

Radix Sorting

- LSD radix sort
- MSD radix sort
- 3-way radix quicksort
- Suffix sorting

Reference: Chapter 13, Algorithms in Java, 3rd Edition, Robert Sedgewick.

Radix Sorting

Radix sorting.

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

Applications.

- Sorting strings.
- Full text indexing.
- Plagiarism detection.
- Burrows-Wheeler transform. *stay tuned*
- Computational molecular biology.

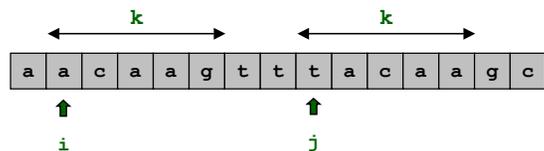
An Application: Redundancy Detector

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: computational molecular biology.

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.



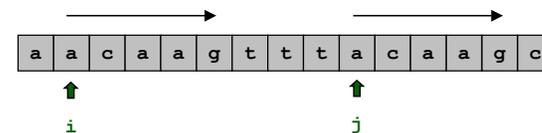
An Application: Redundancy Detector

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: computational molecular biology.

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.



String Sorting

Notation.

- String = variable length sequence of characters.
- W = max # characters per string.
- N = # input strings.
- R = radix (256 for extended ASCII, 65,536 for UNICODE).

Java syntax.

- Array of strings: `String[] a;`
- The i^{th} string: `a[i]`
- The d^{th} character of the i^{th} string: `a[i].charAt(d)`
- Strings to be sorted: `a[lo], ..., a[hi]`

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th character)

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

frequencies

`d = 0;`

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		

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th 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
1	a	d	d	b	2
2	c	a	b	c	5
3	f	a	d	d	6
4	f	e	e	e	8
5	b	a	d	f	9
6	d	a	d	g	11
7	b	e	e		
8	f	e	d		
9	b	e	d		
10	e	b	b		
11	a	c	e		

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th character)
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

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

rearrange

`d = 0;`

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

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th character)
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

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

rearrange

d = 0;

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

13

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th character)
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

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

rearrange

d = 0;

	a	count	temp
0	d a b	a 8	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 11	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

24

Key Indexed Counting

Key indexed counting.

- Count frequencies of each letter. (0th character)
- Compute cumulative frequencies.
- Use cumulative frequencies to rearrange strings.

```
for (int i = L; i <= R; i++)
    a[i] = temp[i - L];
```

copy back

	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 11	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

25

LSD Radix Sort

Least significant digit radix sort.

- Ancient method used for card-sorting.
- Consider digits from right to left:
 - use key-indexed counting to STABLE sort by character

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

27

LSD Radix Sort

Least significant digit radix sort.

- Ancient method used for card-sorting.
- Consider digits from right to left:
 - use key-indexed counting to STABLE sort by character

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

Fixed length strings (length = W)

LSD Radix Sort: Correctness

Proof 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.

Proof 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, later pass won't affect order.

now	sob	caw	ace
for	nob	wad	ago
tip	cab	tag	and
ilk	wad	jam	bet
dim	and	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	yqu	jot	tag
raw	now	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

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.

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.

How to sort on d^{th} character?

- Use 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	few
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	jay
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, int lo, int hi) {
    msd(a, lo, hi, 0);
}

private static void msd(String[] a, int lo, int hi, int d) {
    if (hi <= lo) return;

    // do key-indexed counting sort on digit d
    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] - 1, d+1);
}
    
```

32

String Sorting Performance

	String Sort	Suffix (sec)
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 ^s
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.

^s estimate
* assumes fixed length strings.
[†] probabilistic guarantee.

33

MSD Radix Sort Analysis

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.

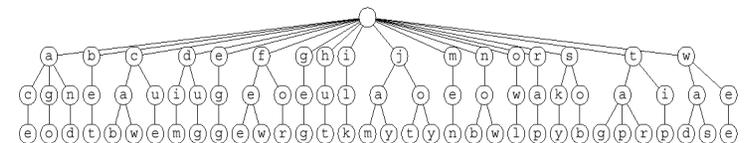
Solution: cutoff to insertion sort for small N.

- Competitive with quicksort for string keys.

34

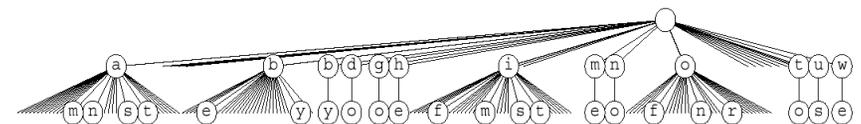
Recursive Structure of MSD Radix Sort

Trie structure to describe recursive calls in MSD radix sort.



Problem: algorithm touches lots of empty nodes ala R-way tries.

- Tree can be as much as 256 times bigger than it appears!

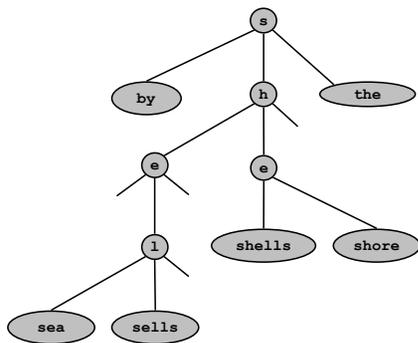


35

Correspondence With Sorting Algorithms

Correspondence between trees and sorting algorithms.

- BSTs correspond to quicksort recursive partitioning structure.
- R-way tries corresponds to MSD radix sort.
- What corresponds to ternary search tries?



36

3-Way Radix Quicksort

- Idea 1.** Use d^{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
jeffrey	conelrad	bracteal
coenobite	actinian	coenobite
conelrad	bracteal	conelrad
secureness	secureness	cumin
cumin	dilatedly	chariness
chariness	inkblot	centesimal
bracteal	jeffrey	cankerous
displease	displease	circumflex
millwright	millwright	millwright
repertoire	repertoire	repertoire
dourness	courness	dourness
centesimal	southeast	southeast
fondler	fondler	fondler
interval	interval	interval
reversionary	reversionary	reversionary
dilatedly	cumin	secureness
inkblot	chariness	dilatedly
southeast	centesimal	inkblot
cankerous	cankerous	jeffrey
circumflex	circumflex	displease

Partition

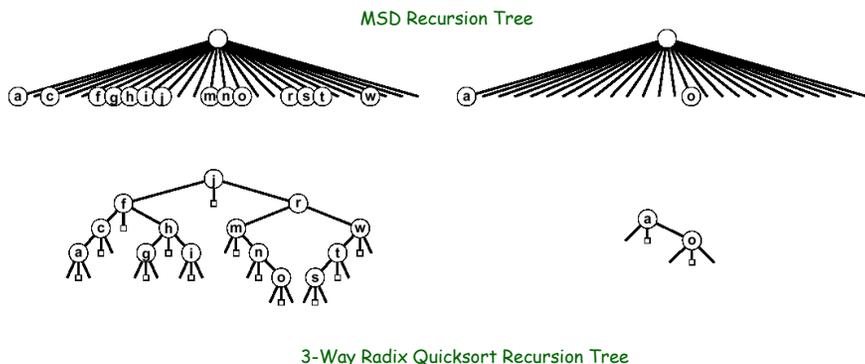
now	gig	ace	ago	a	go						
for	for	bet	bet	a	ce						
tip	dug	dug	and	a	nd						
ilk	ilk	cab	ace	b	et						
dim	dim	dim	c	a	b						
tag	ago	ago	c	a	w						
jot	and	and	c	u	e						
sob	fee	egg	egg								
nob	cue	dug	dug								
sky	caw	caw	dim								
hut	hut	f	e								
ace	ace	f	o	r							
bet	bet	f	e	w							
men	cab	ilk									
egg	egg	gig									
few	few	hut									
jay	j	a	y	j	a	m					
owl	j	o	t	j	a	y					
joy	j	o	y	j	a	y					
rap	j	a	m	o	t						
gig	owl	owl	m	e	n						
wee	wee	now	owl								
was	was	nob	nob								
cab	men	men	now								
wad	wad	f	a	p							
caw	sky	sky	sky	sky							
cue	nob	was	tip	so	b						
fee	sob	sob	sob	t	i	p	t	a	r		
tap	tap	tap	tap	t	a	p	t	a	p		
ago	tag	tag	tag	t	a	g	t	a	g		
tar	tar	tar	tar	t	a	r	t	a	r		
dug	tip	tip	was	w	e	e	w	e	e		
and	now	wee	wee	w	e	e	w	e	e		
jam	rap	wad	w	a	d	w	a	d	w	a	d

Algorithm

37

Recursive Structure of MSD Radix Sort vs. 3-Way Quicksort

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



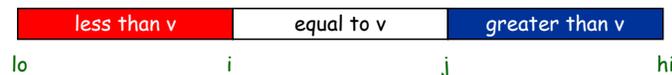
3-Way Radix Quicksort Recursion Tree

38

3-Way Partitioning

3-way partitioning.

- Natural way to deal with equal keys.
- Partition elements into 3 parts:
 - elements between i and j equal to partition element v
 - no larger elements to left of i
 - no smaller elements to right of j



Dutch national flag problem.

- Not easy to implement efficiently. (Try it!)
- Not done in practical sorts before mid-1990s.
- Incorporated into Java system sort, C qsort.

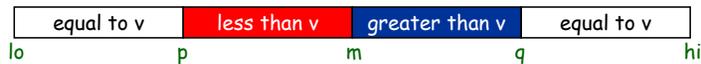


39

3-Way Partitioning

Elegant solution to Dutch national flag problem.

- Partition elements into 4 parts:
 - no larger elements to left of m
 - no smaller elements to right of m
 - equal elements to left of p
 - equal elements to right of q



- Afterwards, swap equal keys into center.

All the right properties.

- Not much code.
- In-place.
- Linear if keys are all equal.
- Small overhead if no equal keys.

40

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, p = lo-1, q = hi;
    char v = a[hi].charAt(d);

    while (i < j) {
        while (a[++i].charAt(d) < v); // repeat until pointers cross
        while (v < a[--j].charAt(d)); // find i on left and j on right to swap
        if (j == lo) break;
        if (i > j) break;
        exch(a, i, j);
        if (a[i].charAt(d) == v) { p++; exch(a, p, i); } // swap equal chars
        if (a[j].charAt(d) == v) { q--; exch(a, j, q); } // to left or right
    }
    if (p == q) {
        if (v != '\0') quicksortX(a, lo, hi, d+1); // special case for
        return; // all equal chars
    }
    if (a[i].charAt(d) < v) i++;
    for (int k = lo; k <= p; k++, j--) exch(a, k, j); // swap equal ones
    for (int k = hi; k >= q; k--, i++) exch(a, k, i); // back to middle
    quicksortX(a, lo, j, d);
    if ((i == hi) && (a[i].charAt(d) == v)) i++; // sort 3 pieces
    if (v != '\0') quicksortX(a, j+1, i-1, d+1); // recursively
    quicksortX(a, i, hi, d);
}
```

41

Significance of 3-Way Partitioning

Equal keys omnipresent in applications when purpose of sort is to bring records with equal keys together.

- Finding collinear points.
- Sort population by age.
- Remove duplicates from mailing list.
- Sort job applicants by college attended.

Typical application.

- Huge file.
- Small number of key values.
- Randomized 3-way quicksort is LINEAR time. (Try it!)

Theorem. Quicksort with 3-way partitioning is OPTIMAL.

Proof. Ties cost to entropy. Beyond scope of 226.

42

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 .

43

String Sorting Performance

	String Sort	Suffix Sort
	Worst Case	Moby Dick
Brute	$W N^2$	36,000 [§]
Quicksort	$W N \log N^\dagger$	9.5
LSD *	$W(N + R)$	-
MSD	$W(N + R)$	395
MSD with cutoff	$W(N + R)$	6.8
3-Way Radix Quicksort	$W N \log N^\dagger$	2.8

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

§ estimate
* fixed length strings only
† probabilistic guarantee

44

Suffix Sorting: Worst Case Input

Length of longest match small.

- 3-way radix quicksort rules!

Length of longest match very long.

- 3-way radix quicksort is quadratic.
- Two copies of Moby Dick.

Can we do better?

- $\Theta(N \log N)$?
- $\Theta(N)$?

Observation. Must find longest repeated substring WHILE suffix sorting to beat N^2 .

```

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

Input: "abcdefghiabcdefghi"

45

Suffix Sorting in $N \log N$ Time: Key Idea

```

0 babaaaabcbabaaaa0
1 abaaaabcbabaaaa0b
2 baaaabcbabaaaa0ba
3 aaaabcbabaaaa0bab
4 aaabcbabaaaa0baba
5 aabcbabaaaa0babaa
6 abcbabaaaa0babaaa
7 bcbabaaaa0babaaaa
8 cbabaaaa0babaaaab
9 babaaaa0babaaaabc
10 abaaaa0babaaaabcb
11 baaaa0babaaaabcbab
12 aaaa0babaaaabcbab
13 aaaa0babaaaabcbaba
14 aa0babaaaabcbabaa
15 aa0babaaaabcbabaaa
16 a0babaaaabcbabaaaa
17 0babaaaabcbabaaaa

17 0babaaaabcbabaaaa
16 a0babaaaabcbabaaaa
15 aa0babaaaabcbabaaa
14 aaa0babaaaabcbabaa
3 aaaabcbabaaaa0bab
12 aaaa0babaaaabcbab
13 aaaa0babaaaabcbaba
4 aaabcbabaaaa0baba
5 aabcbabaaaa0babaa
1 abaaaabcbabaaaa0b
10 abaaaa0babaaaabcb
6 abcbabaaaa0babaaa
2 baaaabcbabaaaa0ba
11 baaaa0babaaaabcbab
0 babaaaabcbabaaaa0
9 babaaaa0babaaaabc
7 bcbabaaaa0babaaaa
8 cbabaaaa0babaaaab
    
```

Input: "babaaaabcbabaaaa"

46

Suffix Sorting in $N \log N$ Time

Manber's MSD algorithm.

- Phase 0: sort on first character using key-indexed sorting.
- Phase n: given list of suffixes sorted on first n characters, create list of suffixes sorted on first 2n 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.

47

String Sorting Performance

	String Sort	Suffix Sort (seconds)	
	Worst Case	Moby Dick	AesopAesop
Brute	$W N^2$	36,000 [§]	3,990 [§]
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.

§ estimate
* fixed length strings only
† probabilistic guarantee
‡ suffix sorting only