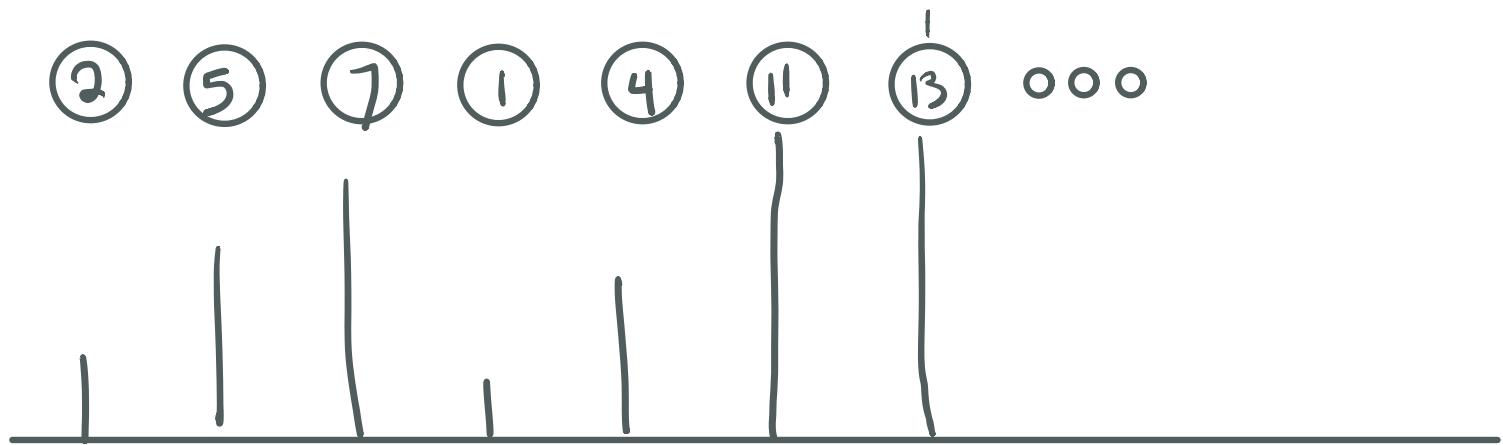


Space efficient quantile selection

Input: stream $a_1, \dots, a_n \in U$ where U has order $<$



e.g. numerical data
names w/ alphabetic order
grades

allowed multiple passes

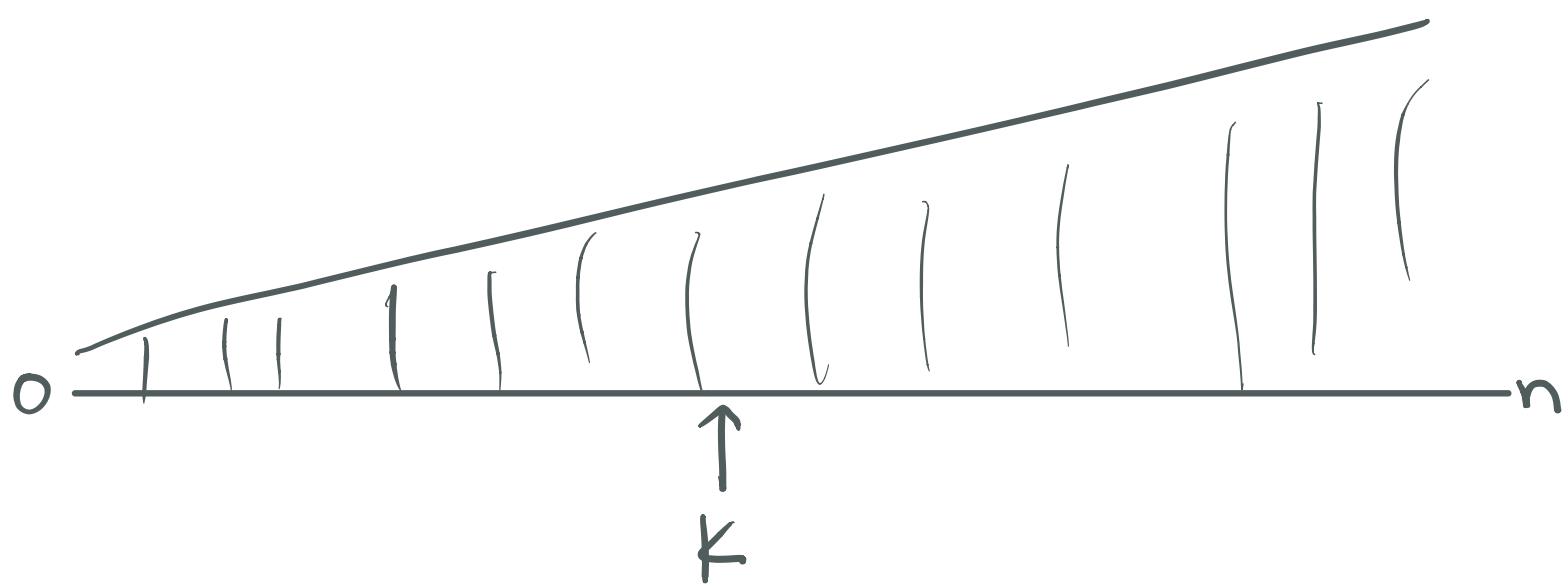
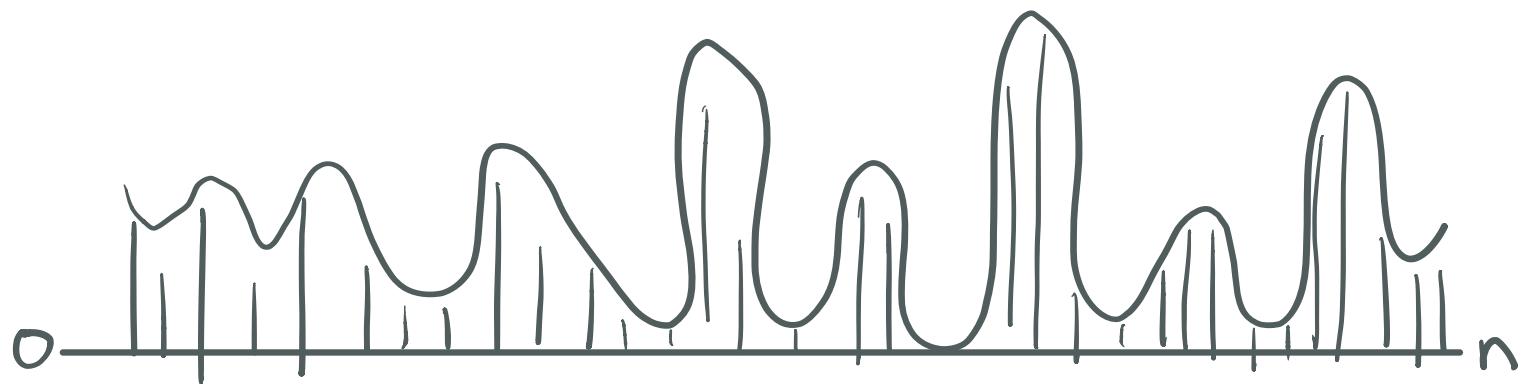
Goal: return the median w/ minimum:

(a) # passes (b) space

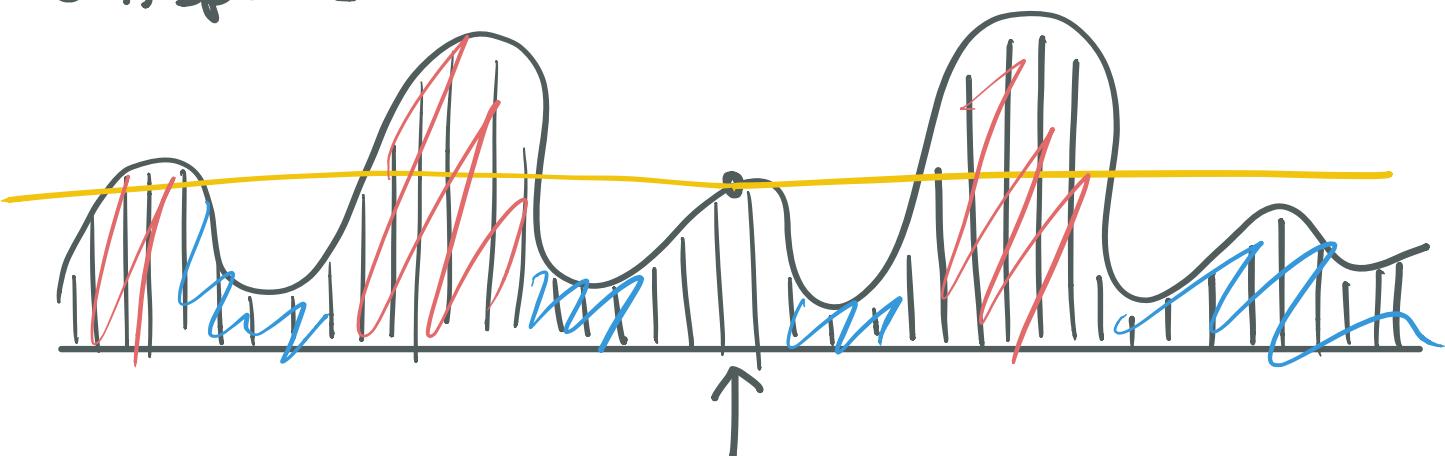
more generally: "quantile queries"

select rank k element
(k th largest)

1 pass: Input

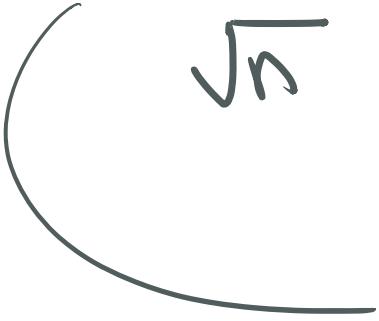


$O(1)$ space



<u>Passes</u>	<u>Space</u>	
1	$O(n)$	sort and select
$\alpha(\log n)$	$O(1)$	quickselect (random pivot)
P	$\tilde{O}(n^{1/p})$	Munro Paterson 1981

2 \sqrt{n}

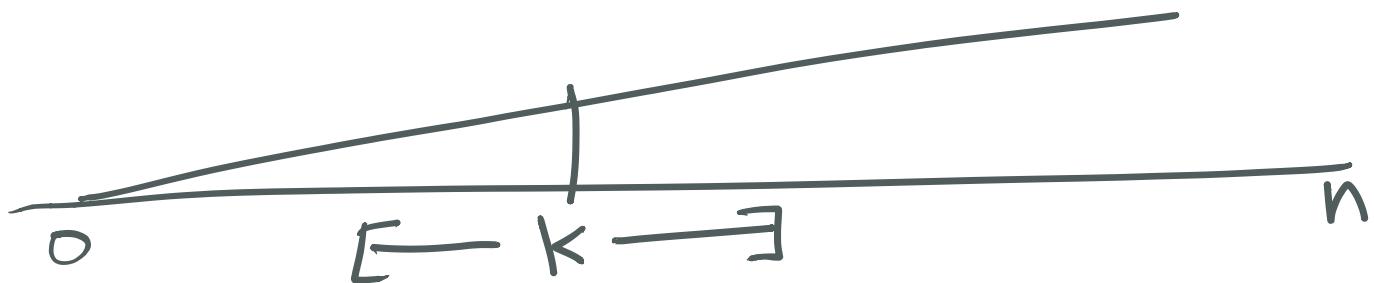


quantile summaries

Approximations

given rank $k \in [n]$ and param $\epsilon > 0$,

return element w/ rank $k \pm \epsilon n$



Sampling:

for median:

sample $l = O\left(\frac{\log \frac{1}{\delta}}{\epsilon^2}\right)$ elements

return median of sample

for rank $k = dn$

sample l

return rank dl out of the sample

Deterministic?

Quantiles

space-efficient

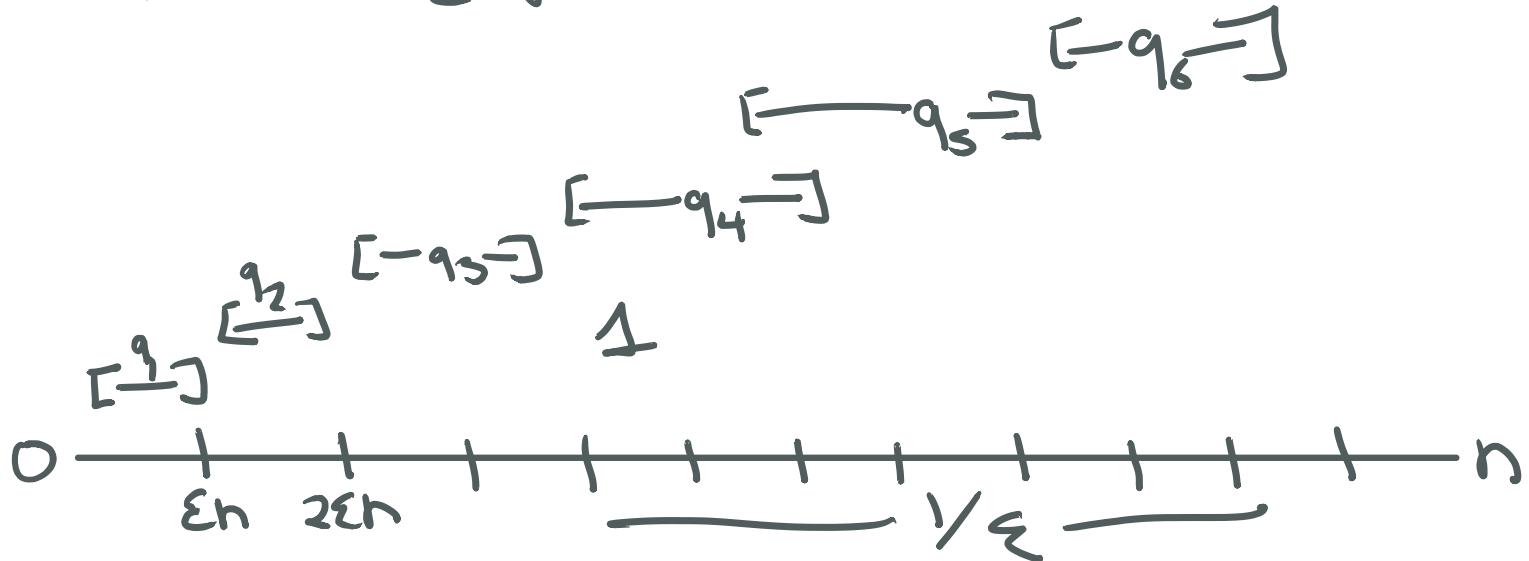
mergable

answer ϵ -approximate quantile queries

ℓ elements $q_1 < q_2 < \dots < q_\ell \in S$

along w/ intervals $I(q_i)$ $\text{rank}(q_i) \in I(q_i)$

need $\geq \frac{1}{\epsilon}$ points



By tracking min/max specially, assume

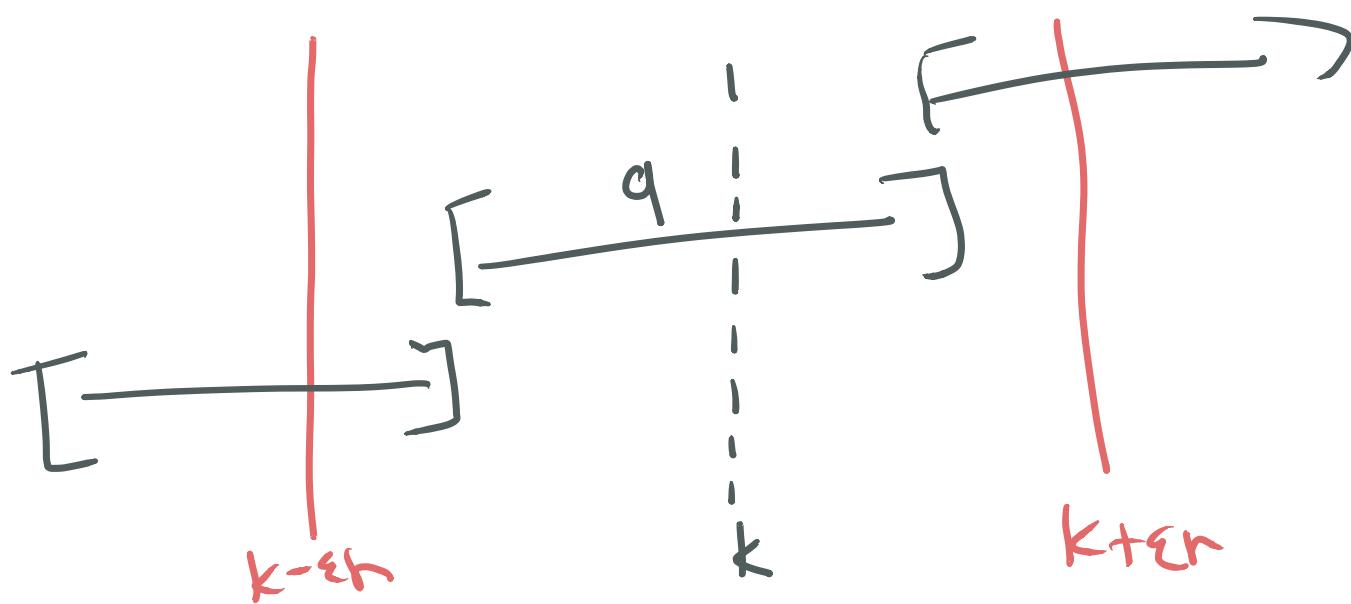
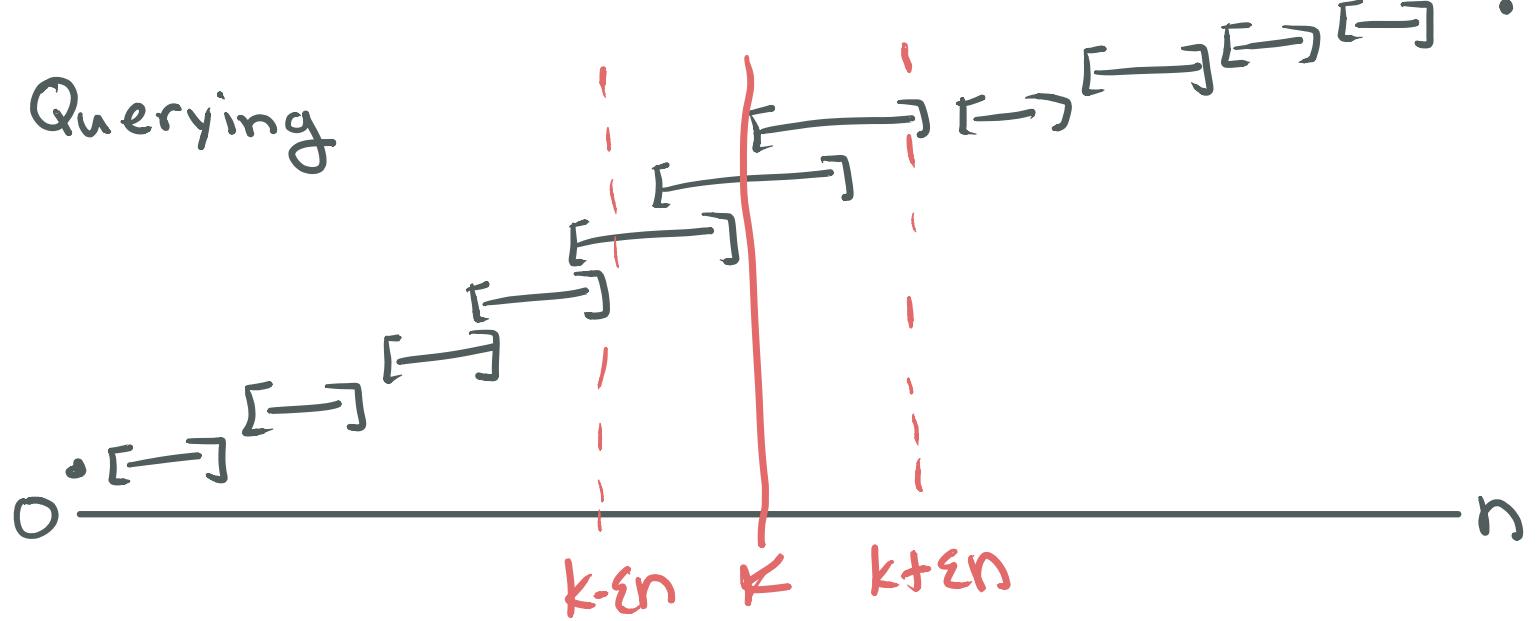
rank(q_1) = 1, rank(q_2) = n.

[] [] [] []

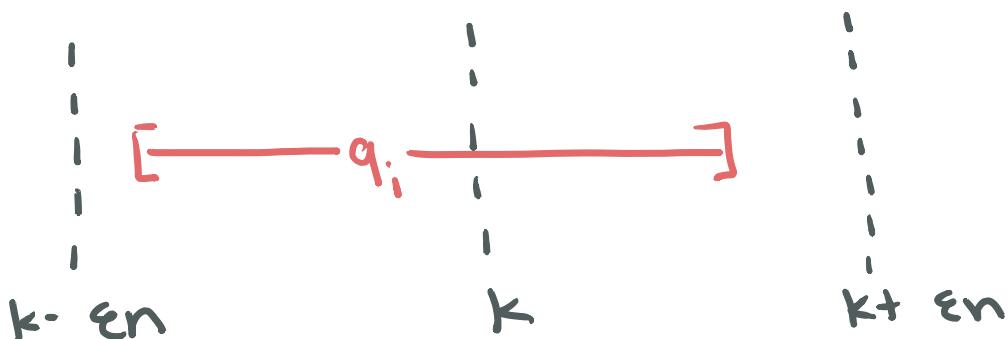
q! [] [] [] [] [] [] []

$$O \longrightarrow n$$

Querying

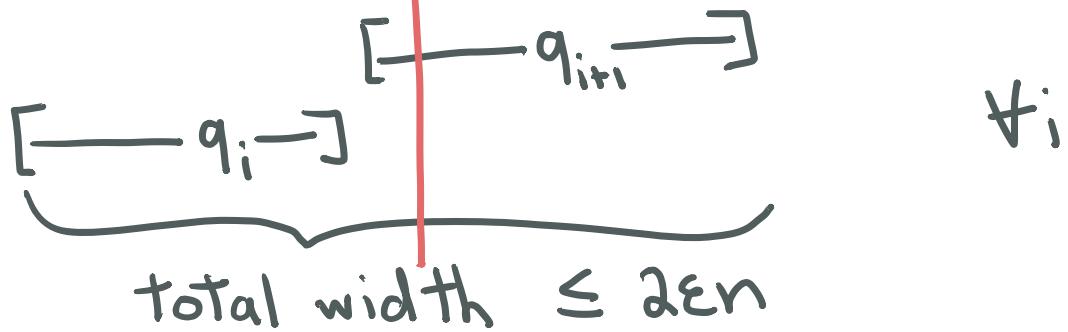


If $I(q_i) \subseteq [k - \epsilon n, k + \epsilon n]$ for some q_i ,
then return q_i .



how to ensure such q_i exists $\neq k$?

lemma



Then every query k contains an interval $I(q_i)$.

Proof two cases:

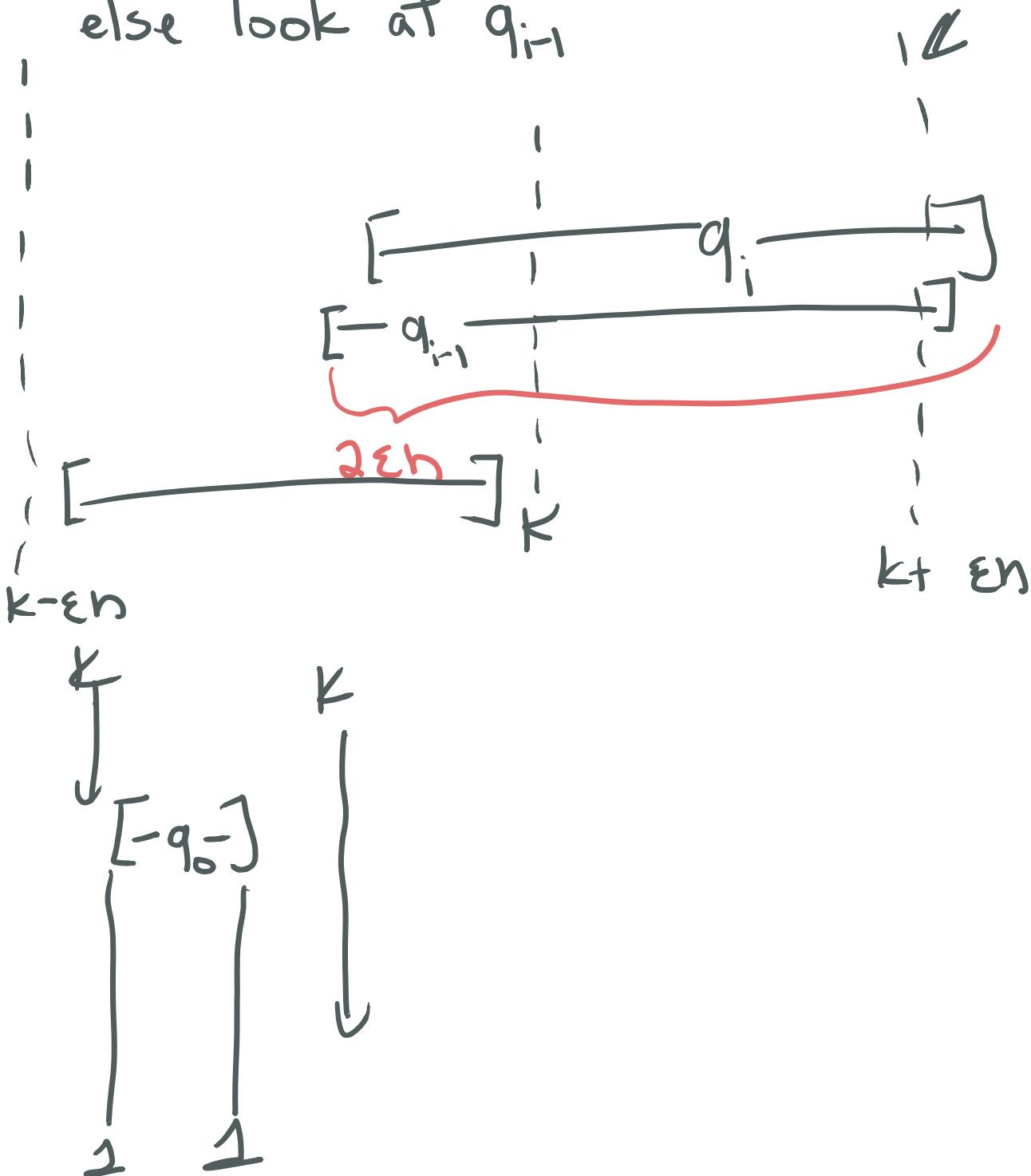
$k \in I(q_i)$ for some q_i

$k \notin I(q_i) \quad \forall q_i$

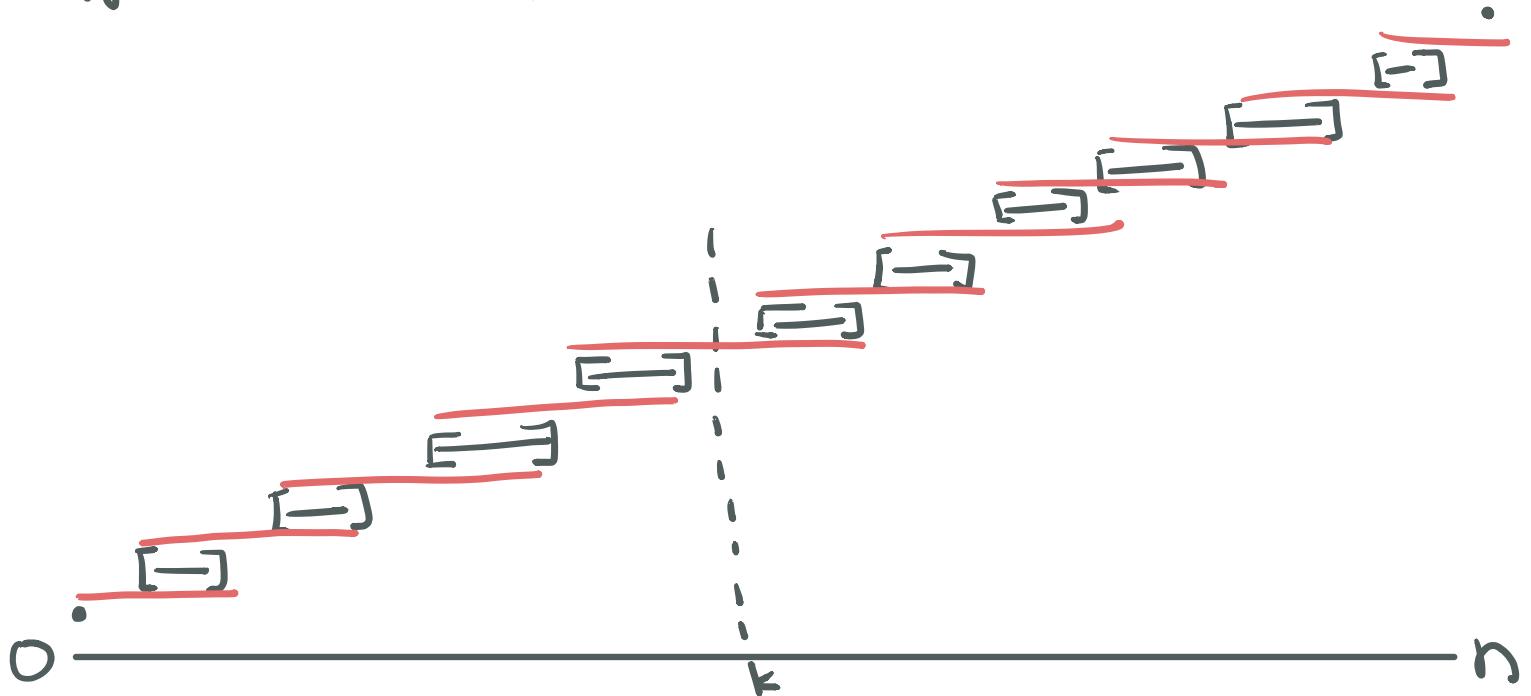
Suppose $k \in I(q_i)$ for some i

if $I(q_i) \subseteq [k - \varepsilon_n, k + \varepsilon_n]$ then done

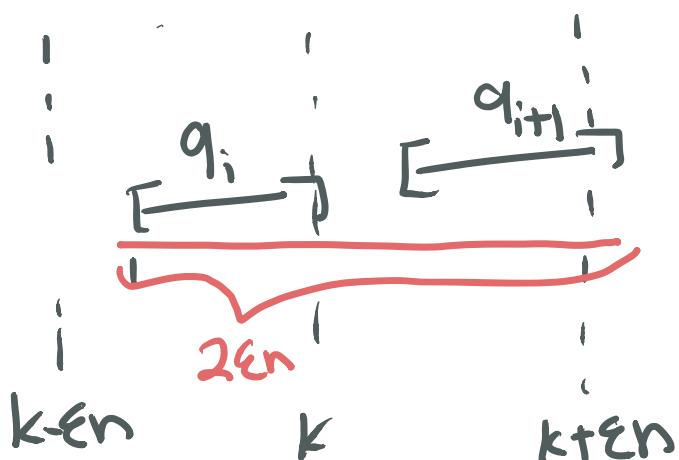
else look at q_{i-1}



suppose $k \notin I(q_i) + i$



the "combined intervals" cover $[n]$.
pick 1 covering k .



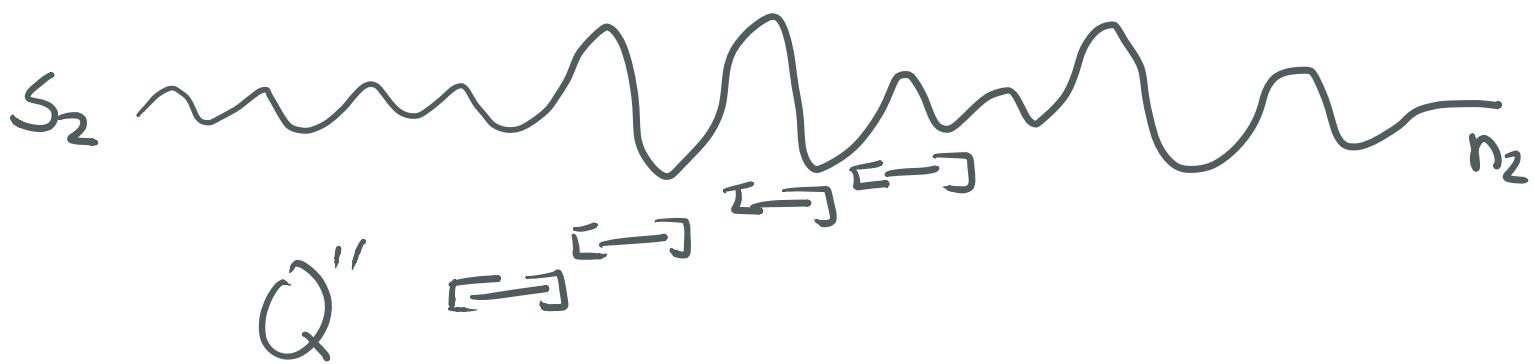
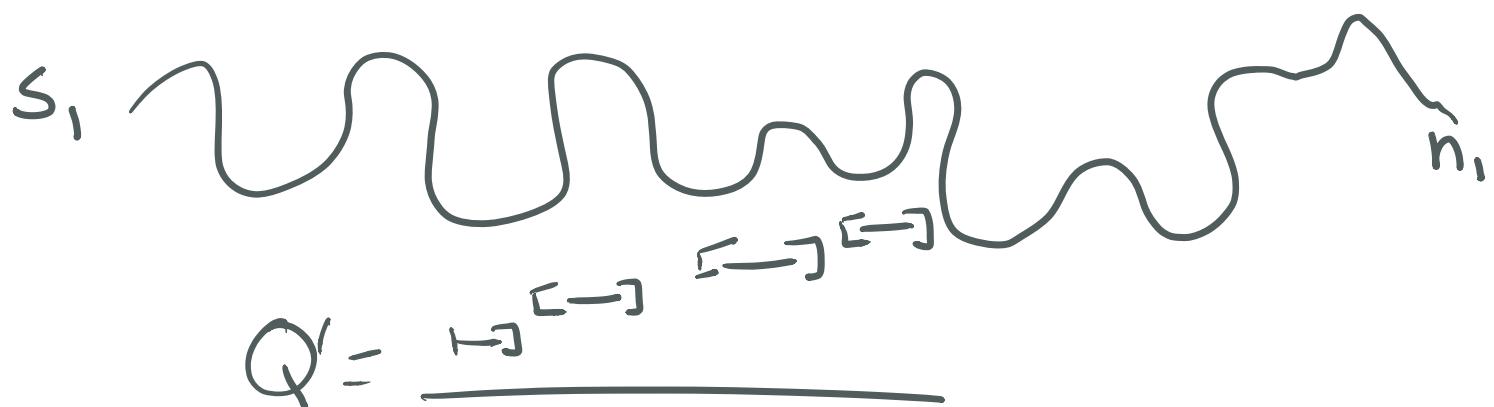
one of the intervals must lie inside

Key invariant: any two consecutive intervals have width $\leq 2\epsilon n$.

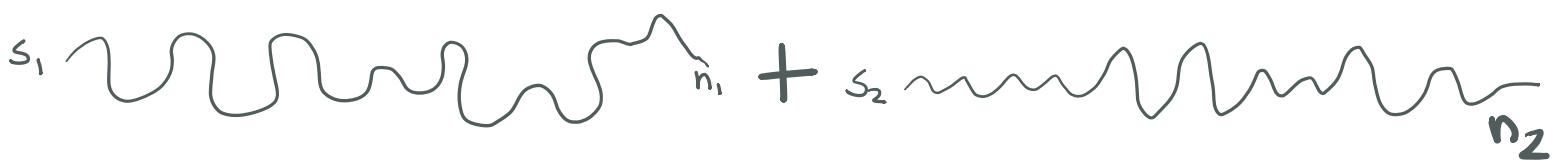
" ϵ -APX quantile summary"

Merging given two ϵ -APX quantile

summaries over 2 streams, want ϵ -APX
summary over combined stream



want to combine Q', Q'' to get summary of



$$"Q' + Q''" = \{q_1'', \dots, q_c'', I''(q_1''), \dots, I''(q_c'')\}$$

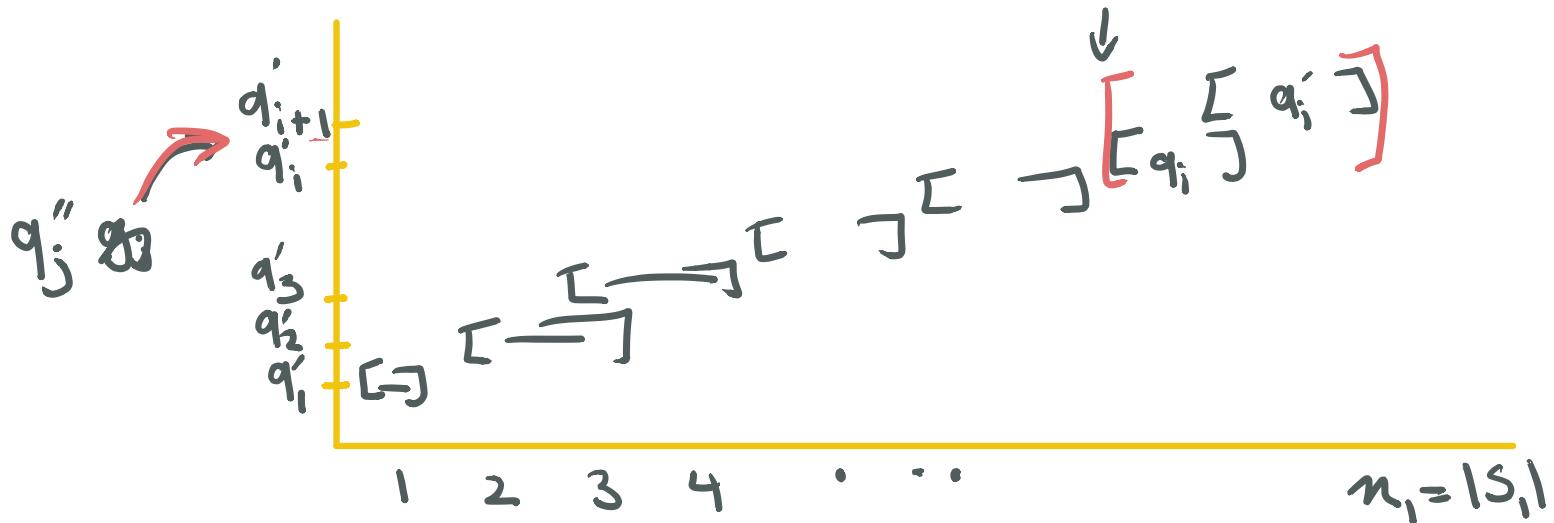
denote:

$$Q = \{q'_1, \dots, q'_e, I'(q'_1), \dots, I'(q'_e)\}$$

$$Q'' = \{q''_1, \dots, q''_m, I''(q''_1), \dots, I''(q''_m)\}$$

let $q''_j \in Q''$. $I''(q''_j)$ bounds rank q''_j wrt S_2

goal: bound rank q''_j wrt $S_1 + S_2$.



rank of q''_j in $S_1 \left\{ \begin{array}{l} \geq \min I'(q'_i) \\ \leq \max I'(q'_{i+1}) \end{array} \right.$

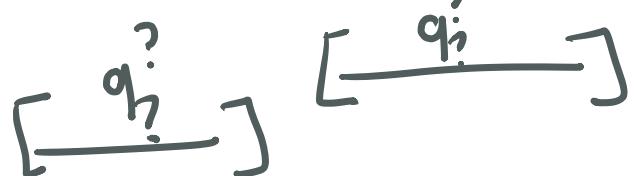
$$\begin{aligned} & \approx \min \underline{I'(q'_i)} + \min \overline{I''(q''_j)} \\ & \leq \text{rank}(q''_j \text{ in } S_1 + S_2) \\ & \leq \max \underline{I'(q'_{i+1})} + \max \overline{I''(q''_j)} \end{aligned}$$

$$\begin{aligned} \text{set } I''(q''_j) = & [\min I'(q'_i) + \min I''(q''_j), \\ & \max I''(q'_{i+1}) + \max I''(q''_j)] \end{aligned}$$

$$Q'' = \{q'_1, \dots, q'_e, q''_1, \dots, q''_m, \text{ w/ intervals } I''\}$$

To show Q'' is ϵ -APX, need to show
"2 ϵ_n width" property.

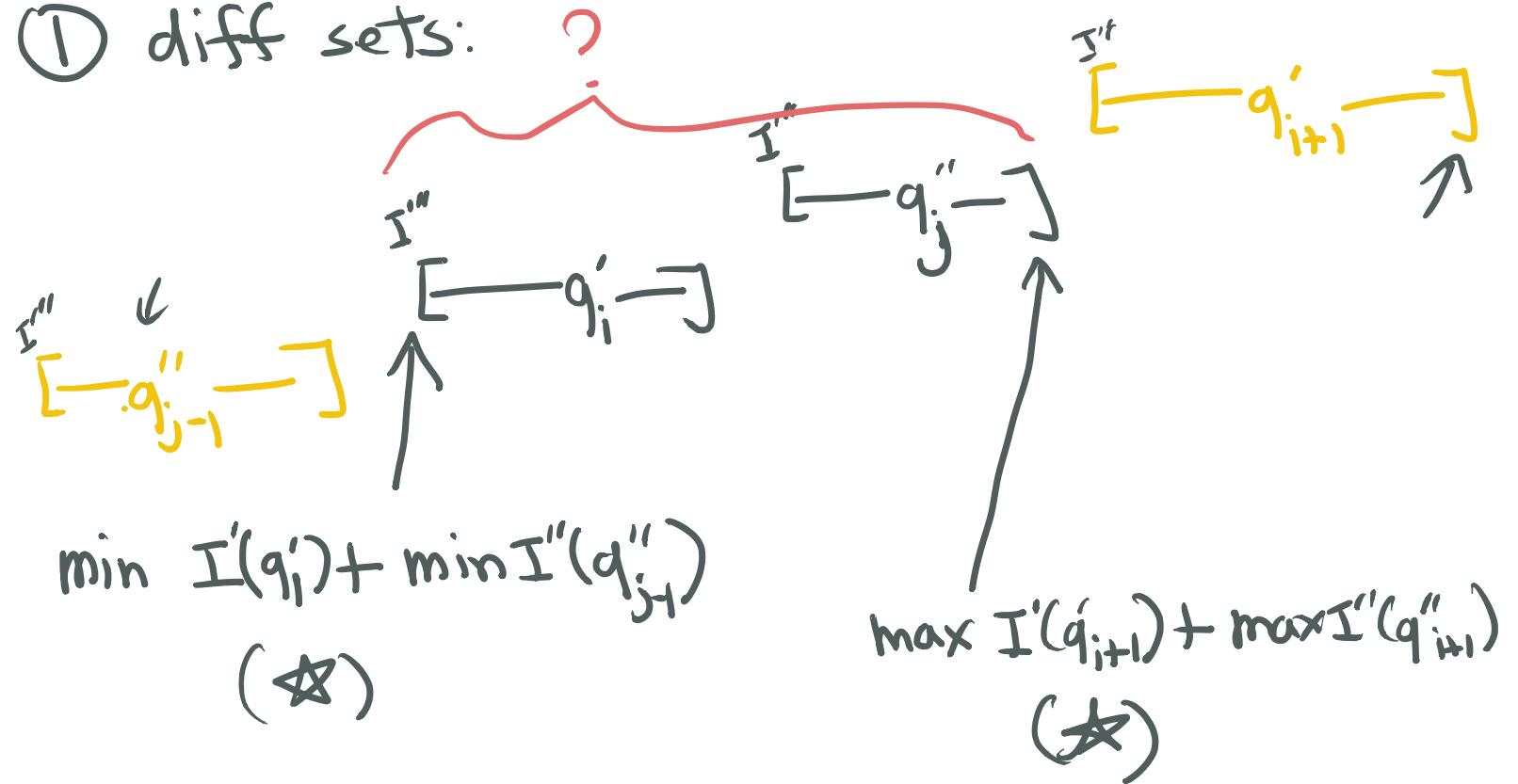
Take two consecutive intervals in Q_3 .



Two cases:

- ① elements from diff sets
- ② elements from same sets

① diff sets:



$$\begin{aligned}
 (\textcircled{*}) - (\textcircled{**}) &= \frac{\max I'(q'_{i+1}) + \max I''(q''_{j+1})}{\underline{\leq 2\epsilon n_1}} + \frac{\min I'(q'_i) + \min I''(q''_{j-1})}{\underline{+ 2\epsilon n_2} \cancel{+ 2\epsilon n_2}} \\
 &= 2\epsilon(n_1 + n_2)
 \end{aligned}$$

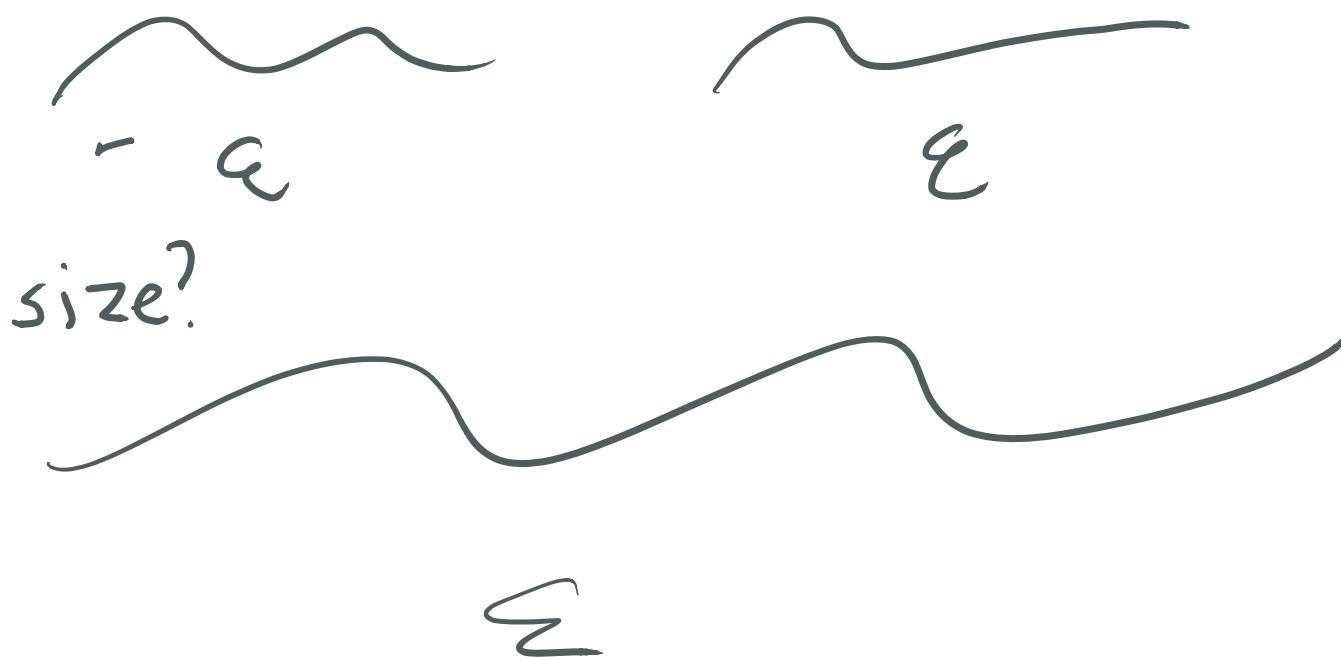
② same sets

$$\begin{array}{c} \text{---} \\ \downarrow \\ \boxed{-q_j''} \end{array} \quad \begin{array}{c} \overset{\text{---}}{I''} \\ \boxed{-q_i'} \end{array} \quad \begin{array}{c} \overset{\text{---}}{I'''} \\ \boxed{-q_{i+1}'} \end{array} \quad \begin{array}{c} \text{---} \\ \uparrow \\ \boxed{-q_{j+1}''} \end{array}$$
$$\min I'(q_i') + \min I''(q_j'')$$
$$\max I'(q_{i+1}') + \max I''(q_{j+1}'')$$

$$\frac{\max I'(q_{i+1}') + \max I''(q_{j+1}'') - \min I'(q_i') - \min I''(q_j'')}{2\epsilon n_1 + 2\epsilon n_2} = \frac{2\epsilon(n_1 + n_2)}{2\epsilon(n_1 + n_2)}$$

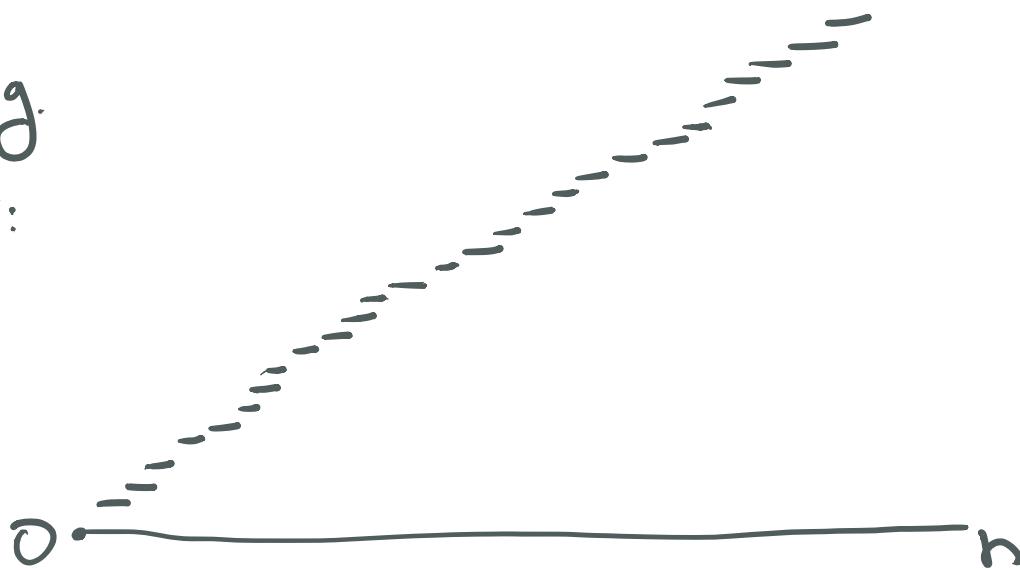


This shows that merging 2 ϵ -APX QS's
gives ϵ -APX QS of combined streams



Pruning.

Input:

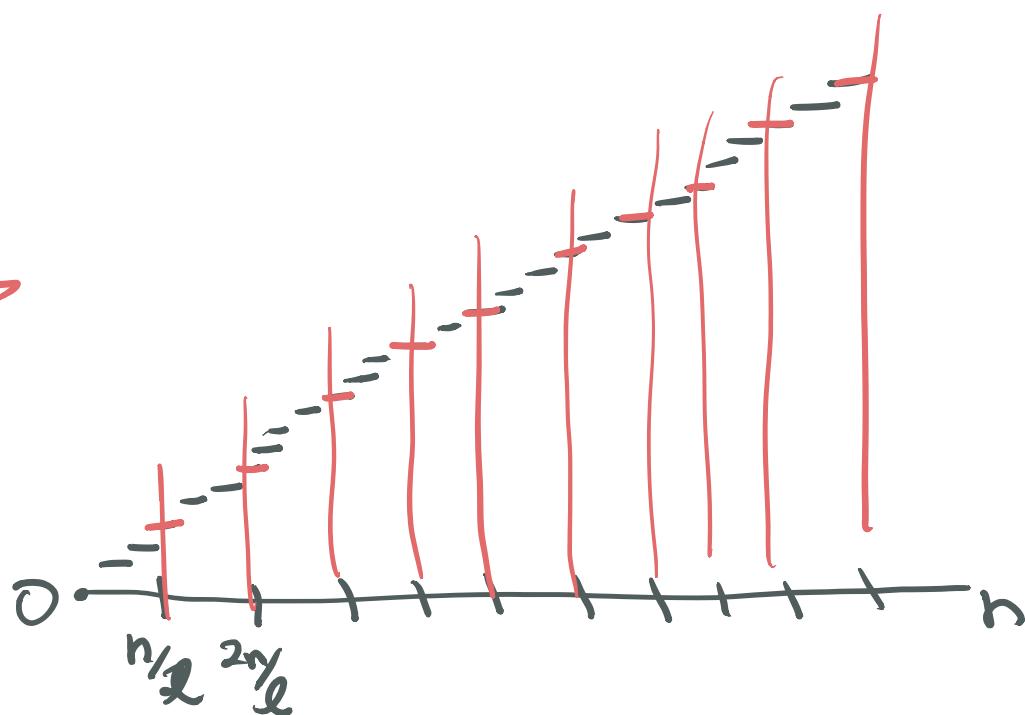


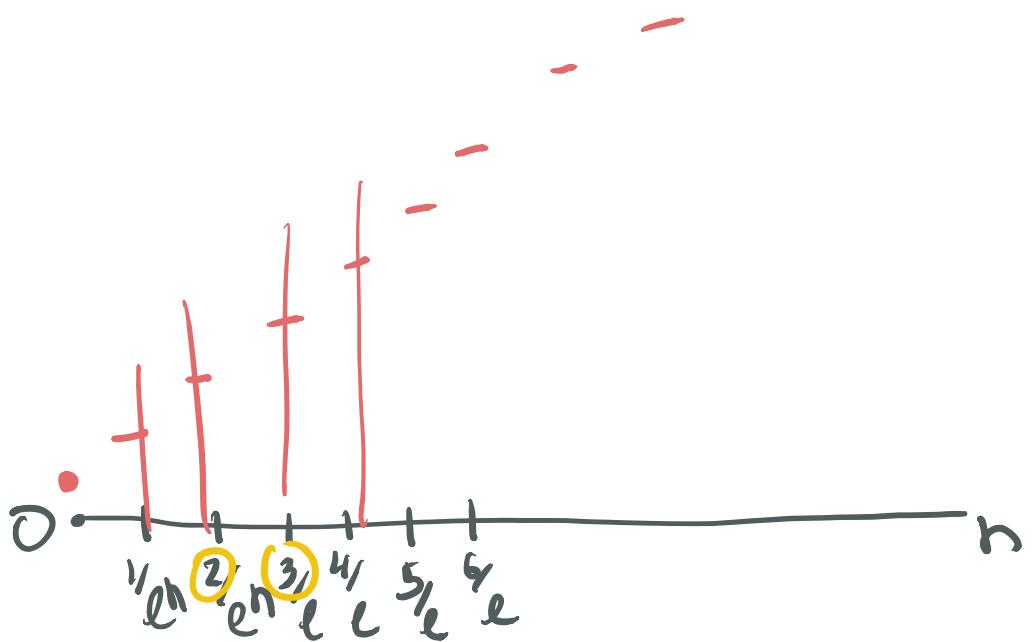
ϵ -approximate quantile summary w/ too many points

Goal: sparser summary that's still very good

& keep

l queries



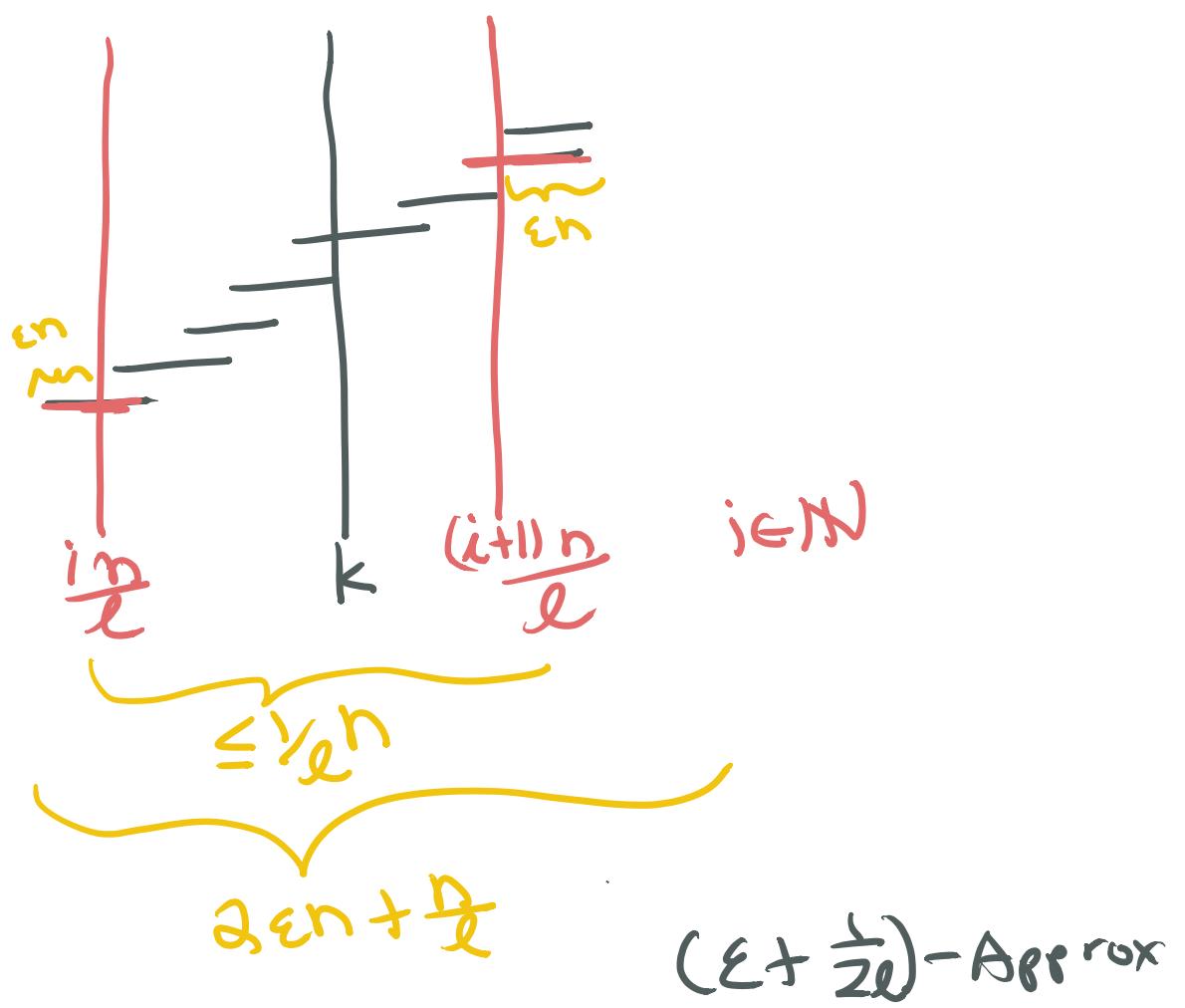


Claim: resulting quantile is $(\varepsilon + \frac{1}{2\varepsilon})$ -APX

Proof

suppose we query a rank k

look at
original
summary



Recap: we can

combine ϵ -APX quantile summaries

to get ϵ -APX quantile summary of whole thing

sparsify ϵ -APX quantile summary to

$(\epsilon + \frac{1}{\alpha k})$ -APX quantile summary w/ k points

Remains to address:

how to make one at all??

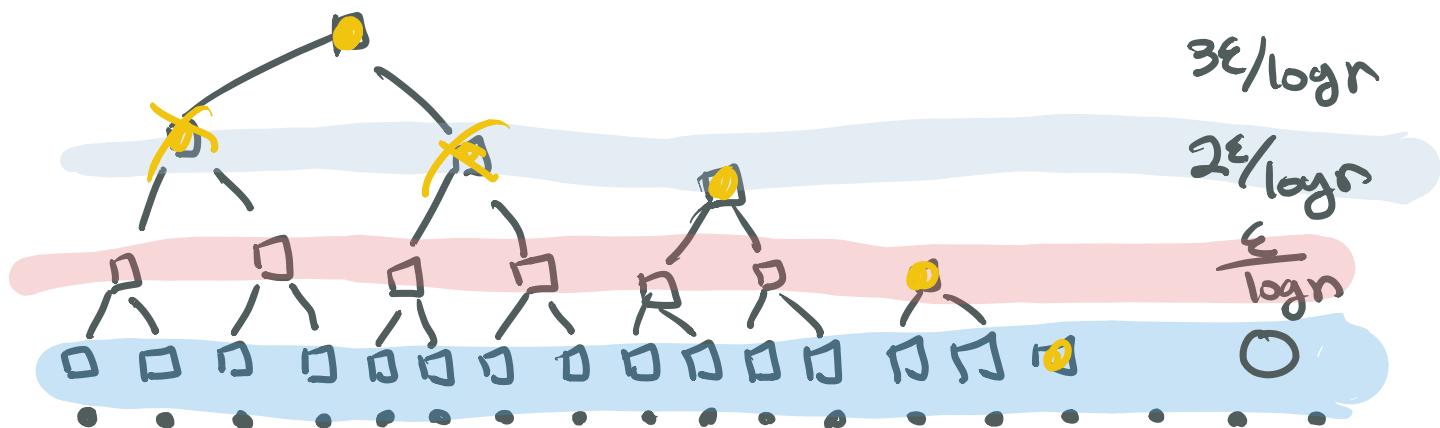
what if $n=1$?

Take the point

I claim that's all we need!

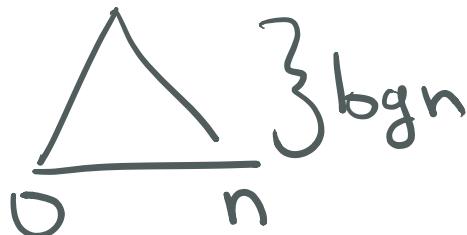


APX



take $k = \frac{\epsilon}{\log n}$ $\frac{\log n}{2\epsilon}$

at the root,



ϵ -approximate quantiles

Space?



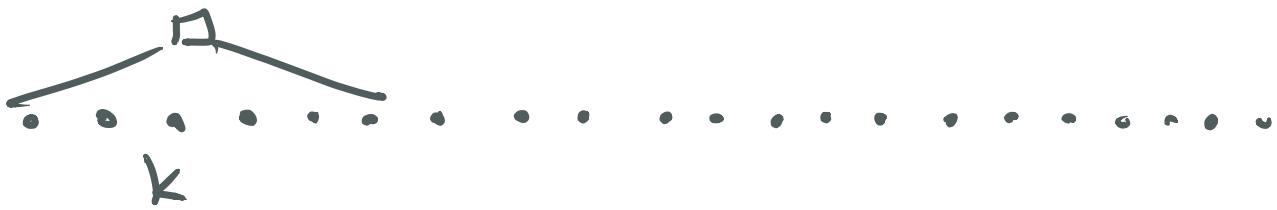
only keep "root summaries"

Theorem

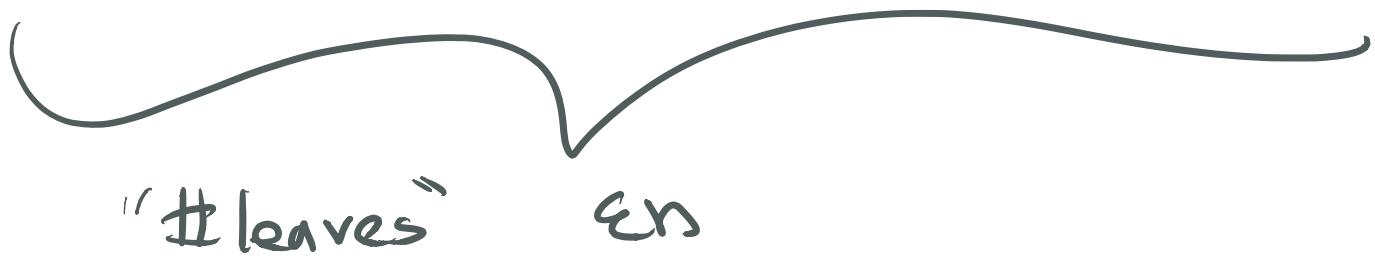
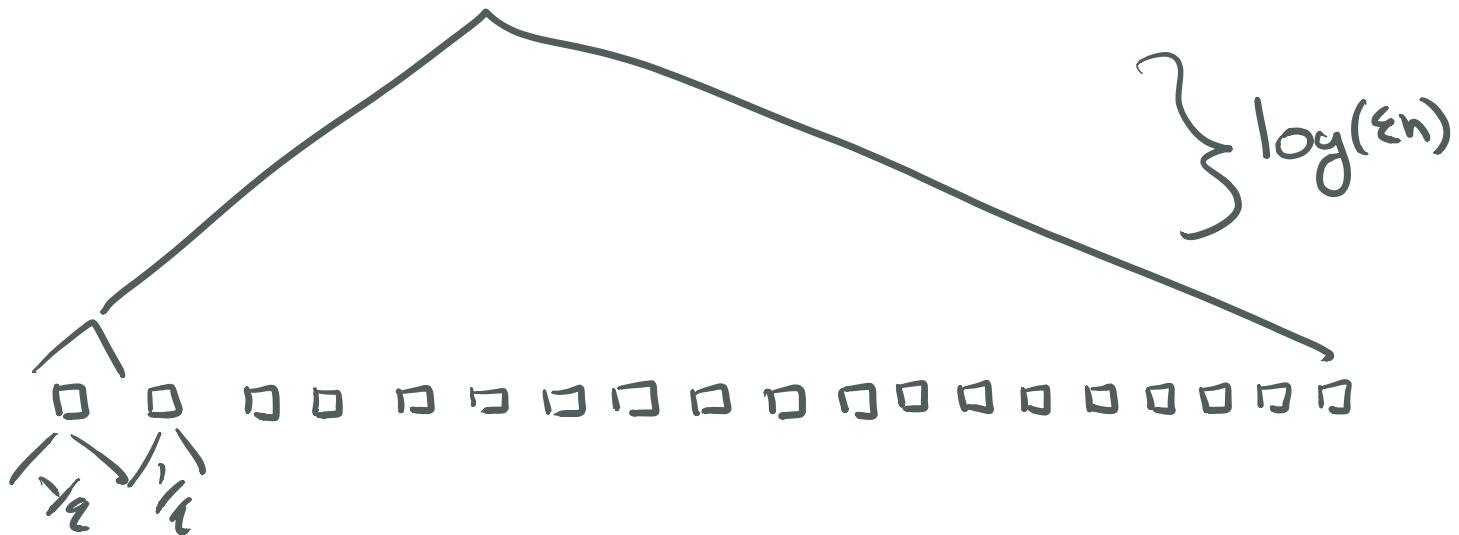
- 1-pass
- $O(\log^2(n)/\epsilon)$ space
- deterministic
- ϵ -APX quantile over stream

idea: mergability + dyadic intervals trick

slightly better:



have first level contain γ_ε points



Theorem ++

- 1-pass
- $O(\log^2(\epsilon n)/\epsilon)$ space
- deterministic
- ϵ -APX quantile over stream

Even better?

Khanna-Greenwald [2001]:

$\frac{1}{\epsilon} \log(\epsilon n)$ space

- more sophisticated quantile summary, merging
- interval trick

Finding the median (and other ranks) in p passes

Fix $p=2$ for simplicity.

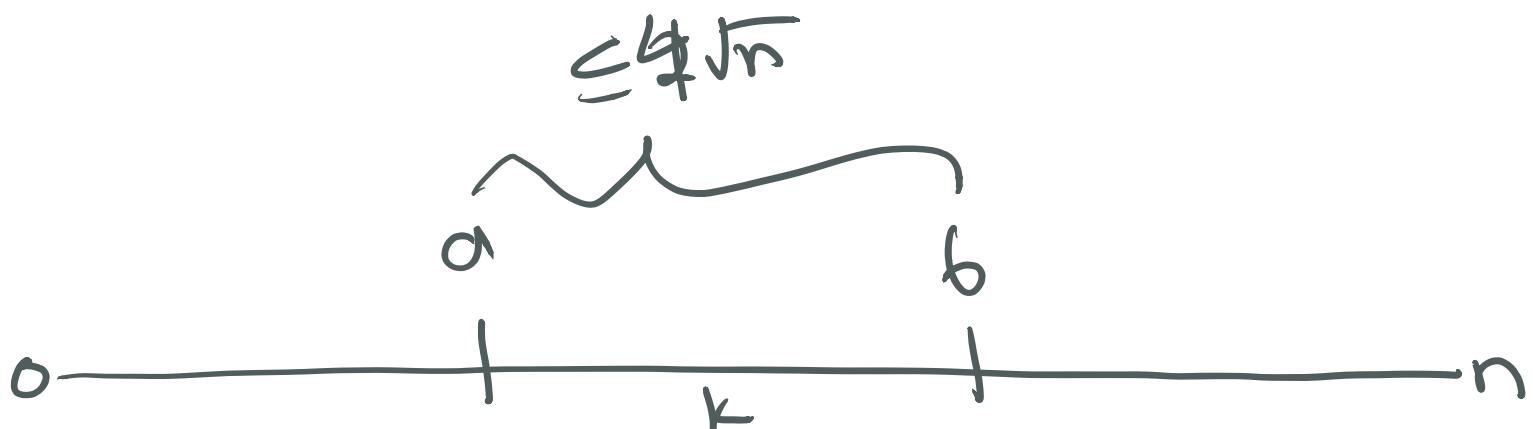
goal: $O(\sqrt{n} \text{ polylog}(n))$ space

suppose we are querying rank k .

1st pass: build ϵ -APX quantile summary

for $\epsilon = \frac{1}{\sqrt{n}}$ ($\sqrt{n} \log(n)$ space w/ GK)

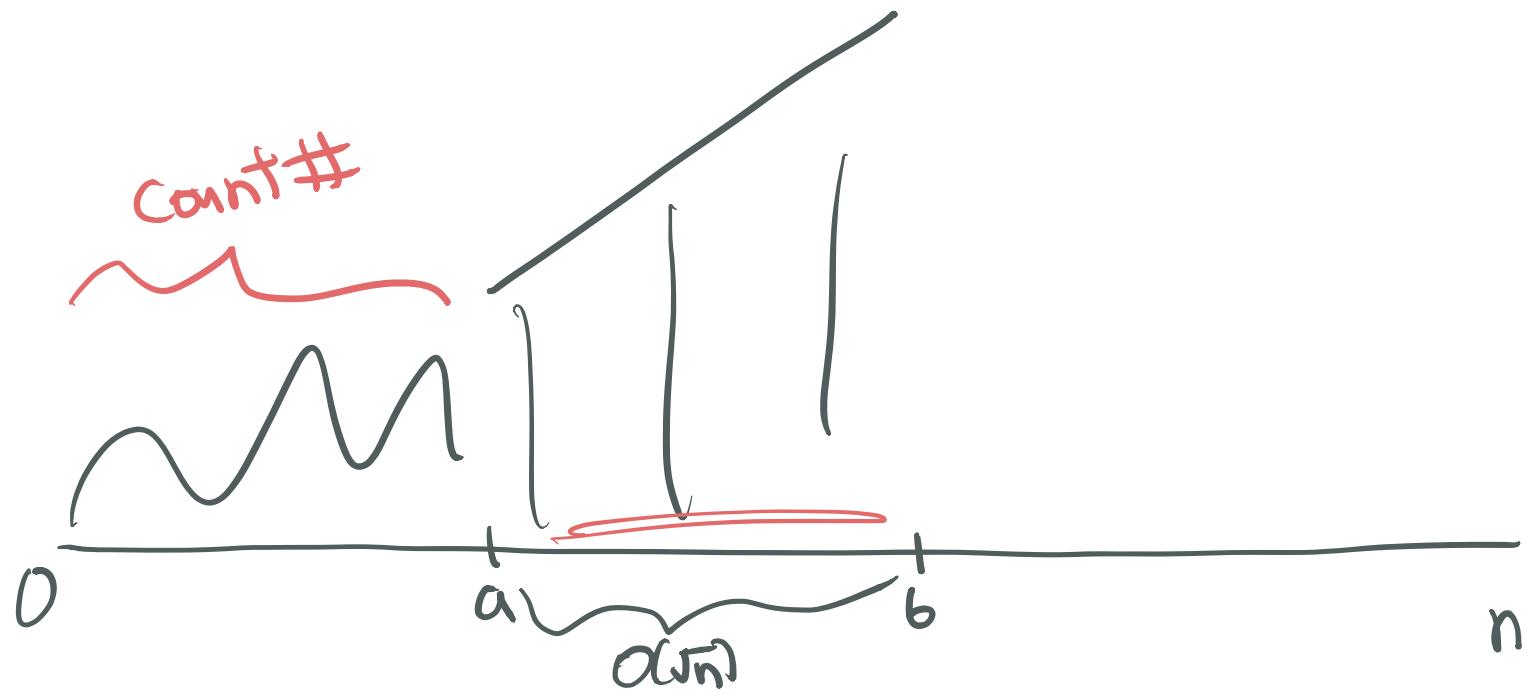
query $k - \sqrt{n}, k + \sqrt{n} \Rightarrow a, b$



$$k - 2\sqrt{n} \leq \text{rank}(a) \leq k \leq \text{rank}(b) \leq k + 2\sqrt{n}$$

$a = \text{query}(k - \sqrt{n})$

2nd pass:



k over all

Take the $(k - \#\{ < a \})$ th in the sorted set

for general p:

make $\frac{1}{n^{1/p}}$ -APX quantile summaries
and filter.

After p passes, down to $n^{1/p}$ elements.
sort and select.