

## 6. Strings

**String.** Sequence of characters.

**Ex.** Natural languages, Java programs, genomic sequences, ...

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

### Using Strings in Java

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

**Substring.** Extract a contiguous list of characters from a string.



```
String s = "strings";           // s = "strings"
char c = s.charAt(2);          // c = 'r'
String t = s.substring(2, 6);  // t = "ring"
String u = s + t;              // u = "stringsring"
```

### Implementing Strings In Java

**Memory.** 40 + 2N bytes for a virgin string!

↙ could use byte array instead of String to save space

```
public final class String implements Comparable<String> {
    private char[] value; // characters
    private int offset;   // index of first char into 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);
    }
    ...
}
java.lang.String
```

## String vs. StringBuilder

**String.** [immutable] Fast substring, slow concatenation.

**StringBuilder.** [mutable] Slow substring, fast (amortized) append.

```
public static String reverse(String s) {
    String rev = "";
    for (int i = s.length() - 1; i >= 0; i--)
        rev += s.charAt(i);
    return rev;
}
```

quadratic time

```
public static String reverse(String s) {
    StringBuilder rev = new StringBuilder();
    for (int i = s.length() - 1; i >= 0; i--)
        rev.append(s.charAt(i));
    return rev.toString();
}
```

linear time

## Radix Sorting

---

Reference: Chapter 13, Algorithms in Java, 3<sup>rd</sup> Edition, Robert Sedgewick.

## Longest Common Prefix

**Longest common prefix.** Given two strings, find the common prefix that is as long as possible.

p	r	e	f	i	x		
0	1	2	3	4	5	6	7
p	r	e	f	e	t	c	h

```
public static String lcp(String s, String t) {
    int n = Math.min(s.length(), t.length());
    for (int i = 0; i < n; i++) {
        if (s.charAt(i) != t.charAt(i))
            return s.substring(0, i);
    }
    return s.substring(0, n);
}
```

## Radix Sorting

**Radix sorting.**

- Specialized sorting solution for strings.
- Same ideas for bits, digits, etc.

**Applications.**

- Bioinformatics.
- Sorting strings.
- Full text indexing.
- Plagiarism detection.
- Burrows-Wheeler transform. [see data compression]

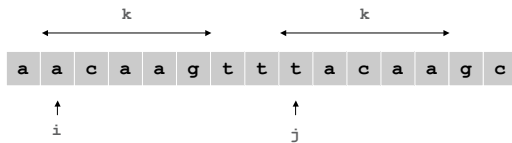
## An Application: Redundancy Detector

### Longest repeated substring.

- Given a string of  $N$  characters, find the longest repeated substring.
- Ex: `a a c a a g t t t a c a a g c`
- Application: bioinformatics.

### Dumb brute force.

- Try all indices  $i$  and  $j$ , and all match lengths  $k$ , and check.
- $O(W N^3)$  time, where  $W$  is length of longest match.



9

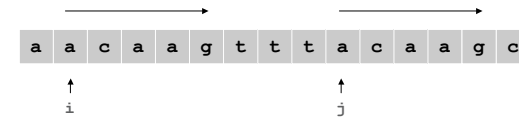
## An Application: Redundancy Detector

### Longest repeated substring.

- Given a string of  $N$  characters, find the longest repeated substring.
- Ex: `a a c a a g t t t a c a a g c`
- Application: bioinformatics.

### Brute force.

- Try all indices  $i$  and  $j$  for start of possible match, and check.
- $O(W N^2)$  time, where  $W$  is length of longest match.



10

## An Application: Redundancy Detector

### Longest repeated substring.

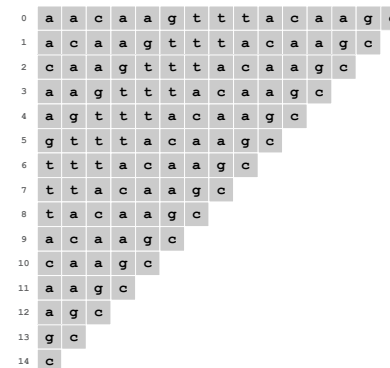
- Given a string of  $N$  characters, find the longest repeated substring.
- Ex: `a a c a a g t t t a c a a g c`
- Application: bioinformatics.

### Suffix sort solution.

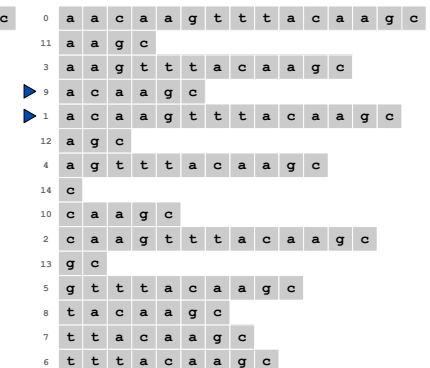
- Form  $N$  **suffixes** of original string.
- Sort to bring longest repeated substrings together.
- $O(W N \log N)$  time, where  $W$  is length of longest match.

## Suffix Sorting

### suffixes



### sorted suffixes



11

12

## Suffix Sorting: Java Implementation

```
public class LRS {
    public static void main(String[] args) {
        String s = StdIn.readAll();
        int N = s.length();

        String[] suffixes = new String[N];
        for (int i = 0; i < N; i++)
            suffixes[i] = s.substring(i, N);

        Arrays.sort(suffixes);

        String lrs = "";
        for (int i = 0; i < N - 1; i++) {
            String x = lcp(suffixes[i], suffixes[i+1]);
            if (x.length() > lrs.length()) lrs = x;
        }
        System.out.println(lrs);
    }
}
```

read input  
create suffixes (linear time)  
sort suffixes  
find longest match

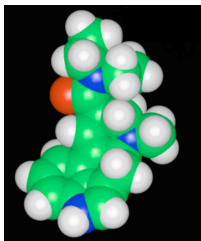
```
% java LRS < moby dick.txt
,- Such a funny, sporty, gamy, jesty, joky, hoky-poky lad, is the Ocean, oh! Th
```

## String Sorting Performance

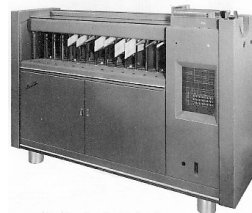
	String Sort	Suffix (sec)
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 <sup>§</sup>
Quicksort	$W N \log N$ <sup>†</sup>	9.5

W = max length of string  
N = number of strings  
§ estimate  
† probabilistic guarantee  
1.2 million for Moby Dick

## LSD Radix Sort



Lysergic Acid Diethylamide, Circa 1960



Card Sorter, Circa 1960

## String Sorting Notation

### Notation.

- String = sequence of characters.
- W = max # characters per string.
- N = # input strings.
- R = radix.  
256 for extended ASCII, 65,536 for original Unicode

### Java syntax.

- Array of strings: `String[] a`
- Number of strings: `N = a.length`
- The *i*<sup>th</sup> string: `a[i]`
- The *d*<sup>th</sup> character of the *i*<sup>th</sup> string: `a[i].charAt(d)`
- Strings to be sorted: `a[0], ..., a[N-1]`

## Key Indexed Counting

### Key indexed counting.

- Count frequencies of each letter. [0<sup>th</sup> character]

```
int N = a.length;
int[] count = new int[256+1];
for (int i = 0; i < N; i++) {
    char c = a[i].charAt(d);
    count[c+1]++;
}
```

frequencies

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

18

## Key Indexed Counting

### Key indexed counting.

- Count frequencies of each letter. [0<sup>th</sup> character]
- Compute cumulative frequencies.

```
for (int i = 1; i < 256; i++)
    count[i] += count[i-1];
```

cumulative counts

	a				count			
0	d	a	b		a	0	a	0
1	a	d	d		b	2	b	2
2	c	a	b		c	3	c	5
3	f	a	d		d	1	d	6
4	f	e	e		e	2	e	8
5	b	a	d		f	1	f	9
6	d	a	d		g	3	g	12
7	b	e	e					
8	f	e	d					
9	b	e	d					
10	e	b	b					
11	a	c	e					

19

## Key Indexed Counting

### Key indexed counting.

- Count frequencies of each letter. [0<sup>th</sup> character]
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

```
for (int i = 0; i < N; i++) {
    char c = a[i].charAt(d);
    temp[count[c]++] = a[i];
}
```

rearrange

	a				count				temp			
0	d	a	b		a	2	0	a	d	d		
1	a	d	d		b	5	1	a	c	e		
2	c	a	b		c	6	2	b	a	d		
3	f	a	d		d	8	3	b	e	e		
4	f	e	e		e	9	4	b	e	d		
5	b	a	d		f	12	5	c	a	b		
6	d	a	d		g	12	6	d	a	b		
7	b	e	e				7	d	a	d		
8	f	e	d				8	e	b	b		
9	b	e	d				9	f	a	d		
10	e	b	b				10	f	e	e		
11	a	c	e				11	f	e	d		

32

## Key Indexed Counting

### Key indexed counting.

- Count frequencies of each letter. [0<sup>th</sup> character]
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

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

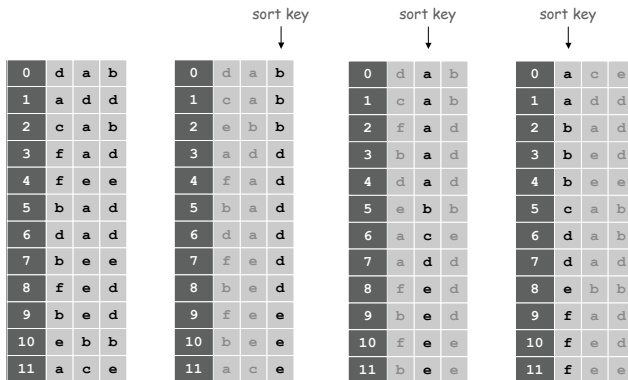
copy back

	a				count				temp			
0	a	d	d		a	2	0	a	d	d		
1	a	c	e		b	5	1	a	c	e		
2	b	a	d		c	6	2	b	a	d		
3	b	e	e		d	8	3	b	e	e		
4	b	e	d		e	9	4	b	e	d		
5	c	a	b		f	12	5	c	a	b		
6	d	a	b		g	12	6	d	a	b		
7	d	a	d				7	d	a	d		
8	e	b	b				8	e	b	b		
9	f	a	d				9	f	a	d		
10	f	e	e				10	f	e	e		
11	f	e	d				11	f	e	d		

33

## Least Significant Digit Radix Sort

**LSD.** Consider digits  $d$  from right to left: stably sort using  $d$ th character as the key via key-indexed counting.



34

## Least Significant Digit Radix Sort

**LSD.** Consider digits  $d$  from right to left: stably sort using  $d$ th character as the key via key-indexed counting.

```
public static void lsd(String[] a) {
    int W = a[0].length();
    for (int d = W-1; d >= 0; d--) {
        // do key-indexed counting sort on digit d
        ...
    }
}
```

Assumes fixed length strings (length =  $W$ )

35

## LSD Radix Sort: Correctness

**Pf 1.** [left-to-right]

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

**Pf 2.** [right-to-left]

- 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.

now	sob	cab	ace
for	ncb	wad	ago
tip	cab	tag	and
ilk	wad	jam	bet
dim	ard	rap	cab
tag	ace	tap	caw
jot	wee	tar	cue
sob	cue	was	dim
nob	fee	caw	dug
sky	tag	raw	egg
hut	egg	jay	fee
ace	gig	ace	few
bet	dug	wee	for
men	ilk	fee	gig
egg	owl	men	hut
few	dim	bet	ilk
jay	jam	few	jam
owl	men	egg	jay
joy	ago	ago	jot
rap	tip	gig	joy
gig	rap	dim	men
wee	tap	tip	nob
was	for	sky	now
cab	tar	ilk	owl
wad	was	and	rap
tap	jot	sob	raw
caw	hut	nob	sky
cue	bet	for	sob
fee	you	jot	tag
raw	ncw	you	tap
ago	few	now	tar
tar	caw	joy	tip
jam	raw	cue	wad
dug	sky	dug	was
you	jay	hut	wee
and	jcy	owl	you

36

## LSD Radix Sort Correctness

**Running time.**  $\Theta(W(N + R))$ .

why doesn't it violate  $N \log N$  lower bound?

**Advantage.** Fastest sorting method for random fixed length strings.

**Disadvantages.**

- Accesses memory "randomly."
- Inner loop has a lot of instructions.
- Wastes time on low-order characters.
- Doesn't work for variable-length strings.
- Not much semblance of order until very last pass.

**Goal.** Find fast algorithm for **variable** length strings.

37

# MSD Radix Sort

## Most significant digit radix sort.

- Partition file into 256 pieces according to first character.
- Recursively sort all strings that start with the same character, etc.

Q. How to sort on d<sup>th</sup> character?

A. Key-indexed counting.

now	a	ce	ac	e	ace
for	a	go	ag	o	ago
tip	a	nd	an	d	and
ilk	b	et	be	t	bet
dim	c	ab	ca	b	cab
tag	c	aw	ca	w	caw
jot	c	ue	cu	e	cue
sob	d	im	di	m	dim
nob	d	ug	du	g	dug
sky	e	gg	eg	g	egg
hut	f	or	fe	w	fee
ace	f	ee	fe	e	fee
bet	f	ew	fo	r	for
men	g	ig	gi	g	gig
egg	h	ut	hu	t	hut
few	i	lk	il	k	ilk
jay	j	am	ja	y	jam
owl	j	ay	ja	m	jam
joy	j	ot	jo	t	jot
rap	j	oy	jo	y	joy
gig	m	en	me	n	men
wee	n	ow	no	w	nob
was	n	ob	no	b	now
cab	o	wl	ow	l	owl
wad	r	ap	ra	p	rap
caw	s	ob	sk	y	sky
cue	s	ky	so	b	sob
fee	t	ip	ta	g	tag
tap	t	ag	ta	p	tap
ago	t	ap	ta	r	tar
tar	t	ar	ti	p	tip
jam	w	ee	wa	d	wad
dug	w	as	wa	s	was
and	w	ad	we	e	wee

## MSD Radix Sort Implementation

```

public static void msd(String[] a) {
    msd(a, 0, a.length, 0);
}

private static void msd(String[] a, int l, int r, int d) {
    if (r <= l + 1) return;

    // key-indexed counting sort on digit d of a[l] to a[r]
    int[] count = new int[256+1];
    ...

    // recursively sort 255 subfiles - assumes '\0' terminated
    for (int i = 0; i < 255; i++)
        msd(a, l + count[i], l + count[i+1], d+1);
}
    
```

inclusive      exclusive

## String Sorting Performance

	String Sort	Suffix (sec)
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 <sup>§</sup>
Quicksort	$W N \log N$ <sup>†</sup>	9.5
LSD*	$W(N + R)$	-
MSD	$W(N + R)$	395

R = radix  
W = max length of string  
N = number of strings

§ estimate  
\* fixed length strings only  
† probabilistic guarantee

↖  
1.2 million for Moby Dick

## MSD Radix Sort: Small Files

### Disadvantages.

- Too slow for small files.
  - ASCII: 100x slower than insertion sort for  $N = 2$
  - Unicode: 30,000x slower for  $N = 2$
- Huge number of recursive calls on small files.

**Solution.** Cutoff to insertion sort for small  $N$ .

**Consequence.** Competitive with quicksort for string keys.

## String Sorting Performance

	String Sort	Suffix (sec)
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 §
Quicksort	$W N \log N †$	9.5
LSD *	$W(N + R)$	-
MSD	$W(N + R)$	395
MSD with cutoff	$W(N + R)$	6.8

R = radix  
 W = max length of string  
 N = number of strings

§ estimate  
 \* fixed length strings only  
 † probabilistic guarantee

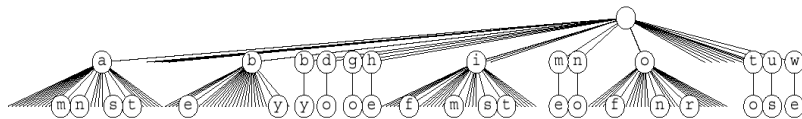
↑ 1.2 million for Moby Dick

42

43

## Recursive Structure: MSD

**Recursive structure.** R-way branching leads to lots of empty calls.



## 3-Way Radix Quicksort

44



### 3-Way Radix Quicksort

**Idea 1.** Use  $d^{\text{th}}$  character to "sort" into 3 pieces instead of 256, and sort each piece recursively.

**Idea 2.** Keep all duplicates together in partitioning step.

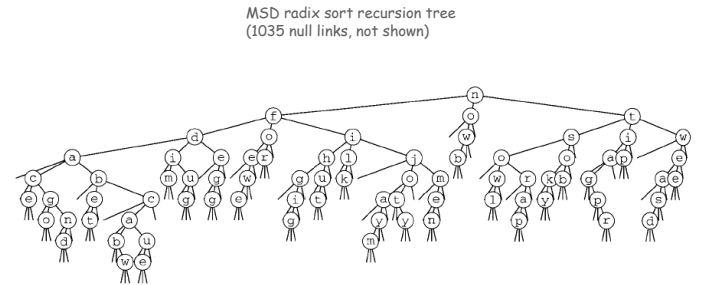
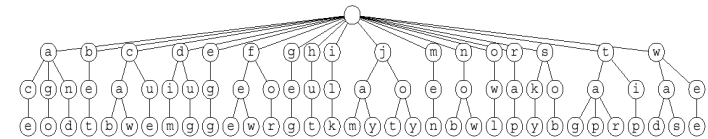
actinian	coenobite	actinian	now	gig	ace	ago	a go
jeffrey	conelrad	bracteal	for	for	bet	bet	a ce
coenobite	coenobite	coenobite	tip	dug	dug	and	a nd
conelrad	bracteal	conelrad	iik	iik	cab	ace	b et
secureness	secureness	conelrad	dim	dim	dim	c iab	
cumin	dilatedly	cumin	tag	ago	ago	c aw	
chariness	inkblot	centesimal	jot	and	and	c iue	
bracteal	jeffrey	cankerous	sob	fee	egg	egg	
displease	displease	circumflex	nob	cue	cue	dug	
millwright	millwright	millwright	sky	caw	caw	dim	
repertoire	repertoire	repertoire	hut	hut	f ee		
dourness	dourness	dourness	ace	ace	f or		
centesimal	southeast	southeast	bet	bet	f ew		
fondler	fondler	fondler	men	cab	i ik		
interval	interval	interval	egg	egg	g ig		
reversionary	reversionary	reversionary	few	few	hut		
dilatedly	cumin	secureness	jay	j ay	jam		
inkblot	chariness	dilatedly	owl	ot	ja y		
southeast	centesimal	inkblot	joy	jo y	jo y		
cankerous	cankerous	jeffrey	rap	jam	jo t		
circumflex	circumflex	displease	gig	owl	owl	m en	
			wee	wee	now	owl	
			was	was	nob	nob	
			cab	men	men	now	
			wad	wad	r ap	sky	sky
			caw	sky	sky	sky	sky
			cue	nob	was	tip	sob
			fee	sob	sob	t ip	ta r
			tap	tap	tap	t ap	ta p
			ago	tag	tag	t ag	ta g
			tar	tar	tar	t ar	ti p
			dug	tip	tip	w ias	
			and	now	wee	w iee	
			jam	rap	wad	w iad	

3-way partition

3-way radix quicksort

### Recursive Structure: MSD vs. 3-Way Quicksort

3-way radix quicksort collapses empty links in MSD tree.



### 3-Way Radix Quicksort

```
private static void quicksortX(String a[], int lo, int hi, int d) {
    if (hi - lo <= 0) return;
    int i = lo-1, j = hi;
    int p = lo-1, q = hi;
    char v = a[hi].charAt(d);

    while (i < j) {
        while (a[++i].charAt(d) < v) if (i == hi) break;
        while (v < a[--j].charAt(d)) if (j == lo) break;
        if (i > j) break;
        exch(a, i, j);
        if (a[i].charAt(d) == v) exch(a, ++p, i);
        if (a[j].charAt(d) == v) exch(a, j, --q);
    }
    if (p == q) {
        if (v != '\0') quicksortX(a, lo, hi, d+1);
        return;
    }
    if (a[i].charAt(d) < v) i++;
    for (int k = lo; k <= p; k++) exch(a, k, j--);
    for (int k = hi; k >= q; k--) exch(a, k, i++);
    quicksortX(a, lo, j, d);
    if ((i == hi) && (a[i].charAt(d) == v)) i++;
    if (v != '\0') quicksortX(a, j+1, i-1, d+1);
    quicksortX(a, i, hi, d);
}
```

### Quicksort vs. 3-Way Radix Quicksort

#### Quicksort.

- $2N \ln N$  string comparisons on average.
- Long keys are costly to compare if they differ only at the end, and this is common case!
- absolutism, absolut, absolutely, absolute.

#### 3-way radix quicksort.

- Avoids re-comparing initial parts of the string.
- Uses just "enough" characters to resolve order.
- $2N \ln N$  character comparisons on average for random strings.
- Sub-linear sort for large  $W$  since input is of size  $NW$ .

**Theorem.** Quicksort with 3-way partitioning is OPTIMAL.

**Pf.** Ties cost to entropy. Beyond scope of 226.

## String Sorting Performance

	String Sort	Suffix (sec)
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 §
Quicksort	$W N \log N \uparrow$	9.5
LSD *	$W(N + R)$	-
MSD	$W(N + R)$	395
MSD with cutoff	$W(N + R)$	6.8
3-way radix quicksort	$W N \log N \uparrow$	2.8

R = radix  
W = max length of string  
N = number of strings

§ estimate  
\* fixed length strings only  
† probabilistic guarantee

↑ 1.2 million for Moby Dick

## Suffix Sorting: Worst Case Input

Length of longest match small.

- Hard to beat 3-way radix quicksort.

Length of longest match very long.

- 3-way radix quicksort is quadratic.
- Ex: two copies of Moby Dick.

Can we do better?  $\Theta(N \log N)$ ?  $\Theta(N)$ ?

Observation. Must find longest repeated substrings while suffix sorting to beat  $N^2$ .

```

abcdefghi
abcdefghiabcdefghi
bcdefghi
bcdefghiabcdefghi
cdefghi
cdefghiabcdefghi
defghi
efghiabcdefghi
efghi
efghiabcdefghi
fghi
fghiabcdefghi
fghi
ghiabcdefghi
fhi
hiabcdefghi
hi
iabcdefghi
i
    
```

Input: "abcdeghiabcdefghi"

50

51

## Suffix Sorting in Linearithmic Time: Key Idea

		sorted
0	babaaaabcbabaaaaa0	17 0bab aaaa bcbabaaa aa
1	abaaaabcbabaaaaa0b	16 a0ba baaa abcbabaa aa
2	baaaabcbabaaaaa0ba	15 aa0b abaa aabcbaba aa
3	aaaabcbabaaaaa0bab	14 aaa0 baba aaabcbab aa
4	aaabcbabaaaaa0baba	3 aaaa bcba baaaaa0b ab
5	aabcbabaaaaa0babaa	12 aaaa a0ba baaaabcb ab
6	abcbabaaaaa0babaaa	13 aaaa 0bab aaaabcbaba
7	bcbabaaaaa0babaaaa	4 aaab cbab aaaaa0ba ba
8	cbabaaaaa0babaaaab	5 aabc baba aaaa0bab aa
9	babaaaaa0babaaaabc	1 abaa aabc babaaaaa 0b
10	abaaaaa0babaaaabcb	10 abaa aaa0 babaaaab cb
11	baaaaa0babaaaabcbaba	6 abcb abaa aaa0baba aa
12	aaaaa0babaaaabcbab	2 baaa abcb abaaaaa0 ba
13	aaaa0babaaaabcbaba	11 baaa aa0b abaaaabc ba
14	aaa0babaaaabcbabaa	0 baba aaab cbabaaaa a0
15	aa0babaaaabcbabaaa	9 baba aaaa 0babaaaa bc
16	a0babaaaabcbabaaaa	7 bcba baaa aa0babaa aa
17	0babaaaabcbabaaaaa	8 cbab aaaa a0babaaa ab

text = "babaaaabcbabaaaaa"

52

## Suffix Sorting in Sub-Quadratic Time

Manber's MSD algorithm.

- Phase 0: sort on first character using key-indexed sorting.
- Phase i: given list of suffixes sorted on first  $2^{i-1}$  characters, create list of suffixes sorted on first  $2^i$  characters
- Finishes after  $\lg N$  phases.

Manber's LSD algorithm.

- Same idea but go from right to left.
- $O(N \log N)$  guaranteed running time.
- $O(N)$  extra space (but need several auxiliary arrays).

Best in theory.  $O(N)$  but more complicated to implement.

53

## String Sorting Performance

	String Sort	Suffix Sort (seconds)	
	Worst Case	Moby Dick	AesopAesop
Brute	$W N^2$	36,000 <sup>§</sup>	3,990 <sup>§</sup>
Quicksort	$W N \log N$ †	9.5	167
LSD *	$W(N + R)$	-	-
MSD	$W(N + R)$	395	memory
MSD with cutoff	$W(N + R)$	6.8	162
3-way radix quicksort	$W N \log N$ †	2.8	400
Manber ‡	$N \log N$	17	8.5

R = radix

W = max length of string.

N = number of strings

↖  
1.2 million for Moby Dick  
191 thousand for Aesop's Fables

§ estimate

\* fixed length strings only

† probabilistic guarantee

‡ suffix sorting only