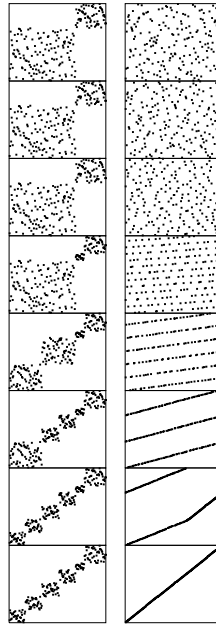


## COS 226 Lecture 6: Radix sorting

- Bits and digits
- Binary Quicksort
- MSD radix sort
- Three-way radix Quicksort
- LSD radix sort
- Sorting in linear time



6.1

## Extracting digits with macros

```
#define bitsword 32
#define bitsbyte 8
#define bytesword 4
#define R (1 << bitsbyte)
#define digit(A, B)
    ((A >> (bitsword-(B+1)*bitsbyte)) & (R-1))
```

**Ex: Single-byte access: bitsbyte = 8**

```
x = 0X61626364
digit(x, 2) = (x >> 8) & 255 = c
0110 0001 0110 0010 0110 0011 0110 0100  x
0000 0000 0110 0001 0110 0010 0110 0011  x >> 8
0000 0000 0000 0000 0000 0000 1111 1111  255 (R-1)
0000 0000 0000 0000 0000 0000 0110 0011  c
```

**Ex: Single-bit access: bitsbyte = 1**

```
digit(x, 11) = (x >> 20) & 1 = 0
0110 0001 0110 0010 0110 0011 0110 0100  x
0000 0000 0000 0000 0000 0110 0001 0110  x >> 20
0000 0000 0000 0000 0000 0000 0000 0001  1 (R-1)
```

6.3

## Bits and digits

Extracting bits is easy in C

Radix: base of number system

Power of 2 radix: groups of bits

- binary (radix-2): 1 bit at a time
- hexadecimal (radix-16): 4 bits at a time
- ascii (radix-256): 8 bits at a time

```
bin  01100001011000100110001101100100
hex   6  1  6  2  6  3  6  4
ascii  a  b  c  d
```

6.2

## Binary Quicksort

Partition file into two pieces

- all keys with first bit 0
- all keys with first bit 1

Sort two pieces recursively

Equivalent to partitioning on the VALUE  $2^{(\text{bitsword}-w+i)}$

- instead of some key in the file.

Bad partition if all keys have same leading bit

- one subfile of size N
- one empty subfile
- BUT keys one bit shorter

Worst case: one pass per key bit

6.4

## Binary Quicksort code

```
quicksortB(int a[], int l, int r, int w)
{ int i = l, j = r;
  if (r <= l || w > bitsword) return;
  while (j != i)
  {
    while (digit(a[i], w) == 0 && (i < j)) i++;
    while (digit(a[j], w) == 1 && (j > i)) j--;
    exch(a[i], a[j]);
  }
  if (digit(a[r], w) == 0) j++;
  quicksortB(a, l, j-1, w+1);
  quicksortB(a, j, r, w+1);
}
```

6.5

## Binary Quicksort issues

Problems:

- leading 0 bits
- cost of inner loop  
(could be advantage if carefully done)

Worst case: all keys equal

- $32N$  passes on a 32-bit machine
- $64N$  passes on a 64-bit machine

Good way to avoid quadratic worst case of quicksort

Random bits?

- should sort out after  $\lg N$  bits examined

Nonrandom bits?

- take bigger chunks

6.7

## Binary Quicksort example

A	00001	A	00001	A	00001	A	00001	A	00001	A	00001
S	10011	E	0101	E	0101	A	00001	A	00001	A	00001
O	01111	O	01111	A	00001	E	00101	E	00101	E	00101
R	10010	L	01100	E	00101	E	00101	E	00101	E	00101
T	10100	M	01101	G	00111	G	00111	G	00111	G	00111
I	01001	I	01001	I	01001	I	01001	I	01001	I	01001
N	01110	N	01110	N	01110	N	01110	L	01100	L	01100
G	00111	G	00111	M	01101	M	01101	M	01101	M	01101
E	00101	E	00101	L	01100	L	01100	N	01110	N	01110
X	11000	X	00001	O	01111	O	01111	O	01111	O	01111
A	00001	X	11000	S	10011	S	10011	P	10000	P	10000
M	01101	T	10100	T	10100	R	10010	R	10010	R	10010
P	10000	P	10000	P	10000	P	10000	S	10011	S	10011
L	01100	R	10010	R	10010	T	10100	T	10100	T	10100
E	00101	S	10011	X	11000	X	11000	X	11000	X	11000

6.6

## MSD radix sort

Partition file into R buckets

- all keys with first byte 0
- all keys with first byte 1
- all keys with first byte 2
- ...
- all keys with first byte R-1

Sort R pieces recursively

Take  $R=2^{\text{bitsbyte}}$

Tradeoff

- large R: space for buckets (too many empty buckets)
- small R: too many passes (too many keys per bucket)

Upper bound on running time:  $(\text{bytesword}) * (N + R)$

(Worst case: all keys equal)

6.8



## MSD radix sort code

```
#define bin(A) l+count[A]
void radixMSD(Item a[], int l, int r, int w)
{ int i, j, count[R+1];
  if (w > bytesword) return;
  if (r-l <= M) { insertion(a,l,r); return; }
  for (j = 0; j < R; j++) count[j] = 0;
  for (i = 1; i <= r; i++)
    count[digit(a[i], w) + 1]++;
  for (j = 1; j < R; j++)
    count[j] += count[j-1];
  for (i = 1; i <= r; i++)
    b[l+count[digit(a[i], w)]] = a[i];
  for (i = 1; i <= r; i++) a[i] = b[i];
  radixMSD(a, l, bin(0)-1, w+1);
  for (j = 0; j < R-1; j++)
    radixMSD(a, bin(j), bin(j+1)-1, w+1);
}
```

6.13

## LSD radix sort

Ancient (older than computers) method

- used for card-sorting

Consider digits from right to left

- use key-indexed counting (has to be stable)

Running time:  $N \cdot (\text{bitsword} / \text{bitsbyte})$

Disadvantage:

- doesn't work for variable-length keys
- totally out of order until MSD encountered

6.15

## MSD radix sort potential fatal flaw

each pass ALWAYS takes time proportional to  $N+R$

- initialize the buckets
- scan the keys

Ex: (ASCII bytes)  $R = 256$

- 100 times slower than insertion sort for  $N = 2$

Ex: (UNICODE)  $R = 65536$

- 30,000 times slower than insertion sort for  $N = 2$

TOO SLOW FOR SMALL FILES

RECURSIVE PROGRAM WILL CALL ITSELF

FOR A HUGE NUMBER OF SMALL FILES

Solution: cutoff to insertion sort

6.14

## LSD radix sort code

```
void radixLSD(Item a[], int l, int r)
{
  int i, j, w, count[R+1];
  for (w = bytesword-1; w >= 0; w--)
  {
    for (j = 0; j < R; j++) count[j] = 0;
    for (i = 1; i <= r; i++)
      count[digit(a[i], w) + 1]++;
    for (j = 1; j < R; j++)
      count[j] += count[j-1];
    for (i = 1; i <= r; i++)
      b[count[digit(a[i], w)]]] = a[i];
    for (i = 1; i <= r; i++) a[i] = b[i];
  }
}
```

6.16

### LSD radix sort example

now	sob	cab	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	jct	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

6.17

### Two proofs for LSD radix sort

#### Left-right

- if two keys differ on first bit  
o-1 sort puts them in proper relative order
- if two keys agree on first bit  
stability keeps them in proper relative order

#### Right-left

- if the bits not yet examined differ  
doesn't matter what we do now
- if the bits not yet examined agree  
later pass won't affect their order

6.19

### Binary LSD radix sort example

Cannot use Quicksort-style partitioning

- o-1 sort has to be stable
- stable inplace o-1 sort? (possible, but not easy)

A	00001	R	10010	T	10100	X	11000	P	10000	A	00001
S	10011	T	10100	X	11000	P	10000	A	00001	A	00001
O	01111	N	01110	P	10000	A	00001	A	00001	E	00101
R	10010	X	11000	L	01100	I	01001	R	10010	E	00101
T	10100	P	10000	A	00001	A	00001	S	10011	G	00111
I	01001	L	01100	I	01001	R	10010	T	10100	I	01001
N	01110	A	00001	E	00101	S	10011	E	00101	L	01100
G	00111	S	10011	A	00001	T	10100	E	00101	M	01101
E	00101	O	01111	M	01101	L	01100	G	00111	N	01110
X	11000	I	01001	E	00101	E	00101	X	11000	O	01111
A	00001	G	00111	R	10010	M	01101	I	01001	P	10000
M	01101	E	00101	N	01110	E	00101	L	01100	R	10010
P	10000	A	00001	S	10011	N	01110	M	01101	S	10011
L	01100	M	01101	O	01111	O	01111	N	01110	T	10100
E	00101	E	00101	G	00111	G	00111	O	01111	X	11000

6.18

### Linear sorting method

LSD radix sort!

To sort N 64-bit keys take bitsbyte=16

- 4N steps, linear extra memory (plus 2^16)

Does not violate  $N \lg N$  lower bound because

- comparisons are not used

LSD radix sort liabilities

- inner loop has a lot of instructions
- accesses memory "randomly"
- wastes time on low-order bits

Therefore, use just "enough" bits

6.20

## LSD-MSD hybrid

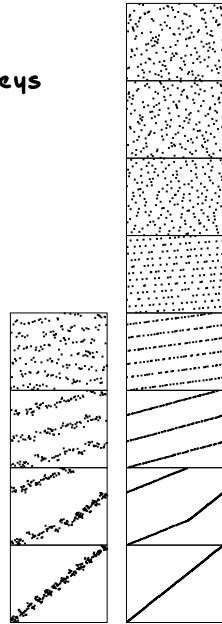
MSD radix sort also linear

Use LSD-MSD hybrid for random keys

- (assume fixed-size keys)
- use  $(\lg N)/2 < \text{bitsbyte} < \lg N$

Three passes

- LSD radix sort on 2nd byte
- LSD radix sort on 1st byte
- insertion sort to clean up



6.21

## Sorting strings

**PROBLEM:**

- long key strings costly to compare when they differ only at the end
- [this is the common case!]

```
absolutism
absolut
absolutely
absolute
```

**SOLUTION:** 3-way radix Quicksort

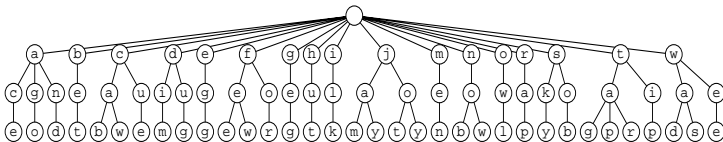
- Use three-way partitioning on key characters
- Recurse and pass current char index

6.23

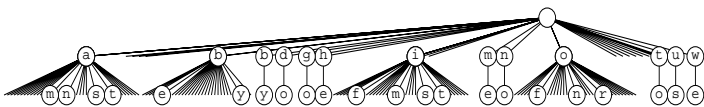
## Recursive structure of MSD radix sort

Tree structure to describe recursive call

Paths in tree give keys



**Problem:** algorithm touches empty nodes



Tree can be as much as M times bigger than they seem

6.22

## 3-way radix Quicksort partitioning

actinian	coenobite	actinian
jeffrey	conelrad	bracteal
coenobite	actinian	coenobite
conelrad	bracteal	conelrad
secureness	secureness	cumin
cumin	dilatedly	chariness
chariness	inkblot	centesimal
bracteal	jeffrey	cankurous
displease	displease	circumflex
millwright	millwright	millwright
repertoire	repertoire	repertoire
dourness	dourness	dourness
centesimal	southeast	southeast
fondler	fondler	fondler
interval	interval	interval
reversionary	reversionary	reversionary
dilatedly	cumin	secureness
inkblot	chariness	dilatedly
southeast	centesimal	inkblot
cankurous	cankurous	jeffrey
circumflex	circumflex	displease

6.24

### 3-way radix Quicksort code

```

#define ch(A) digit(A, D)
void quicksortX(Item a[], int l, int r, int D)
{
    int i, j, k, p, q; int v;
    if (r-l <= M) { insertion(a, l, r); return; }
    v = ch(a[r]); i = l-1; j = r; p = l-1; q = r;
    while (i < j)
    {
        while (ch(a[++i]) < v) ;
        while (v < ch(a[--j])) if (j == l) break;
        if (i > j) break;
        exch(a[i], a[j]);
        if (ch(a[i])==v) { p++; exch(a[p], a[i]); }
        if (v==ch(a[j])) { q--; exch(a[j], a[q]); }
    }
    if (p == q)
        { if (v != '\0') quicksortX(a, l, r, D+1);
          return; }
    if (ch(a[i]) < v) i++;
    for (k = l; k <= p; k++, j--) exch(a[k], a[j]);
    for (k = r; k >= q; k--, i++) exch(a[k], a[i]);
    quicksortX(a, l, j, D);
    if ((i == r) && (ch(a[i]) == v)) i++;
    if (v != '\0') quicksortX(a, j+1, i-1, D+1);
    quicksortX(a, i, r, D);
}

```

6.25

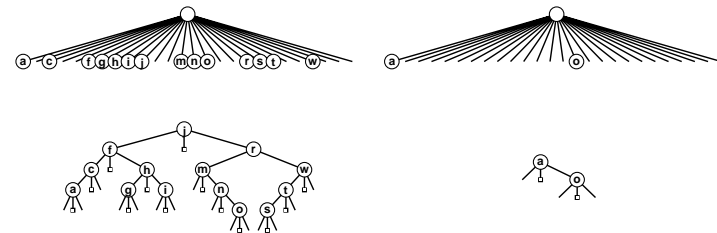
### Another sublinear sort

Three-way radix quicksort is **SUBLINEAR**

$N$  records with  $w$ -byte keys

- Bytes of data:  $Nw$
  - Bytes examined by sort:  $2 N \ln N$
- Ex: 100000 keys, 100 bytes per key
- 10 million bytes of data
  - algorithm examines 2.3 million bytes  
1/5 of the data

Corresponds to collapsing null links in MSD trees



6.27

### 3-way radix Quicksort example

```

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|ab
tag  ago  ago  c|aw
jot  and  and  c|ue
sob  fee  egg  egg
nob  cue  cue  dug
sky  caw  caw  dim
hut  hut  f|ee
ace  ace  f|or
bet  bet  f|ew
men  cab  ilk
egg  egg  gig
few  few  hut
jay  j|ay  ja|m
owl  j|ot  ja|y
joy  j|oy  jo|y
rap  j|am  jo|t
gig  owl  owl  m|en
wee  wee  now  owl
was  was  nob  nob
cab  men  men  now
wad  wad  r|ap
caw  sky  sky  sky  sky
cue  nob  was  tip  sob
fee  sob  sob  sob  t|ip  ta|r
tap  tap  tap  tap  t|ap  ta|p
ago  tag  tag  tag  t|ag  ta|g
tar  tar  tar  tar  t|ar  ti|p
dug  tip  tip  w|as
and  now  wee  w|ee
jam  rap  wad  w|ad

```

6.26