

APPROXIMATION ALGORITHMS

SUMMARIES FOR ORDERED DATA

RASMUS PAGH

UNIVERSITY OF COPENHAGEN



TODAY

- MOTIVATING EXAMPLE
 - SAMPLING SOLUTION
 - Q-DIGEST
 - DYADIC COUNT-MIN
- } DATA FROM FIXED DOMAIN
- KLL SUMMARY (ANY ORDERED DOMAIN)

MOTIVATING EXAMPLE

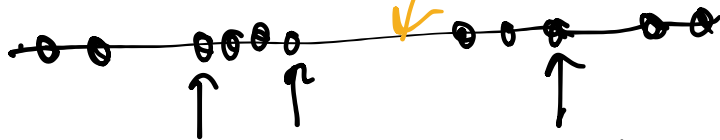
CONSIDER A SET S OF TIME STAMPS (E.G. LINK CLICKS)
WANT TO APPROXIMATE NUMBER OF TIME STAMP IN
A GIVEN INTERVAL.

CLASSIC APPROACH:

$$\text{rank}(x) = |\{y \in S \mid y \leq x\}|$$

QUANTILES

N ELEMENTS OF S IN ORDER



CAN APPROXIMATE
FROM QUANTILES
UP TO ERROR ϵN

STORE EVERY ϵN 'TH ELEMENT

("QUANTILES", E.G. 10%, 20%, ..., 90%)

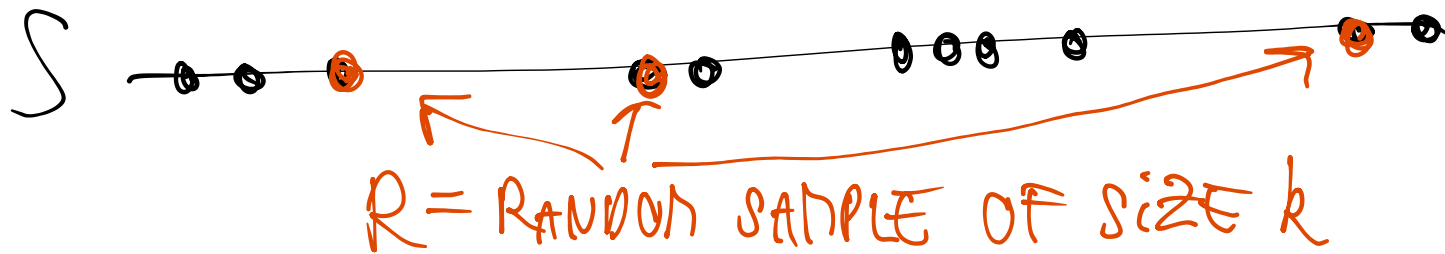
REQUIRES SPACE N TO COMPUTE, $\frac{1}{\epsilon}$ TO STORE

TWO SETTINGS

NEXT \rightarrow • $S \subseteq \{0, \dots, U-1\}$, $U \in \mathbb{N}$

LATER \rightarrow • S IS ANY ORDERED SET

WHY NOT JUST USE SAMPLING?



$$E[|\{y \in R \mid y \leq x\}|] = \frac{k}{N} |\{y \in S \mid y \leq x\}| = \frac{k}{N} \text{rank}(x)$$

$$\text{rank}(x) \approx \frac{N}{k} |\{y \in R \mid y \leq x\}|$$

EXERCISE:

DERIVE THE SIZE OF
R NEEDED TO MAKE THE
ERROR $\leq \epsilon N$ FOR ALL X

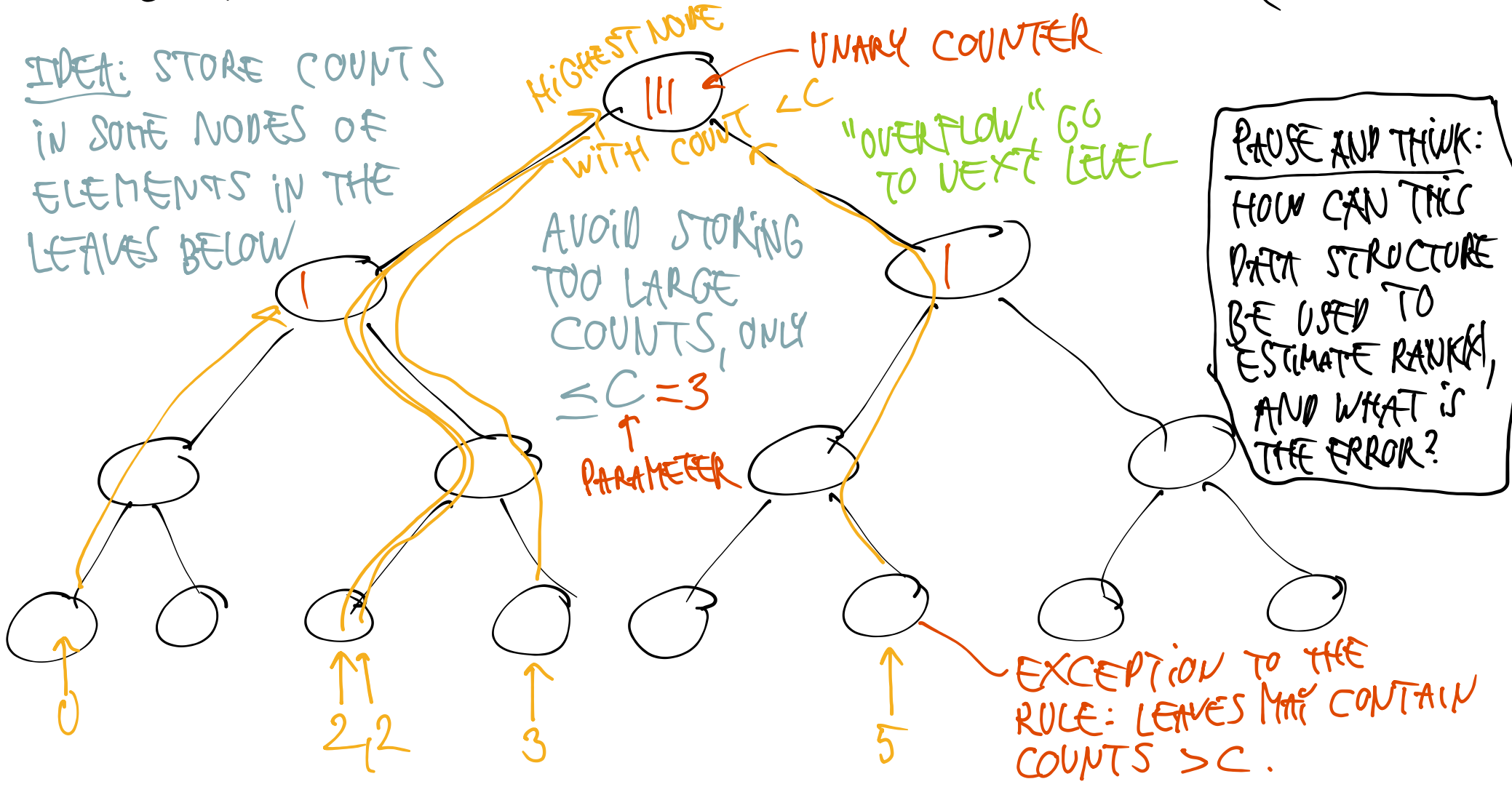
Q-DIGEST

ASSUME V IS A POWER OF 2 (WITHOUT LOSS OF GEN)

$$S \subseteq \{0, \dots, V-1\}$$

COMPLETE BINARY TREE WITH V LEAVES ($V=8$)

IDEA: STORE COUNTS IN SOME NODES OF ELEMENTS IN THE LEAVES BELOW



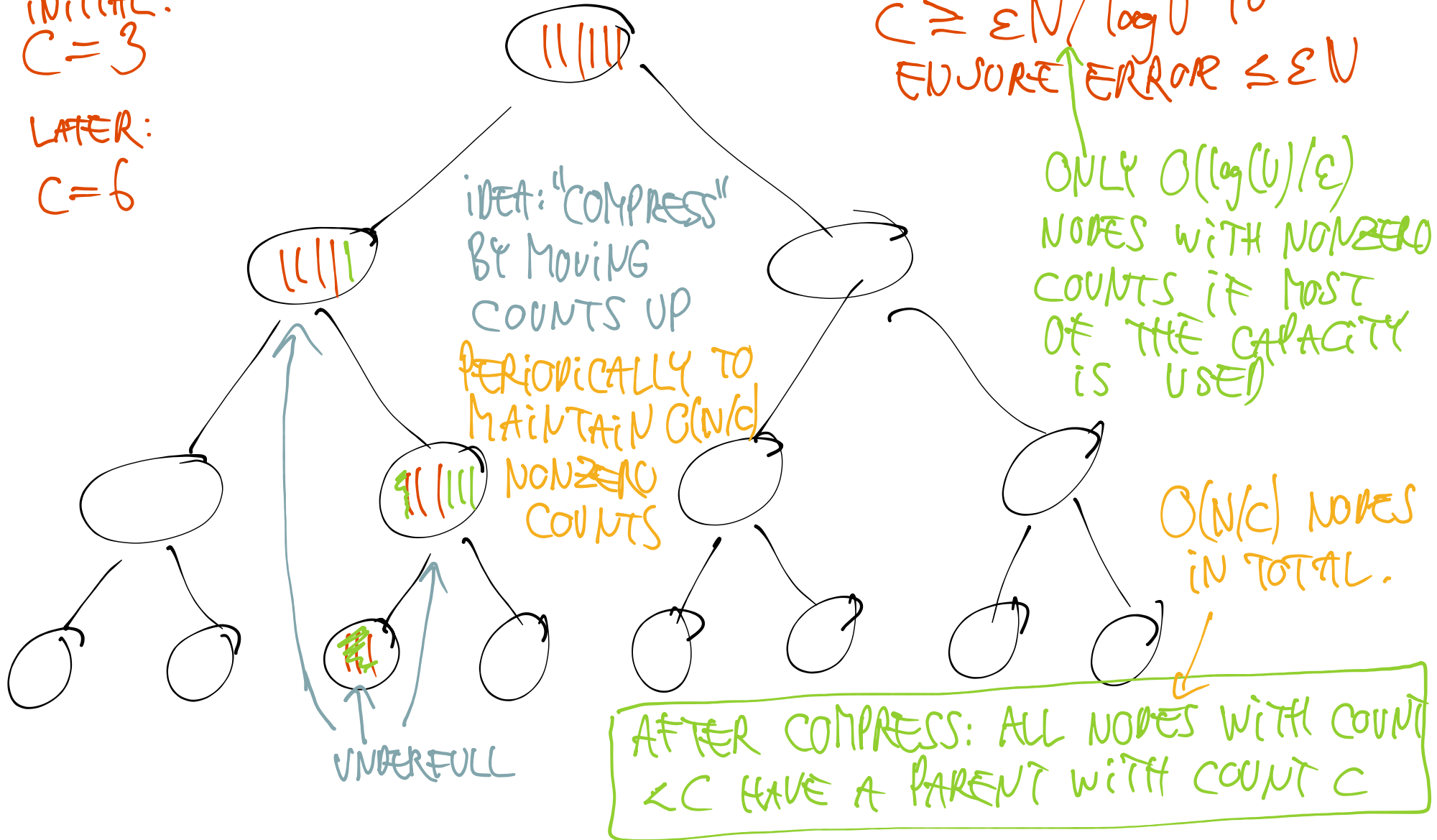
PAUSE AND THINK:
HOW CAN THIS DATA STRUCTURE BE USED TO ESTIMATE RANKS, AND WHAT IS THE ERROR?

EXCEPTION TO THE RULE: LEAVES MAY CONTAIN COUNTS $> C$.

REDUCING THE SIZE OF Q-DIGEST

INITIAL:
C=3

LATER:
C=6



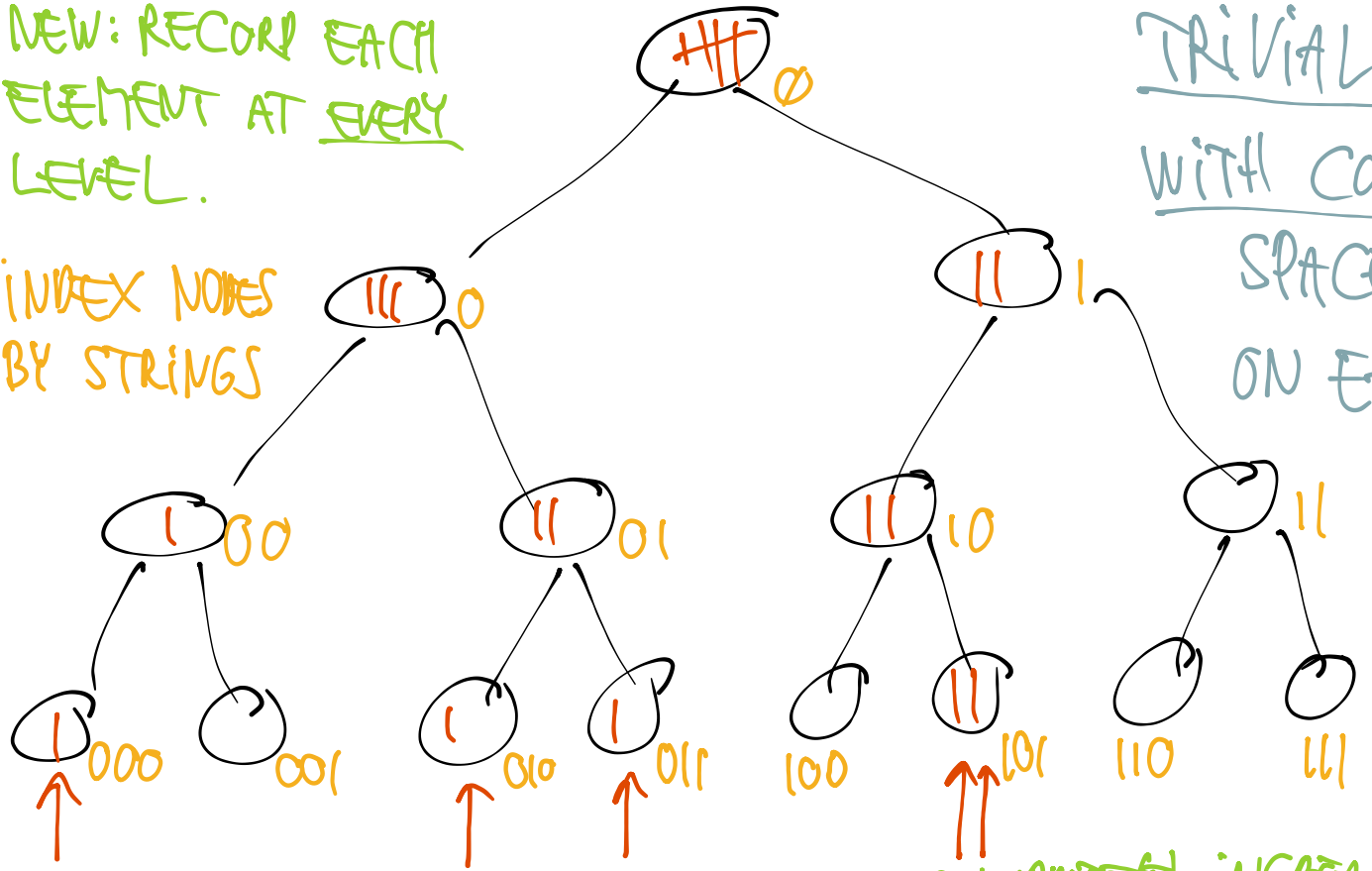
DYADIC COUNT-MIN SKETCH

← IN BOOK COUNT SKETCH IS USED, BUT IT WORKS THE SAME WAY

DIFFERENT APPROACH TO APPROXIMATING COUNTS
IN NODES OF COMPLETE BINARY TREE: USE A SKETCH!

NEW: RECORD EACH ELEMENT AT EVERY LEVEL.

INDEX NODES BY STRINGS



TRIVIAL SPACE: #DISTINCT STRINGS

WITH COUNT-MIN: CAN CHOOSE SPACE TO GET ERROR $\leq \frac{\epsilon N}{\log U}$ ON EACH ESTIMATE WITH PROBABILITY $1 - \epsilon/\log U$.

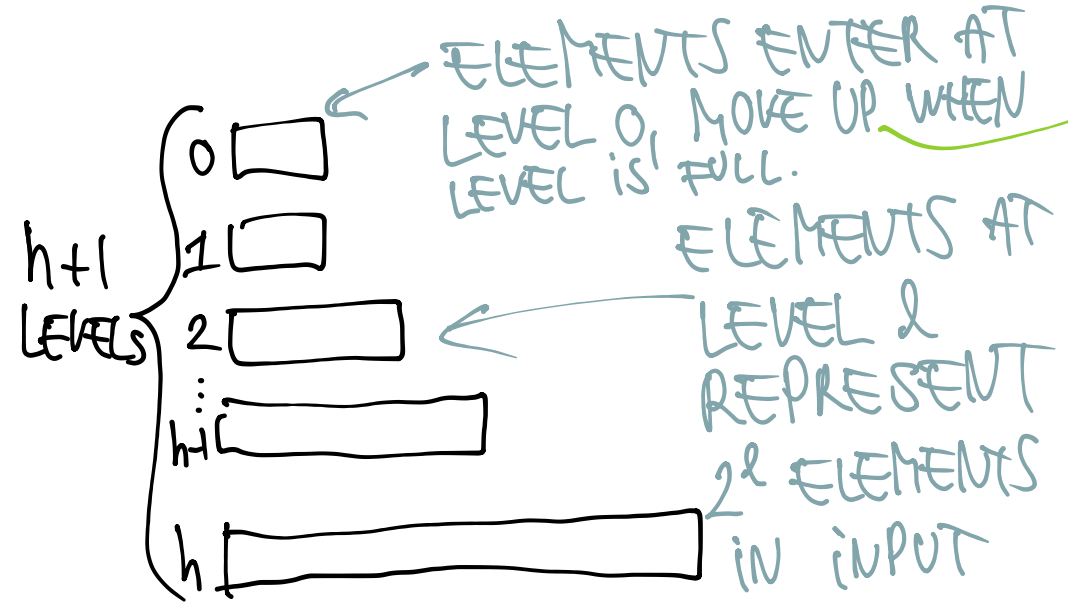
SPACE: $|w^*| \leq N \log U$

$O\left(\frac{1}{\epsilon} \log^2 U \log\left(\frac{\log U}{\delta}\right)\right)$
(CAN BE IMPROVED)

ON UPDATE(S), INCREASE COUNTS OF '101, 10, 1, 0' (ALL PREFIXES OF BIN. REP.)

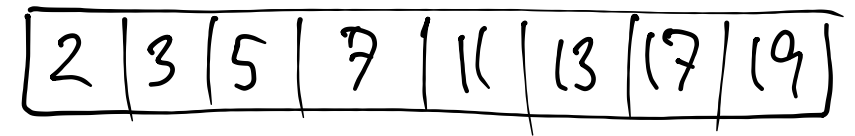
KARNIN-LANG-LIBERTY (KLL)

- WORKS WITH ANY ORDERED DOMAIN (NO "log" DEPENDENCE)
- DATA STRUCTURE IS A LIST OF ORDERED LISTS "BUFFERS"



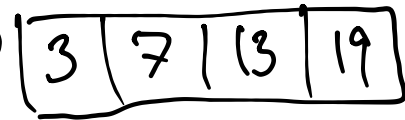
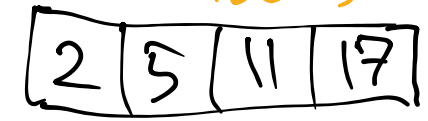
"COMPRESS", DISCARDING HALF OF THE ELEMENTS

COMPRESS BY EXAMPLE:



CHOOSE BETWEEN

ODD-NUMBERED

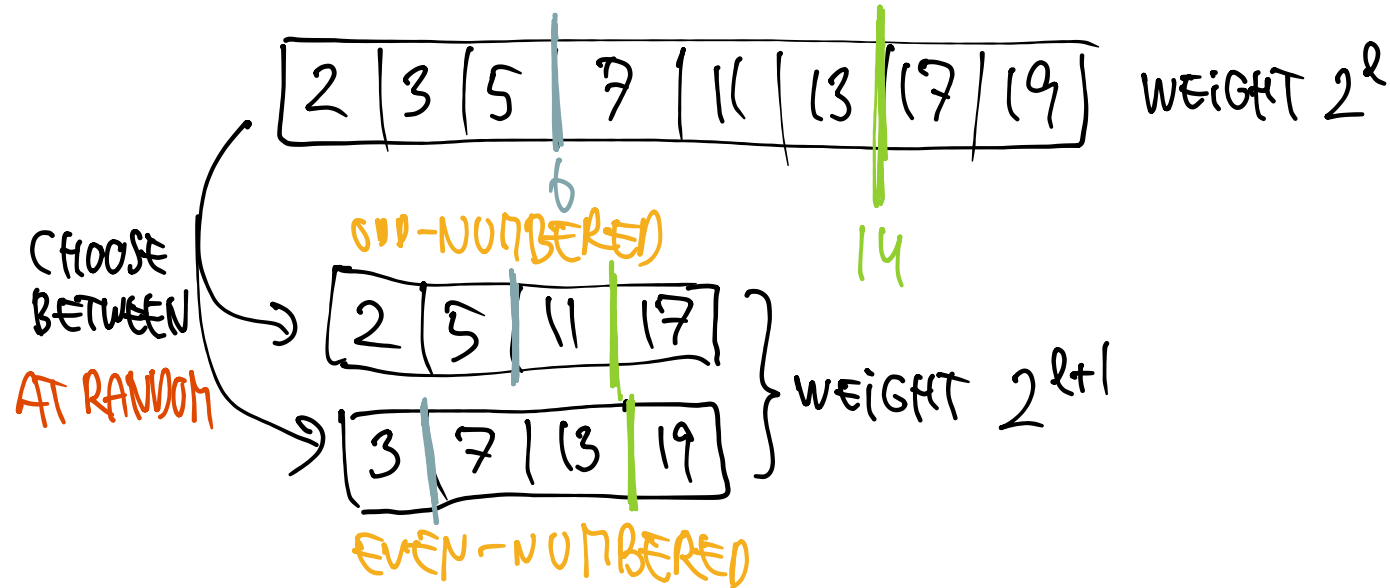


EVEN-NUMBERED

FOR SOME $c \in (\frac{1}{2}, 1)$
 EACH LEVEL GROWS BY
 \approx A FACTOR $\frac{1}{c}$
 (EXCEPT LOWEST LEVELS WHICH HAVE SIZE 2)

HOW COMPRESS AFFECTS RANK

COMPRESS BY EXAMPLE:



CASE 1: RANK IS "EVEN", $2i \cdot 2^2$ IN ORIGINAL LIST

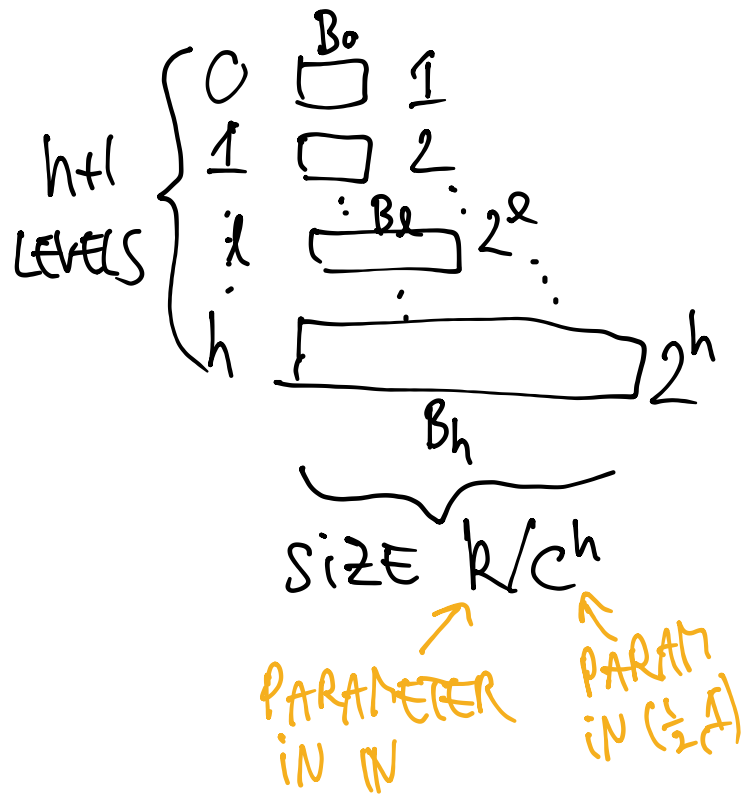
\rightarrow RANK BECOMES $i \cdot 2^{2+1}$ IN BOTH ODD- AND EVEN-NUMBERED CASES

CASE 2: RANK IS "ODD", $(2i+1) \cdot 2^2$ IN ORIGINAL LIST

\rightarrow RANK IS $i \cdot 2^{2+1}$ IN EVEN-NUMBERED CASE,
AND $(i+1) \cdot 2^{2+1}$ IN ODD-NUMBERED CASE

CORRECT ON AVERAGE EXPECTATION

KLL, THE GENERAL PICTURE



UNBIASED ESTIMATOR FOR RANK(x):

$$\tilde{r}(x) = \sum_{i=0}^h |\{y \in B_i \mid y \leq x\}| \cdot 2^i$$

COMPACTION OPERATIONS AT LEVEL l

$$\tilde{r}(x) - \text{RANK}(x) = \sum_{l=0}^{h-1} \sum_{i=1}^{m_l} 2^l X_{i,l}$$

$X_{i,l} = \begin{cases} +1 & \text{if POS. ERR} \\ -1 & \text{if NEG. ERR} \end{cases}$
 IN COMPACTION

ADDITIVE CHERNOFF-HOFFMANN BOUND:

GIVEN INDEPENDENT X_1, \dots, X_n WITH $|X_i| \leq a_i$ FOR $X = \sum_i X_i$ AND $E[X] = 0$ IT HOLDS THAT

$$\Pr[|X| > k] \leq 2 \exp\left(\frac{-k^2}{2 \sum_i a_i^2}\right).$$

(SPECIAL CASE OF FACT 1.9 FROM BOOK)

$$\tilde{r}(X) - \text{RANK}(X) = \sum_{l=0}^{h-1} \sum_{i=1}^{m_l} 2^l X_{i,l}$$

"a_i"

COMPACTION OPERATIONS

$$\sum a_i^2 = \sum_{l=0}^{h-1} \sum_{i=1}^{m_l} 2^{2l} = \sum_{l=0}^{h-1} m_l \cdot 2^{2l} = O\left(\sum_{l=0}^{h-1} \left(\frac{2}{c}\right)^{h-l} \cdot 2^{2l}\right) \approx O(2^{2h})$$

CLAIM 1

DOMINATED BY LAST TERM

ENOUGH TO CHOOSE

$k = 2\sqrt{\log(1/\delta)}/\epsilon$
TO GET ERROR PROB. δ .

$$\Pr[|X| > \epsilon N] \leq 2 \exp\left(-\frac{(\epsilon N)^2}{2^{2h+1}}\right)$$

CLAIM 2 \rightarrow

$$\leq 2 \exp\left(-\frac{(\epsilon N)^2}{4(N/R)^2}\right)$$

$$= 2 \exp(-k^2/(4\epsilon^2))$$

ADDITIVE CHERNOFF-HOFFMANN BOUND:

GIVEN INDEPENDENT X_1, \dots, X_n WITH $|X_i| \leq a_i$,
FOR $X = \sum X_i$ AND $\mathbb{E}[X] = 0$ IT HOLDS THAT

$$\Pr[|X| > k] \leq 2 \exp\left(-\frac{k^2}{2 \sum a_i^2}\right).$$

PAUSE AND THINK:
WHY ARE ALL RANKS ACCURATE WITHIN $\pm \epsilon N$ WITH PROB. δ/ϵ ? (HINT: UNION BOUND)

CLAIM 1:

$$m_l = O\left(\left(\frac{2}{c}\right)^{h-l}\right)$$

$h-l=0$. NO CONTRACTIONS
AT LEVEL h .

INDUCTION STEP. (INFORMAL)

SUPPOSE $m_l = O\left(\left(\frac{2}{c}\right)^{h-l}\right)$, THEN

THE OUTPUT OF THESE CONTRACTIONS
HAS SIZE $m_l \cdot \underbrace{c^{h-l}}_{\text{BUFFER SIZE}} \cdot \underbrace{k/2}_{\text{BOUND ON \# CONTRACTIONS AT LEVEL } l+1} = m_l \cdot \frac{c}{2} \cdot \underbrace{c^{h-(l+1)} k}_{\text{BUFFER } l+1 \text{ SIZE}}$

$$m_{l+1} \leq m_l \cdot \frac{c}{2} \stackrel{\text{IND. HYP}}{=} O\left(\left(\frac{2}{c}\right)^{h-(l+1)}\right)$$

BOUND ON
CONTRACTIONS
AT LEVEL
 $l+1$

BOUND ON
CONTRACTIONS
AT LEVEL
 $l+1$

CLAIM 2

$$2^{2h} \leq 4(N/k)^2$$

- LEVEL $h-1$ ELEMENTS
CORRESPOND TO 2^{h-1} INPUT ELEM.
- AT LEAST kc INPUT ELEMENTS
HAVE BEEN AT LEVEL $h-1$

$$kc \cdot 2^{h-1} \leq N$$

$$\Downarrow 2^h \leq \frac{2}{c} N/k$$

$$\Downarrow 2^{2h} \leq 16(N/k)^2$$

CAN BE IMPROVED
TO 4, SEE BOOK

KLL WRAP-UP

• SPACE USAGE:

$$\sum_{l=0}^h c^{h-l} k = c^h k \sum_{l=0}^h c^{-l} = O(k)$$

GEOMETRICALLY INCREASING,
DOMINATED BY LAST TERM

THEORETICAL SPACE
USAGE CAN BE IMPROVED
TO $O\left(\frac{1}{\epsilon} \log \log\left(\frac{1}{\epsilon}\right)\right)$ BUT
MUCH MORE COMPLICATED
AND NOT MERGEABLE

• MERGING:

GIVEN TWO KLL SUMMARIES, CAN MERGE
BY MERGING LEVEL $l=0, 1, 2, \dots, h$ (IN THIS ORDER),
DOING COMPACTIONS AS NEEDED.

- TIME USAGE: IN EXPECTATION, EACH INPUT IS
THROWN OUT AFTER BEING IN 2 COMPACTIONS
→ EXPECTED CONSTANT TIME.