

15-210

PARALLEL AND SEQUENTIAL
ALGORITHMS AND DATA
STRUCTURES

LECTURE 18

QUICKSORT ANALYSIS AND SORTING LOWER BOUNDS

SYNOPSIS

- Quicksort
- Work and Span Analysis of Randomized Quicksort
- Lower Bound for Comparison-based Sorting
- Lower Bound for Merging

QUICKSORT

- Originally invented and analyzed by Hoare in 1960's.
- I strongly urge to watch Jon Bentley on “Three beautiful Quicksorts” at
 - ▶ www.youtube.com/watch?v=QvgYAQzg1z8.

SEQUENTIAL QUICKSORT

```
int i, j;
for( i = low, j = high - 1; ; )
{
    while( a[ ++i ] < pivot );
    while( pivot < a[ --j ] );
    if( i >= j )
        break;
    swap( a, i, j );
}
// Restore pivot
swap( a, i, high - 1 );
quicksort( a, low, i - 1 ); // Sort small elements
quicksort( a, i + 1, high ); // Sort large elements
```

QUICKSORT

```
1  fun quicksort( $S$ ) =  
2      if  $|S| = 0$  then  $S$   
3      else let  
4           $p = \text{pick a pivot from } S$   
5           $S_1 = \langle s \in S \mid s < p \rangle$   
6           $S_2 = \langle s \in S \mid s = p \rangle$   
7           $S_3 = \langle s \in S \mid s > p \rangle$   
8           $(R_1, R_3) = (\text{quicksort}(S_1) \parallel \text{quicksort}(S_3))$   
9      in  
10          $\text{append}(R_1, \text{append}(S_2, R_3))$   
11     end
```

QUICKSORT

```
1  fun quicksort( $S$ ) =  
2      if  $|S| = 0$  then  $S$   
3      else let  
4           $p = \text{pick a pivot from } S$   
5           $C = \langle s \in S \mid (s, \text{compare}(p, s)) \rangle$   
6           $S_1 = \langle s \mid (s, \text{LESS}) \in C \rangle$   
7           $S_2 = \langle s \mid (s, \text{EQUAL}) \in C \rangle$   
8           $S_3 = \langle s \mid (s, \text{GREATER}) \in C \rangle$   
9           $(R_1, R_3) = (\text{quicksort}(S_1) \parallel \text{quicksort}(S_3))$   
10         in  
11          $\text{append}(R_1, \text{append}(S_2, R_3))$   
12     end
```

QUICKSORT

- Each call to Quicksort either makes
 - ▶ No recursive calls (base case), or
 - ▶ Two recursive calls
- Call tree is a binary
- Depth the call tree determines the span of the algorithm.

PICKING THE PIVOT

- Always pick the first element
 - ▶ Worst case $O(n^2)$ work.
 - ▶ In practice, **almost sorted inputs are not uncommon.**
- Pick the median of 3 elements (e.g., first, middle and last elements)
 - ▶ could possible divide evenly
 - ▶ worst case is still bad
- **Pick an element at random**
 - ▶ we hope this divides evenly in expectation
 - ▶ leading to expected $O(n \log n)$ work and $O(\log^2 n)$ span.

PICKING THE PIVOT

- Pick first element
 - ▶ Worst case $O(n^2)$ work.
 - ▶ Expected $O(n \log n)$ work
 - ★ Averaged over all possible orderings.
 - ▶ Works well on the average
 - ▶ Slow on some, possibly common, cases.
- Pick a random element
 - ▶ Expected worst-case $O(n \log n)$ work.
 - ★ For input in **any** order, the expected work is $O(n \log n)$
 - ▶ No input has expected $O(n^2)$ work.
 - ▶ With a small probability, we could be unlucky and have $O(n^2)$ work.

RANDOMIZED QUICKSORT

- Assign a uniformly random priority to each number in $[0, 1]$.

```
1 fun quicksort(S) =  
2   if |S| = 0 then S  
3   else let  
4       p = pick as pivot the highest priority element from S  
5        $S_1 = \{s \in S \mid s < p\}$   
6        $S_2 = \{s \in S \mid s = p\}$   
7        $S_3 = \{s \in S \mid s > p\}$   
8        $(R_1, R_3) = (\text{quicksort}(S_1) \parallel \text{quicksort}(S_3))$   
9       in  
10      append( $R_1$ , append( $S_2$ ,  $R_3$ ))  
11  end
```

- Once the priorities are assigned, the algorithm is deterministic.

RANDOMIZED QUICKSORT

- Count comparisons made!
 - ▶ Almost all the work is comparisons.
$$X_n = \# \text{ of comparisons } \textit{quicksort} \text{ makes on input of size } n$$
- Find $\mathbf{E}[X_n]$ for any input sequence S
- Notation:
 - ▶ Let $T = \textit{sort}(S)$
 - ▶ T_i and T_j refer to elements in the final sorted order and $i < j$ and $T_i \leq T_j$.
 - ▶ p_i refers to priority chosen for T_i .
 - ▶ $A_{i,j} = 1$ if T_i and T_j were ever compared during the sort.

ANALYZING QUICKSORT

- Crucial point is how to model $A_{i,j}$.
- In any one call to `quicksort`, there are three cases:
 - ▶ Pivot p is either T_i or $T_j \Rightarrow A_{i,j} = 1$
 - ▶ $T_i < p < T_j \Rightarrow T_i \in S_1, T_j \in S_3, A_{i,j} = 0$
 - ▶ Either $p < T_i$ or $p > T_j \Rightarrow T_i, T_j \in S_1$ or $T_i, T_j \in S_3$
- If two elements are compared in a `quicksort` call, they will **never** be compared again in any other call!

ANALYZING QUICKSORT

$$X_n \leq \sum_{i=1}^{n-1} \sum_{j=i+1}^n A_{ij}$$

- The non-optimized code compares each element to pivot 3 times, **the optimized version compares once.**

1 ...
2 ...
3 ...

$$C = \langle s \in S \mid (s, \text{compare}(p, s)) \rangle$$

- By linearity of expectation

$$\mathbf{E}[X_n] \leq \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbf{E}[A_{ij}] = \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbf{Pr}[A_{ij} = 1]$$

ANALYZING QUICKSORT

- Consider first when the pivot is one of T_i, T_{i+1}, \dots, T_j
- T_i and T_j are compared $\Leftrightarrow p_i$ or p_j is the highest priority among $\{p_i, p_{i+1}, \dots, p_j\}$.
 - ▶ Assume $T_k, i < k < j$ has higher priority.
 - ▶ For any subdivision $\dots, T_i, \dots, T_k, \dots, T_j$, T_k will become a pivot and separate T_i and T_j
 - ▶ T_i and T_j will never be compared!

ANALYZING QUICKSORT

$$\begin{aligned}\mathbf{E}[A_{ij}] &= \mathbf{Pr}[A_{ij} = 1] \\ &= \mathbf{Pr}[p_i \text{ or } p_j \text{ is the maximum in } \{p_i, \dots, p_j\}] \\ &= \frac{2}{j-i+1} \text{ (Why ?)}\end{aligned}$$

- $j - i + 1$ elements between p_i and p_j and each is equally likely to be the maximum.
- We want either p_i or p_j , hence $\frac{2}{j-i+1}$
- T_i is compared to T_{i+1} with probability 1.

ANALYZING QUICKSORT

$$\begin{aligned}\mathbf{E}[X_n] &\leq \sum_{i=1}^{n-1} \sum_{j=i+1}^n \mathbf{E}[A_{ij}] \\&= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} \\&= \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} \quad (\text{change variables}) \\&\leq 2 \sum_{i=1}^n H_n \\&\leq 2 \cdot n \cdot H_n \in O(n \log n)\end{aligned}$$

ANALYZING QUICKSORT

- Indirectly, average work for basic deterministic quicksort is $O(n \log n)$.
 - ▶ Just shuffle data randomly and apply the basic algorithm
 - ▶ \equiv to picking random priorities

ALTERNATIVE ANALYSIS

- Write a recurrence for the number of comparisons:

$$X(n) = X(Y_n) + X(n - Y_n - 1) + n - 1$$

- Random variable Y_n is the size of S_1 .

$$\begin{aligned}\mathbf{E}[X(n)] &= \mathbf{E}[X(Y_n) + X(n - Y_n - 1) + n - 1] \\ &= \mathbf{E}[X(Y_n)] + \mathbf{E}[X(n - Y_n - 1)] + n - 1 \\ &= \frac{1}{n} \sum_{i=0}^{n-1} (\mathbf{E}[X(i)] + \mathbf{E}[X(n - i - 1)]) + n - 1\end{aligned}$$

ALTERNATIVE ANALYSIS

$$\begin{aligned}\mathbf{E}[X(n)] &= \frac{1}{n} \sum_{i=0}^{n-1} (\mathbf{E}[X(i)] + \mathbf{E}[X(n-i-1)]) + n - 1 \\ &= \frac{2}{n} \sum_{i=0}^{n-1} \mathbf{E}[X(i)] + n - 1\end{aligned}$$

- With telescoping, this also solves as $O(n \log n)$

EXPECTED SPAN

- S is split into $L(ess)$, $E(qual)$ and $(g)R(eater)$.
- Let $X_n = \max\{|L|, |R|\}$,
- We use `filter` to partition.

$$S(n) = S(X_n) + O(\log n)$$

EXPECTED SPAN

- Let $\overline{S}(n)$ denote $\mathbf{E}[S(n)]$
- We bound $\overline{S}(n)$ by considering $\mathbf{Pr}[X_n \leq 3n/4]$ and $\mathbf{Pr}[X_n > 3n/4]$.
- $\mathbf{Pr}[X_n \leq 3n/4] = 1/2$
 - ▶ As with `SmallestK`, 1/2 of the randomly chosen pivots results in larger partition of at most size $3n/4$ elements.

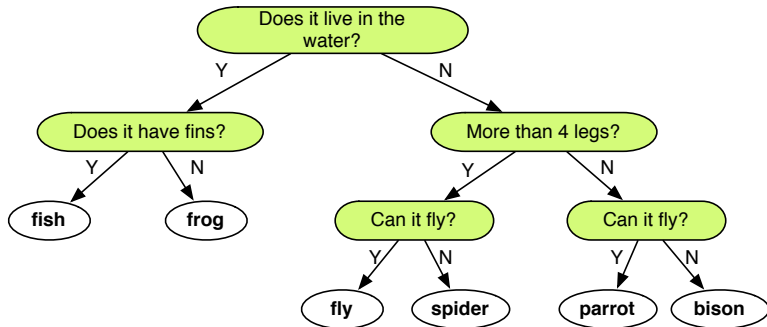
EXPECTED SPAN

$$\begin{aligned}\overline{S}(n) &= \sum_i \mathbf{Pr}[X_n = i] \cdot \overline{S}(i) + c \log n \\ &\leq \mathbf{Pr}[X_n \leq \tfrac{3n}{4}] \overline{S}(\tfrac{3n}{4}) + \mathbf{Pr}[X_n > \tfrac{3n}{4}] \overline{S}(n) + c \cdot \log n \\ &\leq \tfrac{1}{2} \overline{S}(\tfrac{3n}{4}) + \tfrac{1}{2} \overline{S}(n) + c \cdot \log n \\ &\implies (1 - \tfrac{1}{2}) \overline{S}(n) \leq \tfrac{1}{2} \overline{S}(\tfrac{3n}{4}) + c \log n \\ &\implies \overline{S}(n) \leq \overline{S}(\tfrac{3n}{4}) + 2c \log n \\ &\implies \overline{S}(n) \in O(\log^2 n)\end{aligned}$$

LOWER BOUND FOR SORTING

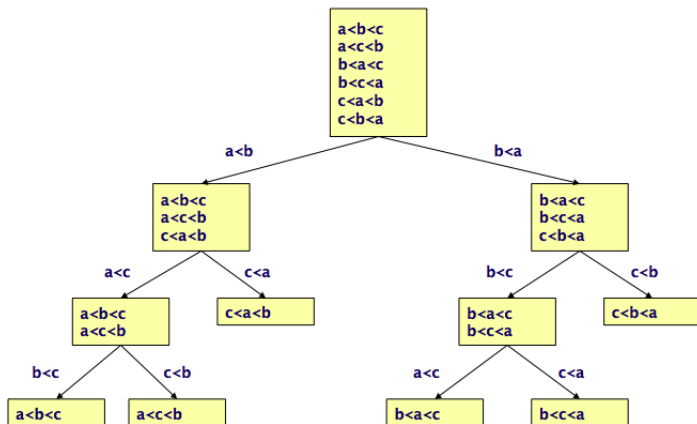
- What is asymptotically the minimum number comparisons any sorting algorithm has to make?
- Lower-bounds apply to problems not to algorithms.
 - ▶ Algorithms provide upper bounds!
- We say **sorting is $\Omega(n \log n)$**
- **No (comparison-based) sorting algorithm has work asymptotically lower than $n \log n$.**

DECISION TREES



- If there are N outcomes, the number of questions is at least $\log_2 N$.

SORTING AS A DECISION PROBLEM



- For n items, how many possible outcomes can there be?
 - $n! \Rightarrow$ we need at least $\log_2(n!)$ “questions”.

SORTING AS A DECISION PROBLEM

$$\begin{aligned}\log(n!) &= \log n + \log(n-1) + \cdots + \log(n/2) + \cdots + \log 1 \\ &\geq \log n + \log(n-1) + \cdots + \log(n/2) \\ &\geq \frac{n}{2} \cdot \log(n/2) \in \Omega(n \log n)\end{aligned}$$

LOWER BOUND FOR MERGING

- We have sorted sequences A , $|A| = n$ and B , $|B| = m$ and $m \leq n$.
 - ▶ Assume all elements are unique.
- All interleavings are possible
- We need to **choose m positions out of $n + m$** to place the elements of B amongst elements of A .
- Finding the right sequence of m positions can be done with at least $\log_2 \binom{n+m}{m}$ comparisons.

LOWER BOUND FOR MERGING

- $\binom{n}{r} \geq \left(\frac{n}{r}\right)^r$
 - ▶ See Lemma in the notes.

$$\log_2 \binom{n+m}{m} \geq \log_2 \left(\frac{n+m}{m}\right)^m = m \log_2 \left(1 + \frac{n}{m}\right)$$