIO720 | 1ICK750 | 3ATW723 | 2RLA629 | 4PGC938 | 4PGC938 | 3CIO720 | 10HV845 | 10HV845 |
| 1ICK750 | 1ICK750 | 4JZY524 | 2RLA629 | 2IYE230 | 10HV845 | 1ICK750 | 10HV845 | 10HV845 |
| 10HV845 | 3CIO720 | 2RLA629 | 3CIO720 | 1ICK750 | 10HV845 | 1ICK750 | 10HV845 | 10HV845 |
| 4JZY524 | 3ATW723 | 2RLA629 | 3CIO720 | 1ICK750 | 10HV845 | 2IYE230 | 2IYE230 | 2IYE230 |
| 1ICK750 | 4JZY524 | 2IYE230 | 3ATW723 | 3CIO720 | 3CIO720 | 4JZY524 | 2RLA629 | 2RLA629 |
| 3CIO720 | 10HV845 | 4PGC938 | 1ICK750 | 3CIO720 | 3CIO720 | 10HV845 | 2RLA629 | 2RLA629 |
| 10HV845 | 10HV845 | 10HV845 | 1ICK750 | 10HV845 | 2RLA629 | 10HV845 | 3ATW723 | 3ATW723 |
| 10HV845 | 10HV845 | 10HV845 | 10HV845 | 10HV845 | 2RLA629 | 10HV845 | 3CIO720 | 3CIO720 |
| 2RLA629 | 4PGC938 | 10HV845 | 10HV845 | 10HV845 | 3ATW723 | 4PGC938 | 3CIO720 | 3CIO720 |
| 2RLA629 | 2RLA629 | 1ICK750 | 10HV845 | 3ATW723 | 2IYE230 | 2RLA629 | 4JZY524 | 4JZY524 |
| 3ATW723 | 2RLA629 | 1ICK750 | 4PGC938 | 4JZY524 | 4JZY524 | 2RLA629 | 4PGC938 | 4PGC938 |

# Summary of the performance of sorting algorithms

Frequency of operations.

| algorithm | guarantee | random | extra space | stable? | operations on keys |
|:---------:|:---------:|:------:|:-----------:|:-------:|:------------------:|
| insertion sort | $N^2/2$ | $N^2/4$ | 1 | yes | `compareTo()` |
| mergesort | N lg N | N lg N | N | yes | `compareTo()` |
| quicksort | 1.39 N lg N * | 1.39 N lg N | c lg N | no | `compareTo()` |
| heapsort | 2 N lg N | 2 N lg N | 1 | no | `compareTo()` |
| LSD † | 2 W N | 2 W N | N + R | yes | `charAt()` |

Q. What if strings do not have same length?

\* probabilistic

† fixed-length W keys

# String sorting challenge 1

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

**Ex.** Account number, date, SS 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 | | |
| ?87-SS-0321 | | |
| KJ-0, 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 challenge 2

**Problem.** Sort 1 million 32-bit integers.

**Ex.** Google interview (or presidential interview).

## Which sorting method to use?

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



**Google CEO Eric Schmidt interviews Barack Obama**

# String Sorts

- key-indexed counting
- LSD string sort
- **MSD string sort**
- 3-way string quicksort
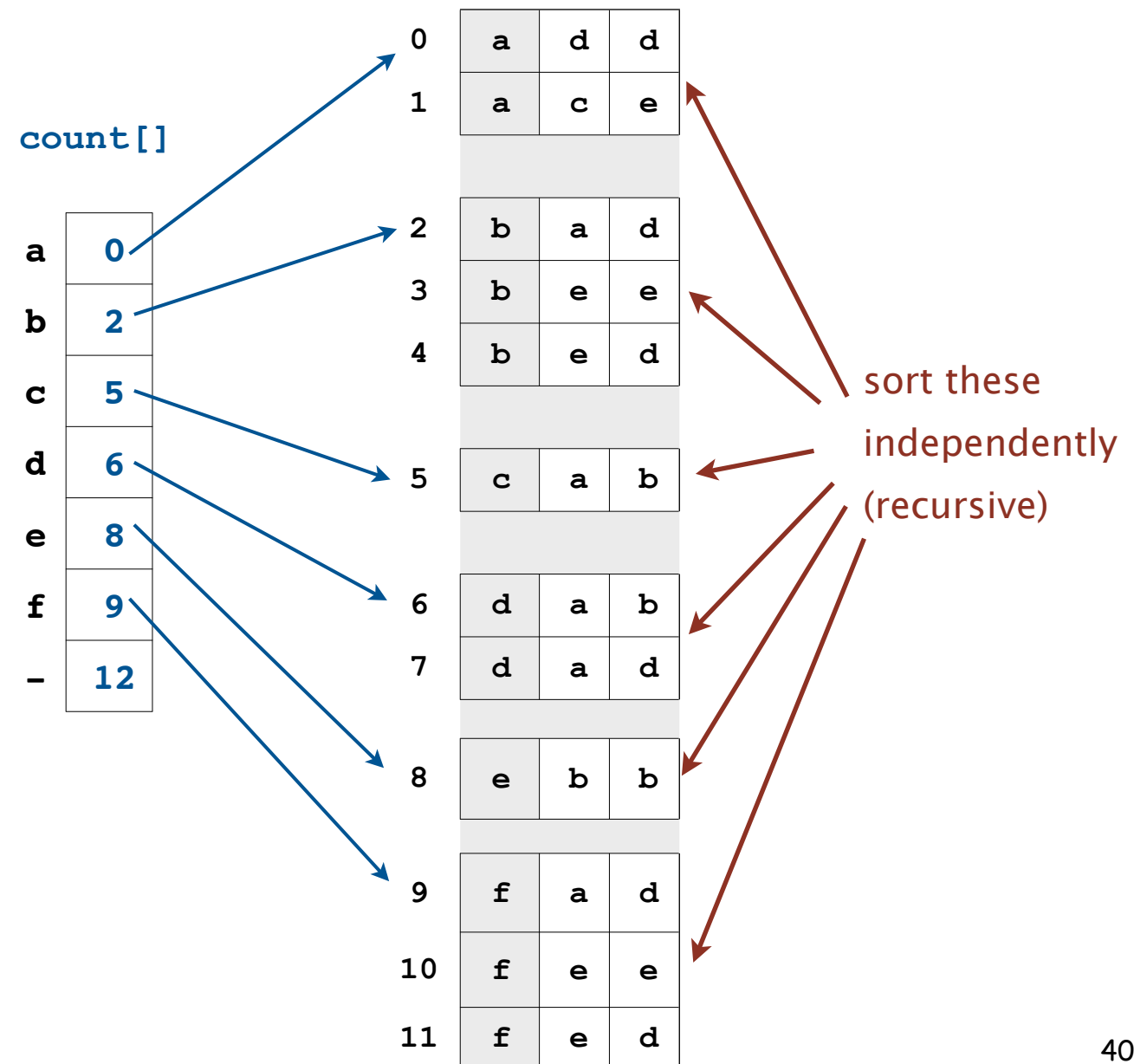- suffix arrays

# Most-significant-digit-first string sort

## MSD string sort.

- Partition file 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).



sort these
independently
(recursive)

sort key

# Variable-length strings

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

why smaller?

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 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: 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*

| | input | | count freq | | trans. | | distribute | | completion | sort calls | | output |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | she | 0 | 0 | 0 | 0 | 0 | are | 0 | 0 | 0 | sort(a, 0, 0); | 0 | are |
| 1 | sells | 1 a | 0 | 1 a | 0 | 1 | by | 1 a | 1 | sort(a, 1, 1); | 1 | by |
| 2 | seashells | 2 b | 1 | 2 b | 1 | 2 | she | 2 b | 2 | sort(a, 2, 1); | | |
| 3 | by | 3 c | 1 | 3 c | 2 | 3 | sells | 3 c | 2 | sort(a, 2, 1); | 2 | sea |
| 4 | the | 4 d | 0 | 4 d | 2 | 4 | seashells | 4 d | 2 | sort(a, 2, 1); | 3 | seashells |
| 5 | sea | 5 e | 0 | 5 e | 2 | 5 | sea | 5 e | 2 | sort(a, 2, 1); | 4 | seashells |
| 6 | shore | 6 f | 0 | 6 f | 2 | 6 | shore | 6 f | 2 | sort(a, 2, 1); | 5 | sells |
| 7 | the | 7 g | 0 | 7 g | 2 | 7 | shells | 7 g | 2 | sort(a, 2, 1); | 6 | sells |
| 8 | shells | 8 h | 0 | 8 h | 2 | 8 | she | 8 h | 2 | sort(a, 2, 1); | 7 | she |
| 9 | she | 9 i | 0 | 9 i | 2 | 9 | sells | 9 i | 2 | sort(a, 2, 1); | 8 | she |
| 10 | sells | 10 j | 0 | 10 j | 2 | 10 | surely | 10 j | 2 | sort(a, 2, 1); | 9 | shells |
| 11 | are | 11 k | 0 | 11 k | 2 | 11 | seashells | 11 k | 2 | sort(a, 2, 1); | 10 | shore |
| 12 | surely | 12 l | 0 | 12 l | 2 | 12 | the | 12 l | 2 | sort(a, 2, 1); | 11 | surely |
| 13 | seashells | 13 m | 0 | 13 m | 2 | 13 | the | 13 m | 2 | sort(a, 2, 1); | | |
| | | 14 n | 0 | 14 n | 2 | | | 14 n | 2 | sort(a, 2, 1); | 12 | the |
| | | 15 o | 0 | 15 o | 2 | | | 15 o | 2 | sort(a, 2, 1); | 13 | the |
| | | 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:  performance

## Number of characters examined.

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

| Random (sublinear) | Non-random with duplicates (nearly linear) | Worst case (linear) |
|---|---|---|
| 1EIO402 | are | 1DNB377 |
| 1HYL490 | by | 1DNB377 |
| 1ROZ572 | sea | 1DNB377 |
| 2HXE734 | seashells | 1DNB377 |
| 2IYE230 | seashells | 1DNB377 |
| 2XOR846 | sells | 1DNB377 |
| 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? | operations on keys |
|-----------|-----------|--------|-------------|---------|---------------------|
| insertion sort | $N^2/2$ | $N^2/4$ | 1 | yes | `compareTo()` |
| mergesort | $N \lg N$ | $N \lg N$ | N | yes | `compareTo()` |
| quicksort | $1.39 N \lg N$ [*] | $1.39 N \lg N$ | $c \lg N$ | no | `compareTo()` |
| heapsort | $2 N \lg N$ | $2 N \lg N$ | 1 | no | `compareTo()` |
| LSD [†] | $2 N W$ | $2 N W$ | $N + R$ | yes | `charAt()` |
| MSD [‡] | $2 N W$ | $N \log_R N$ | $N + D R$ | yes | `charAt()` |

stack depth D = length of longest prefix match

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

# MSD string sort vs. quicksort for strings

## Disadvantages of MSD string sort.

- Accesses memory "randomly" (cache inefficient).
- Inner loop has a lot of instructions.
- Extra space for `count[]`.
- Extra space for `aux[]`.

## Disadvantage of quicksort.

- Linearithmic number of string compares (not linear).
- Has to rescan long keys for compares.

Goal. Combine advantages of MSD and quicksort.