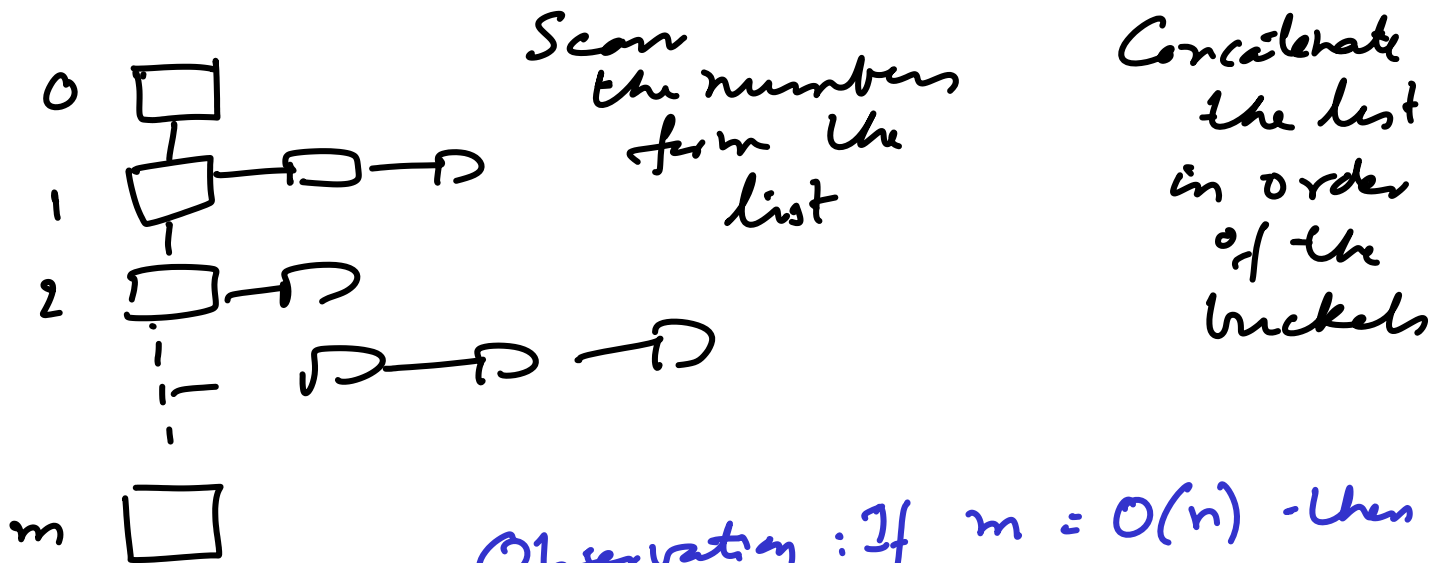


CSL 630 Lecture 6 Aug 11

Sorting strings over an alphabet Σ

n numbers in range $[0..m]$
can be sorted in $O(n+m)$ steps
and $O(n+m)$ space.



Observation: If $m = O(n)$ - then
total time is $O(n)$

But if $m = \omega(n)$

$$\frac{m}{n} \rightarrow \infty$$

$$n \rightarrow \infty$$

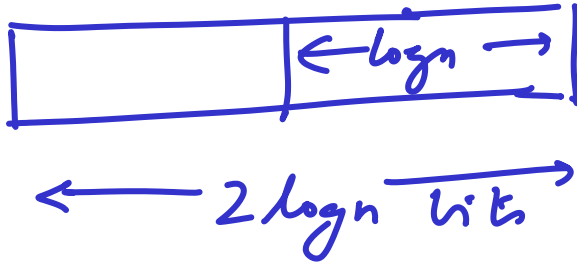
small omega

, then run-time is $\Omega(m)$

To sort numbers in range $[1..n^2]$
we don't apply bucket sort

$$m = n \times n$$

Each number is split into two equal parts, each of size at most



① Then we sort the lower half

② Do a stable sort on the upper half

5 1
2 6
9 5
4 3
1 8
4 5

1st pass →

5 1
4 3
9 5
4 5
2 6
1 8

2nd pass →
Stable

1 8
2 6
4 3
4 5
5 1
9 5

More passes if required

Each pass takes $O(m+n)$ time, so

K passes $\Rightarrow O(K(m+n))$

Since $m = O(n) \Rightarrow O(Kn)$

We can sort n numbers in the range $[1, n^k]$ in $O(k(n))$ steps
As long as k is constant, we can sort in $O(n)$ time

This is not comparisons and therefore we can beat the $\Omega(n \log n)$ bound.

To sort strings that have the same length, the natural choice is Radix sort

Σ : alphabet l : length of each number
 n : strings

If we use radix sort,

$$O\left[\left(\frac{1}{|\Sigma|} + n\right) l\right] = O(nl) \text{ if } |\Sigma| \text{ is } O(1)$$

This is linear-time with input size
Input size $N = nl$ $O(N)$

In lexicographic sorting, when strings are of different lengths, say $l_1, l_2, l_3 \dots l_n$ then we can artificially transform the strings into length l_{max} and then apply radix sort

$$\begin{array}{l}
 d n a \\
 c
 \end{array}
 \Rightarrow
 \begin{array}{l}
 \cancel{d} \cancel{n} \cancel{a} \quad c \text{ --- } \checkmark \\
 \text{---} \cancel{c} \quad d n a
 \end{array}$$

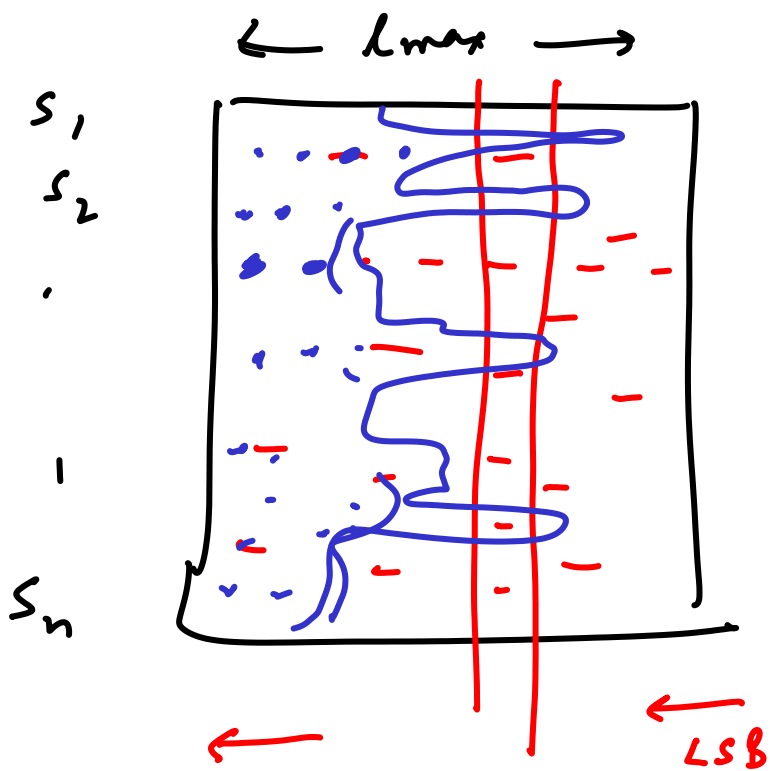
Running time: $O(l_{max} \cdot n)$

Suppose $n-1$ strings of length 1
 1 string of length n

$$N = n-1 + n = 2n-1$$

Time to sort is $n^2!$ Too wasteful

This is the price of "blanks"



When there are blanks, we only need to copy them from the previous pass (stable sort)

We are willing to pay the price for non-blank symbols

Total size of input: $\sum_{i=1}^n l_i = N$

For a specific pass, say j , suppose there are m_j non-blank symbols

$$\sum_j m_j = \sum l_i$$

Goal: Sort in-time in pass j proportional to $O(m_j)$

If we knew that strings involved in the j^{th} pass, then we are done! S^j is the set of j^{th} pass takes time $O(|S^j|) = O(m_j)$

We will some preprocessing
to construct S_j for all j

For each string S_i , let us
construct tuples (j, S_{ij}) where

S_{ij} denotes the j^{th} symbol of string
 S_i . Example c a v e

$(4, e)$ $(3, v)$ $(2, a)$ $(1, c)$

\Rightarrow Total $\sum l_i$ tuples = N

Do a radix sort on the tuples

\Rightarrow we have all the symbols
corresponding to a index clubbed together

Time : First pass : $O(|\Sigma| + N) = O(N)$

Second pass : $(l_{\max} + N) = O(N)$

Overall $O(N)$

c a v e = ¹(1, c) ¹(2, a) ¹(3, v) (4, e)

b a t ²(1, b) ²(2, a) ²(3, t)

a t ³(1, a) ³(2, t)

| | |
|---------------------|--------|
| ¹ (2, a) | (1, a) |
| ² (2, a) | (1, b) |
| ³ (1, a) | (1, c) |
| ² (1, b) | (2, a) |
| ¹ (1, c) | (2, a) |
| ¹ (4, e) | (2, t) |
| . | (3, t) |
| . | (3, v) |
| ² (3, t) | (3, v) |
| ³ (2, t) | (4, e) |
| ¹ (3, v) | |