

CS4800: Algorithms & Data

Jonathan Ullman

Lecture 6:

- Divide-and-Conquer: Inversions, Closest Pair
- Dynamic Programming Warmup (?)

Jan 26, 2018

Counting Inversions

Approximate Sortedness

- Which of these is “more sorted”?

8	3	28	11	17	42
---	---	----	----	----	----

11	3	8	17	28	42
----	---	---	----	----	----

- only one elem out of place
- 5/6 of elems are sorted

Counting Inversions



- A Measure of Sortedness: Number of Inversions

- Inversion: a pair $i < j$ such that $A[i] > A[j]$ (a sorted list has no inversions)
- “Kendall tau distance” / “Bubble sort distance”

8	3	28	11	17	42
1	2	3	4	5	6

11	3	8	17	28	42
1	2	3	4	5	6

Inversions:

(1, 2)

(3, 4)

(3, 5)

3

Inversions:

(1, 2)

(1, 3)

2

Counting Inversions

- A Measure of Sortedness: Number of Inversions
 - Inversion: a pair $i < j$ such that $A[i] > A[j]$
 - “Kendall tau distance”

8	3	28	11	17	42
---	---	----	----	----	----

11	3	8	17	28	42
----	---	---	----	----	----

- Many Applications
 - Collaborative filtering / recommender systems
 - Social choice theory / voting theory

Counting Inversions: Basic Algorithm

8	3	28	17	11	42	2	35
---	---	----	----	----	----	---	----

For all i, j s.t. $i < j$:

check if $A[j] > A[i]$

Output # of inversions

of possible pairs

$$\frac{n(n+1)}{2} = \Theta(n^2)$$

Counting Inversions: D&C

$L, c_L = 2$

Idea: Sort and Count at the same time!

$R, c_R = 3$

8	3	28	17	11	42	2	35
---	---	----	----	----	----	---	----

A

Step 1: Cut list in half

Step 2: $c_L \leftarrow \text{CountInversions}(L)$ $c_R \leftarrow \text{CountInversions}(R)$

Three types of inversions:

① $1 \leq i < j \leq \frac{n}{2}$ ("on the left") $\boxed{c_L}$

② $\frac{n}{2} + 1 \leq i < j \leq n$ ("on the right") $\boxed{c_R}$

③ $1 \leq i \leq \frac{n}{2}$ $\frac{n}{2} + 1 \leq j \leq n$ ("in the middle") $\boxed{c_M}$

$$T(n) = \mathcal{O}(n^2)$$

Compare $\frac{n}{2} \cdot \frac{n}{2}$

items $\frac{n^2}{4} = \mathcal{O}(n^2)$



Step 3: Try all inversions of type ③ and count them

Step 4: return $c_L + c_R + c_M$

Counting Inversions: D&C

$L, c_L = 2$

$R, c_R = 3$



A

Goal: Count inversions
 i, j s.t. $i \in L, j \in R$.

↑ smallest elem on the right



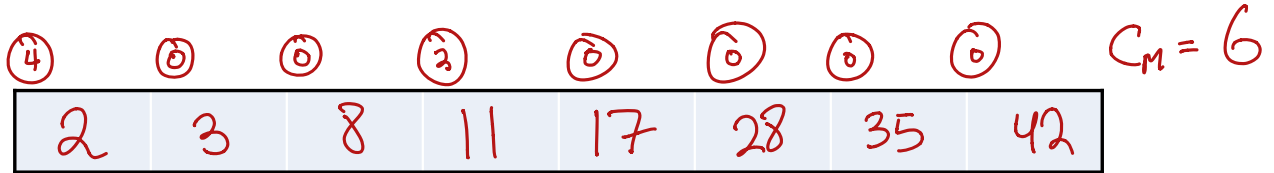
of times an elem of R
 is smaller than an
 elem of L.



L sorted



R sorted



When we moved 2 into place, it was smaller than all elems

Note: Be careful with ties. in L that had not been placed.

Counting Inversions: D&C

8	3	28	17	11	42	2	35
---	---	----	----	----	----	---	----

Claim: Given two sorted lists (L, R) w/ # of inversions (c_L, c_R) , we can produce one sorted list A and compute c_A in $O(n)$ time.

3	8	17	28
---	---	----	----

2	11	35	42
---	----	----	----

--	--	--	--	--	--	--	--

Note: Be careful with ties. Make sure to add L if there is a tie.

Counting Inversions: D&C

MergeAndCount($L[1 \dots \ell], R[1, \dots, r]$):

Let $i, j, k \leftarrow 1, c \leftarrow 0$

time
 $O(\ell) + O(\ell+r)$
 $= O(\ell+r)$

{ For $k = 1, \dots, \ell + r$ // Loop over elts

 If $i > \ell$: // L is empty
 $A[k] = R[j], j \leftarrow j + 1$

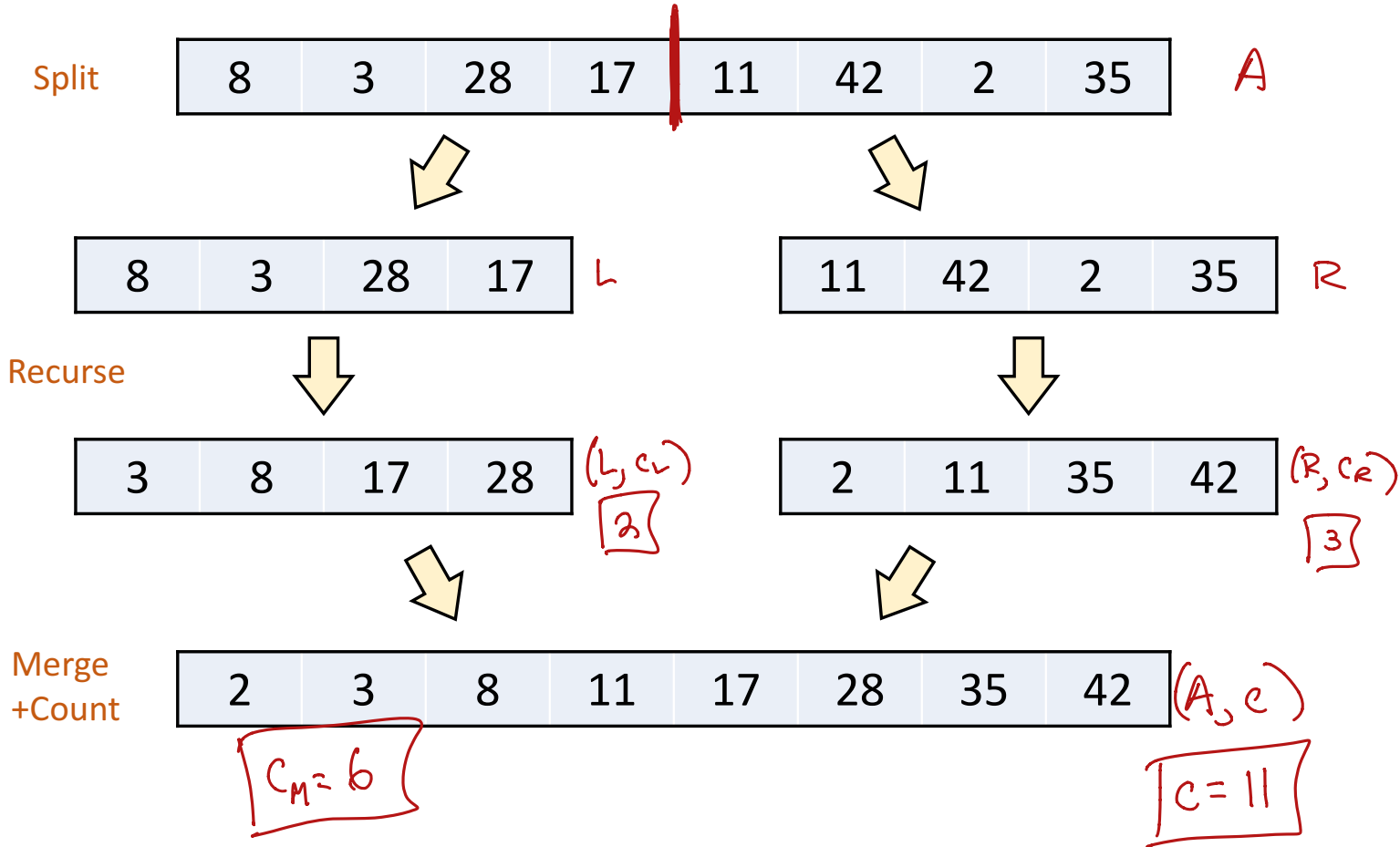
 Elif $j > r$: // R is empty
 $A[k] = L[i], i \leftarrow i + 1$

 Elif $L[i] < R[j]$: // L is smaller
 $A[k] = L[i], i \leftarrow i + 1$

 Else: // R is smaller
 $A[k] = R[j], j \leftarrow j + 1,$
 $c \leftarrow c + (\ell - i + 1)$ // Add inversions

Return (A, c)

Counting Inversions: D&C



Counting Inversions: D&C

SortAndCount($A[1, \dots, n]$):

If $n = 1$, return $(A, 0)$ // Base case

$\ell \leftarrow \lfloor \frac{n}{2} \rfloor$ // Split into two lists
 $L \leftarrow A[1, \dots, \ell], R \leftarrow A[\ell + 1, n]$

$(L, c_L) \leftarrow \text{SortAndCount}(L) \rightarrow T(\frac{n}{2})$ // Recurse

$(R, c_R) \leftarrow \text{SortAndCount}(R) \rightarrow T(\frac{n}{2})$ // Recurse

$(A, c_M) \leftarrow \text{MergeAndCount}(L, R)$ // Merge
 $\rightarrow O(n)$

return $(A, c_L + c_R + c_M)$

$$T(n) = 2 \times T\left(\frac{n}{2}\right) + O(n) \quad T(n) = \Theta(n \log n)$$

Ask the Audience!

- Suppose I wanted to output a list of all inversions (i.e. all pairs $i < j$ such that $A[i] > A[j]$), can our D&C algorithm be modified to output a list of all inversions in $O(n \log n)$ time?

No. There may actually be $\Omega(n^2)$ inversions so we cannot possibly write them down.

Counting Inversions: Fun Fact

- Improved to $O(n\sqrt{\log n})$ in 2010
 - It is “not possible” to sort in less than $n \log n$ time
 - Counting inversions is easier than sorting

Closest Pair

Closest Pair

- General Problem: Given points p_1, \dots, p_n , find the pair p_i, p_j minimizing $d(p_i, p_j)$
 - Foundational problem in computational geometry
 - Closely related to nearest-neighbor search, classification, and clustering problems

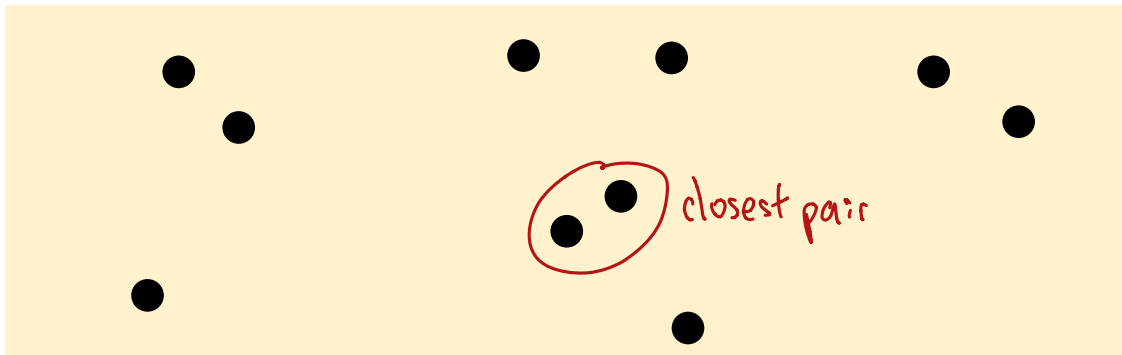
Closest Pair in 2D

Naïve Algorithm: $\Theta(n^2)$ time
(try all pairs)

- Special Case: Closest Pair in 2D
- Given $p_1 = (x_1, y_1), \dots, p_n = (x_n, y_n) \in \mathbb{R}^2$ find the pair p_i, p_j minimizing the Euclidean distance

Assume you can obtain exact distances in $O(1)$ time

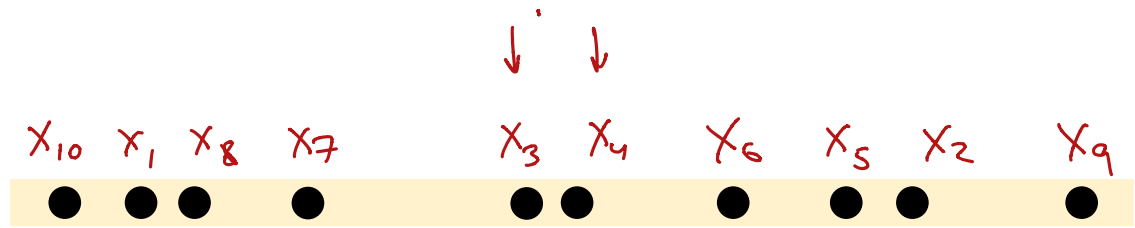
$$d(p_i, p_j) = \sqrt{(x_i - x_j)^2 + (y_i - y_j)^2}$$



Assumption: all the x and y coordinates distinct

Ask the Audience!

- Even More Special Case: Closest Pair in 1D
- Given $x_1, \dots, x_n \in \mathbb{R}$ find the pair x_i, x_j minimizing the distance $|x_i - x_j|$
- Find an $O(n \log n)$ time algorithm



Alg I: • Sort the points

• Consider all pairs $A[i], A[i+1]$ in the sorted list

(Fact: Closest pair i, j is in consecutive places in the sorted order.)

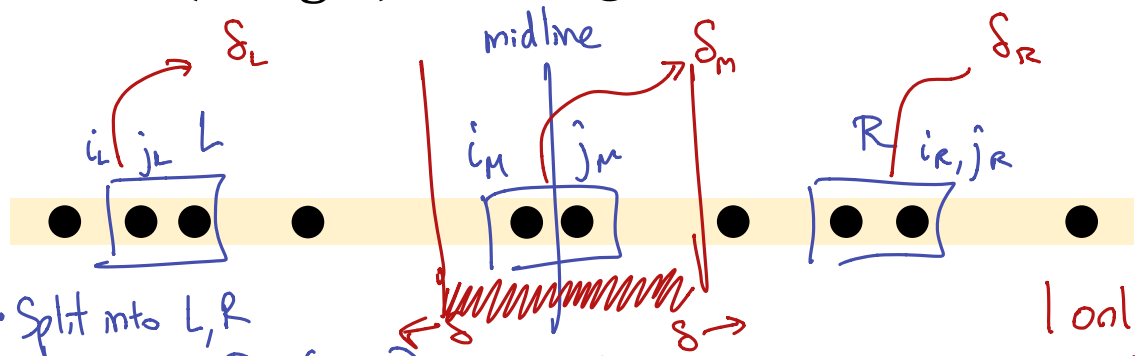
$O(n \log n)$

$n-1$ choices

Ask the Audience!

- Even More Special Case: Closest Pair in 1D
- Given $x_1, \dots, x_n \in \mathbb{R}$ find the pair x_i, x_j minimizing the distance $|x_i - x_j|$
- Find an $O(n \log n)$ time algorithm

Let $\delta = \min\{\delta_L, \delta_R\}$



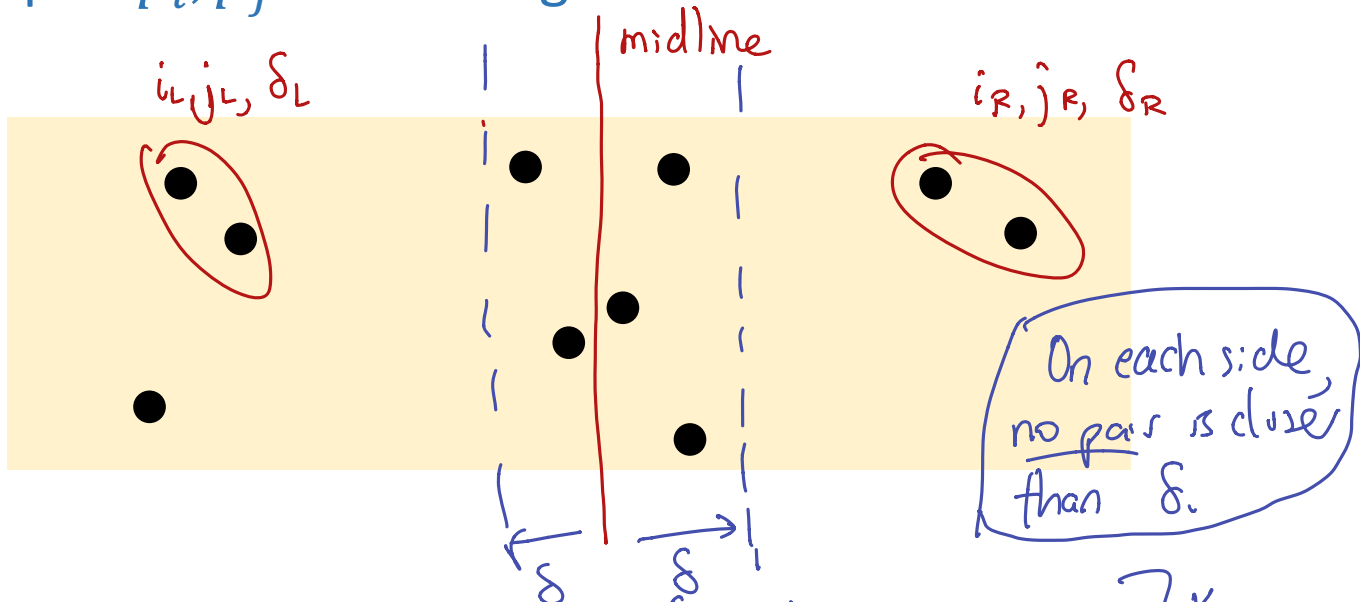
- Alg II:
- Split into L, R
 - Find $(i_L, j_L), (i_R, j_R)$ recursively
 - Consider rightmost i_n in L, leftmost point j_m in R
 - Choose the best of δ_L, δ_R

I only need to consider i_n, j_m at distance $\leq \delta$ from midline.

Closest Pair in 2D

$$\delta = \min \{ \delta_L, \delta_R \}$$

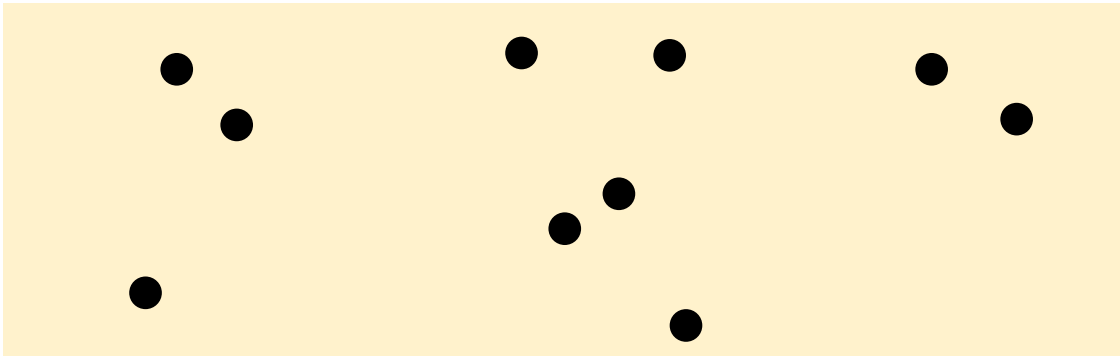
- Given $p_1 = (x_1, y_1), \dots, p_n = (x_n, y_n) \in \mathbb{R}^2$ find the pair p_i, p_j minimizing the Euclidean distance



The only candidates I consider for closest pair are } Key Idea
in the band around the midline.

Closest Pair in 2D

- Given $p_1 = (x_1, y_1), \dots, p_n = (x_n, y_n) \in \mathbb{R}^2$ find the pair p_i, p_j minimizing the Euclidean distance



- Key Idea: Only need to consider a band of width δ on each side around the midline

Closest Pair in 2D

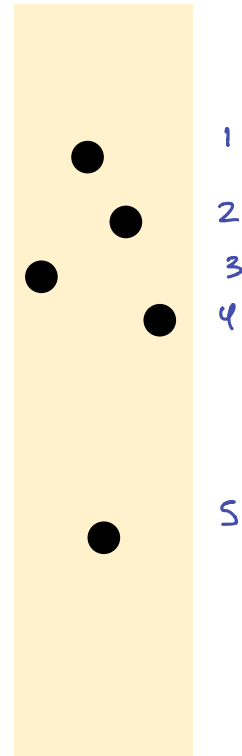
1D: Travis Barker Mohawk

• Sort by y-coord
and look at adjacent
pairs.



2D: Mr. T Mohawk

Idea: Sort by y,
look at adjacent
pairs.

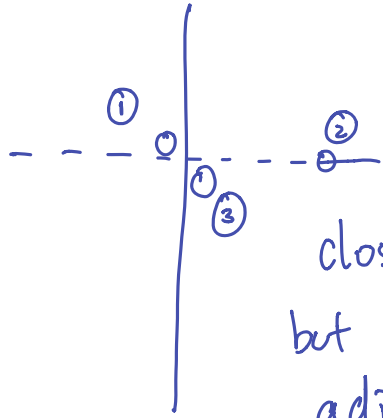


Ask the Audience!

I pity the fool who gets this wrong!

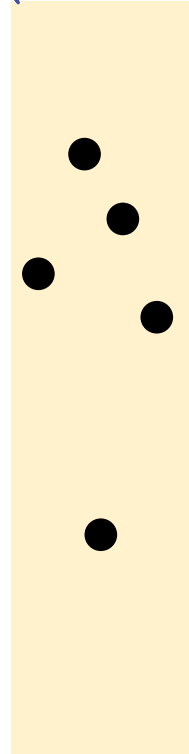


- Suppose I look at the points in the mohawk and sort by the y-coord.
- True/False: The closest pair in the mohawk is **adjacent** in the sorted list?



closest pair is ①, ③
but those are not adjacent

δ ← → δ



No pair in the left/right half is closer than δ

Closest Pair in 2D

→ Very cool, key claim

- Claim: Let $(x_1, y_1), \dots, (x_m, y_m)$ be the points in the mohawk, sorted by y -coordinate and let p_k, p_ℓ be the closest pair, then $|k - \ell| \leq 15$

There are $\leq 15m$ such pairs where $m \leq n$ is the number of points in the mohawk.

Closest Pair in 2D

Closest(P , S_X, S_Y)

Sort all points by x , S_X
Sort all points by y , S_Y

Let q be the middle-element of S_X

Divide P into Left, Right according to q . Scan to get L_Y, R_Y .

Split by x -coord

$\delta_{L,R,j} = \text{MIN}(\text{Closest}(\text{Left}, L_X, L_Y), \text{Closest}(\text{Right}, R_X, R_Y))$

$\delta = \min\{\delta_L, \delta_R\}$
 $(i_L, j_L), \delta_L$ $(i_R, j_R), \delta_R$

Mohawk = { Scan S_Y , add pts that are delta from $q.x$ }

Shave the mohawk
Sorted by y coord

For each point x in Mohawk (in order):

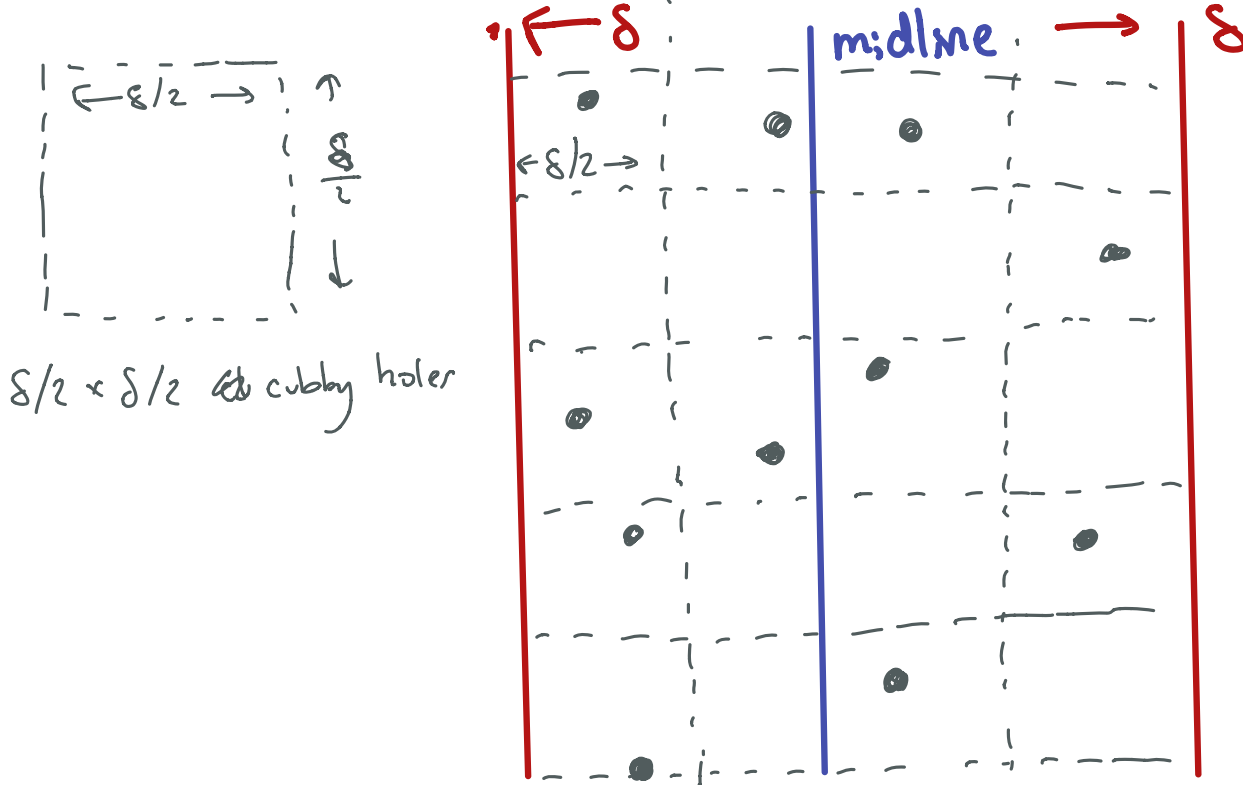
[Compute distance to its next 15 neighbors] only check ≤ 15 other points
Update $\delta_{L,R,j}$ if any pair (x,y) is $< \delta$

Return $(\delta_{L,R,j})$

Best among $(i_L, j_L), (i_R, j_R), \text{Mohawk}$.

Closest Pair in 2D

- Claim: Let $(x_1, y_1), \dots, (x_m, y_m)$ be the points in the mohawk, sorted by y -coordinate and let p_k, p_ℓ be the closest pair, then $|k - \ell| \leq 15$

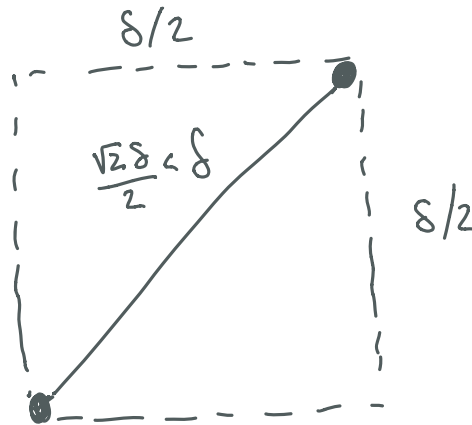


Closest Pair in 2D

- Claim: Let $(x_1, y_1), \dots, (x_m, y_m)$ be the points in the mohawk, sorted by y -coordinate and let p_k, p_ℓ be the closest pair, then $|k - \ell| \leq 15$

Fact ①: At most one point in each cubby hole

- Because no pair on the same side of the midline have distance $< \delta$.



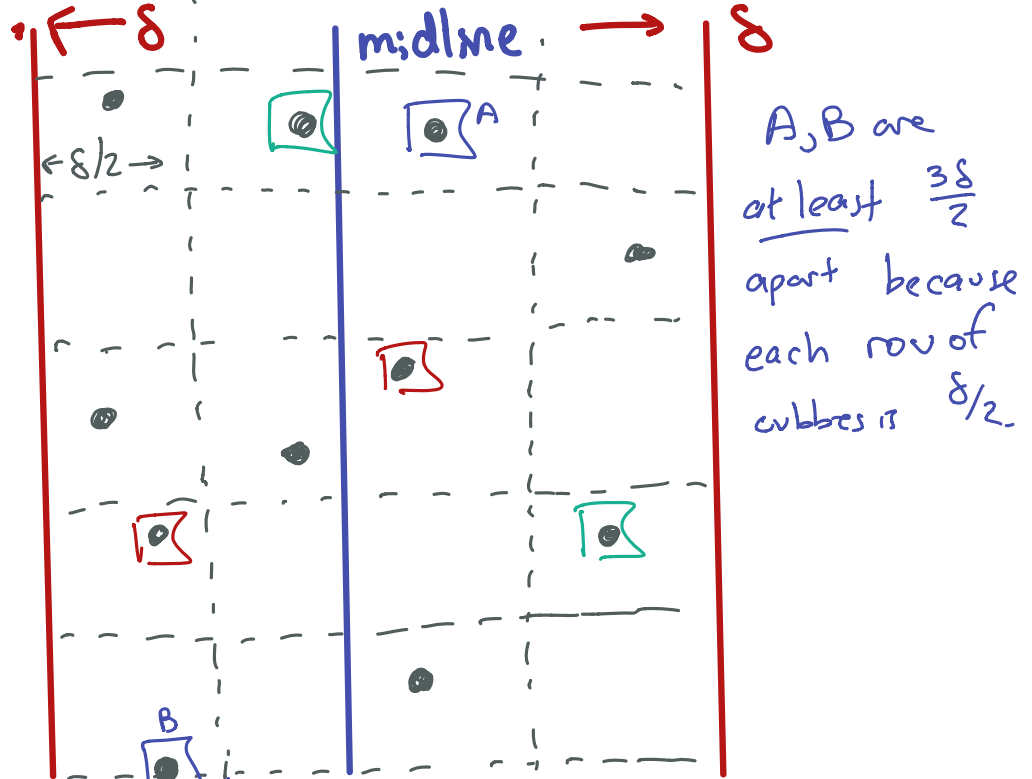
Closest Pair in 2D

- Claim: Let $(x_1, y_1), \dots, (x_m, y_m)$ be the points in the mohawk, sorted by y -coordinate and let p_k, p_ℓ be the closest pair, then $|k - \ell| \leq 15$

Fact ②:

Any pair closer than δ can only be separated by 2 rows of cubbies

Fact ③: If A, B are closest, then only the 15 points in cubbies in the rows btw can be btw them in sorted order.



A, B are at least $\frac{3\delta}{2}$ apart because each row of cubbies is $\delta/2$.