

APPROXIMATION ALGORITHMS

CARDINALITY ESTIMATION

RASMUS PAGH

UNIVERSITY OF COPENHAGEN

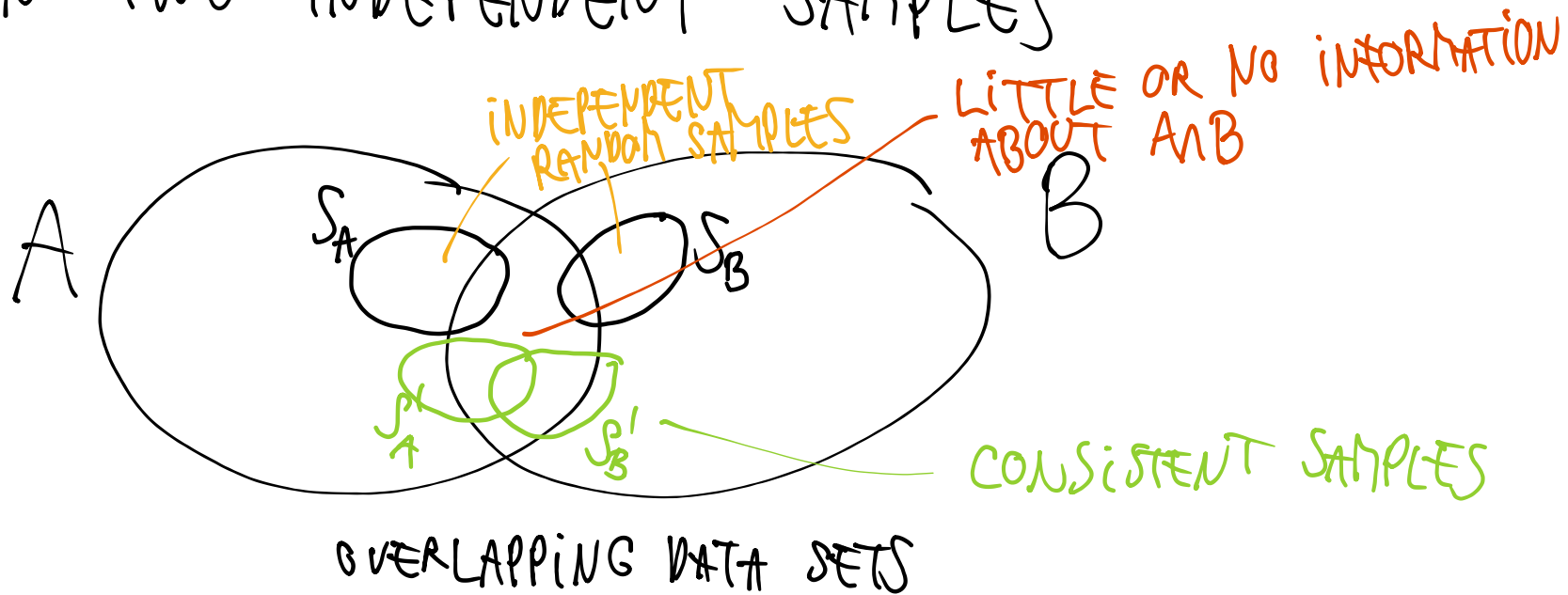


TODAY

- COORDINATED SAMPLING (KRV SUMMARY)
 - APPLICATION: DISTRIBUTED COUNTING
- HYPERLOGLOG CARDINALITY ESTIMATOR
 - APPLICATION: APPROXIMATE NEIGHBORHOOD FUNCTION
- BLOOM FILTERS
 - APPLICATION: APPROXIMATING INTERSECTION SIZE

COORDINATED SAMPLING

ISSUE: UNLIKELY THAT SAME ELEMENT APPEARS
IN TWO INDEPENDENT SAMPLES



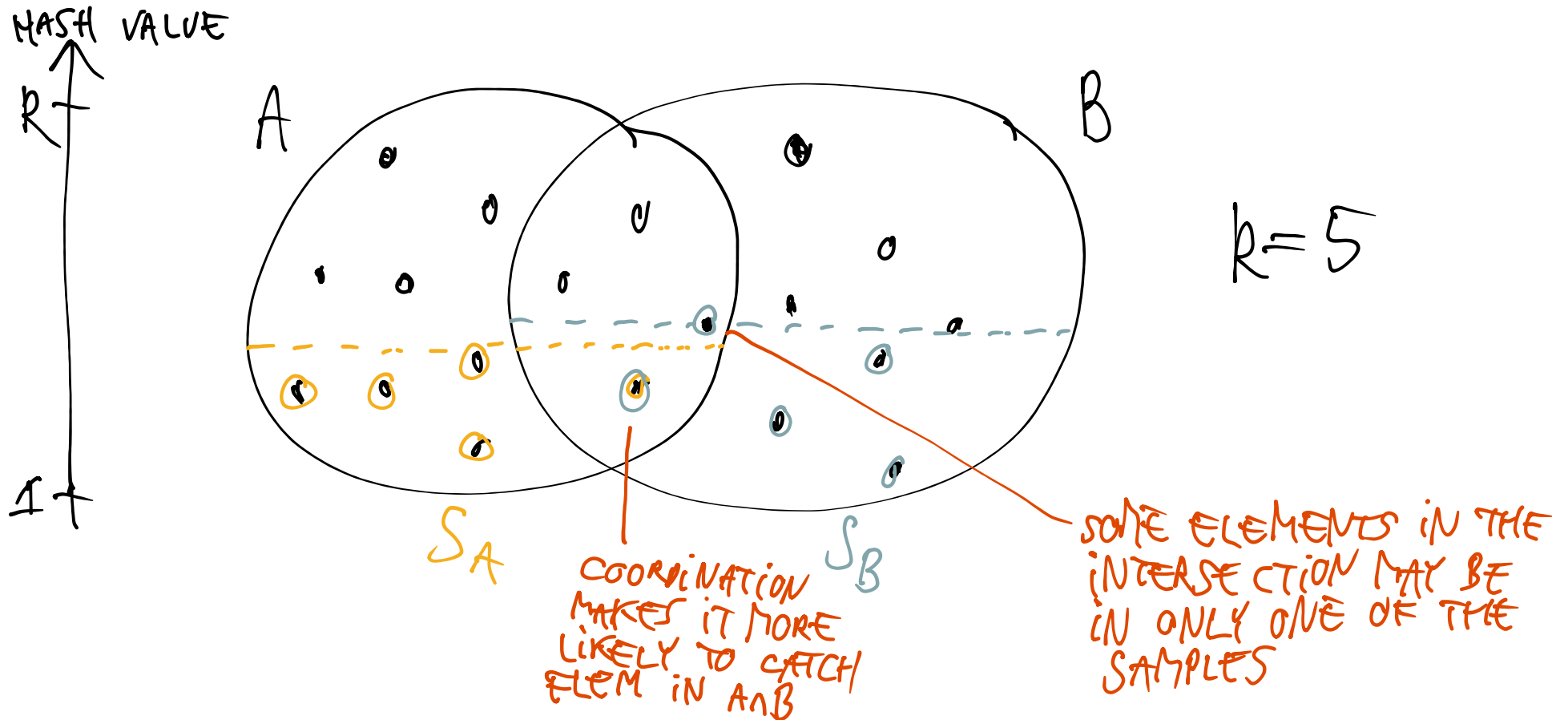
IDEA: MAKE THE DECISION OF WHETHER x SHOULD BE
SAMPLED "CONSISTENT" BY BASING IT ON A RANDOM
HASH VALUE $h(x)$

SAME HASH FUNCTION USED
FOR CHOOSING S_A AND S_B

K-MINIMUM VALUE (KMV) SUMMARY

LARGE VALUE
LIKE 2^{64}

- PICK RANDOM HASH FUNCTION h , MAPPING TO $\{1, \dots, R\}$
- SAMPLE S_A CONSISTS OF THE PAIRS $(x, h(x))$, $x \in A$ THAT HAVE THE k SMALLEST HASH VALUES.



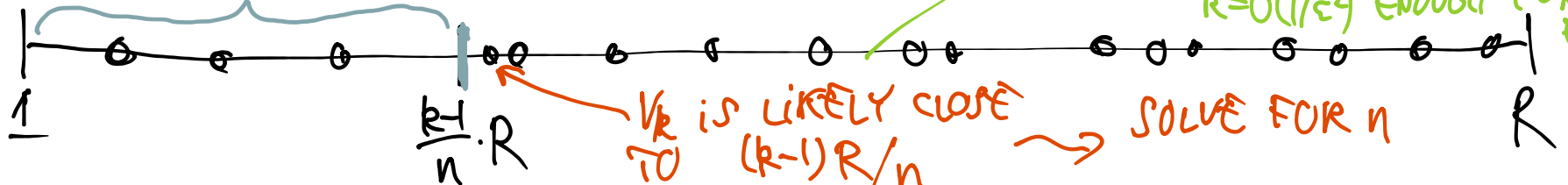
KMV PROPERTIES

- EFFICIENT UPDATES ($O(\log k)$ OR EVEN $O(1)$)
 - INSERTING ELEMENT x MORE THAN ONCE DOES NOT CHANGE S_A .
- EFFICIENT MERGING ($O(k)$ BY COMPUTING MEDIAN OF $S_A \cup S_B$)
- SUPPORTS ESTIMATORS FOR SET CARDINALITIES:
 - LET V_k BE THE LARGEST HASH VALUE IN S_A
 - $\hat{n} = (k-1)R/V_k$ IS A GOOD ESTIMATOR FOR $n=|A|$

INTUITION: CONSIDER THE HASH VALUES $\{h(x) \mid x \in A\}$

- CLOSE TO EVENLY DISTRIBUTED

EXPECT $k-1$ HASH VALUES HERE



FORMAL ARGUMENT:
CHEBYCHEV (SEE BOOK).
 $k = O(1/\epsilon)$ ENOUGH FOR REL. ERROR $\leq \epsilon$

APPLICATION OF KMV

• SETS A_1, A_2, \dots, A_m (VERY LARGE)

• COLLECT KMV SUMMARIES S_{A_1}, \dots, S_{A_m}

• MERGE TO FORM KMV SUMMARY S_A OF $A = \bigcup_i A_i$

• CAN ESTIMATE:

— SIZE $|A|$ WITH RELATIVE ERROR

$\pm O(1/\sqrt{r})$, WITH CONSTANT PROBABILITY $> \frac{3}{4}$

— TAKE MEDIAN OF $O(\log(1/\delta))$ ESTIMATES TO GET ERROR PROB $\leq \delta$

— FOR GIVEN SET Q , CAN ESTIMATE $|A \cap Q|$ FROM $|S_A \cap Q|$
AND v_k

DEPENDS ONLY ON SET
OF ELEMENTS, NOT THE
MULTIPLICITY OF EACH
ELEMENT

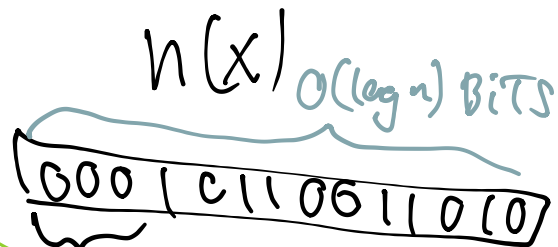
HYPERLOGLOG

- SPACE-EFFICIENT ALTERNATIVE TO KMV SUMMARY
 - DOES NOT KEEP A SAMPLE, ONLY SUPPORTS INSERTION, MERGING, CARDINALITY ESTIMATION.
- IDEA: MANY "CRUDE" ESTIMATORS USING FEW BITS

OBSERVATION. - HASH VALUES CAN BE CRUDELY ORDERED BY THE NUMBER OF LEADING ZEROS IN BIT REP:

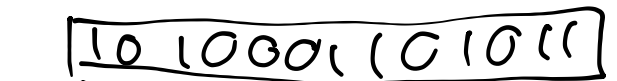
CAN BE STORED IN $\log \log(n) + O(1)$ BITS

IN PRACTICE 6-8



$$2(h(x))=3$$

$h(y)$



$$z(h(y))=0$$

HYPERLOGLOG DETAILS

$2^w \Rightarrow n$
↓
↓

- USE HASH FUNCTIONS h, g WHERE $h(x) \in \{1, \dots, k\}$ AND $g(x) \in \{0, 1\}^w$
- MAINTAIN COUNTERS c_1, \dots, c_k WHERE c_i IS A CROUD COUNTER FOR $A_i = \{x \in A \mid h(x) = i\}$: $c_i = \max_{x \in A_i} (z(g(x)))$.
- MERGING SKETCHES C, C' (SAME h, g):

NUMBER OF LEADING ZEROS

$C_i'' = \max(c_i, c_i')$ ← IDENTICAL TO THE SKETCH THAT WOULD RESULT FROM INSERTING ALL ELEMENTS DIRECTLY

FIRST ATTEMPT

- ESTIMATION: $2^{w-c_i} \approx$ SMALLEST HASH VALUE IN A_i

→ ESTIMATE $|A_i| \approx 2^w / 2^{w-c_i} = 2^{c_i}$

→ ESTIMATE $|A| = \sum_i |A_i| \approx \sum_i 2^{c_i}$

VERY SENSITIVE TO A FEW LARGE COUNTERS

BETTER: TAKE HARMONIC MEAN OF ESTIMATORS 2^{c_i}

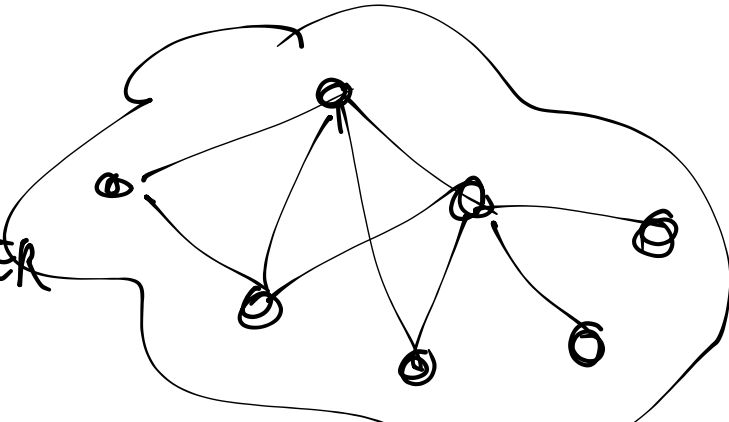
OUTLIER ROBUST AVERAGE REL. ERROR $1.06/\sqrt{k}$

SEE DISCUSSION IN BOOK

HYPERLOGLOG APPLICATION: ANF

SETTING: LARGE, UNDIRECTED GRAPH (V, E) , CONNECTED

TASK: FOR ALL $v \in V$,
COMPUTE THE AVERAGE
DISTANCE TO THE OTHER
VERTICES IN V



$|V| = 10^9$
 $|E| = 10^{11}$ ← "SPARSE", $|E| \ll |V|^2$

BASELINE: RUN BFS $|V|$ TIMES. ← INFEASIBLE ON A SINGLE MACHINE

ALTERNATIVE: SETS!

- S_i^v : SET OF VERTICES AT DISTANCE $\leq i$ FROM v .

- $S_0^v = \{v\}$, $S_{i+1}^v = \left(\bigcup_{\{u,w\} \in E} S_i^w \right) \cup S_i^v$

- $HLL(S_{i+1}^v) = \text{merge}(S_i^v, [S_i^w \mid \{u,w\} \in E])$

- COMPUTE IN TIME $O(k|E|)$ FOR ALL $v \in V$
- USE TO ESTIMATE $|S_i^v|$ FOR $i=1, \dots$?

"APPROXIMATE NEIGHBORHOOD FCT."

FEW ITERATIONS NEEDED
IF GRAPH HAS SMALL
DIAMETER

SPLIT BLOOM FILTER

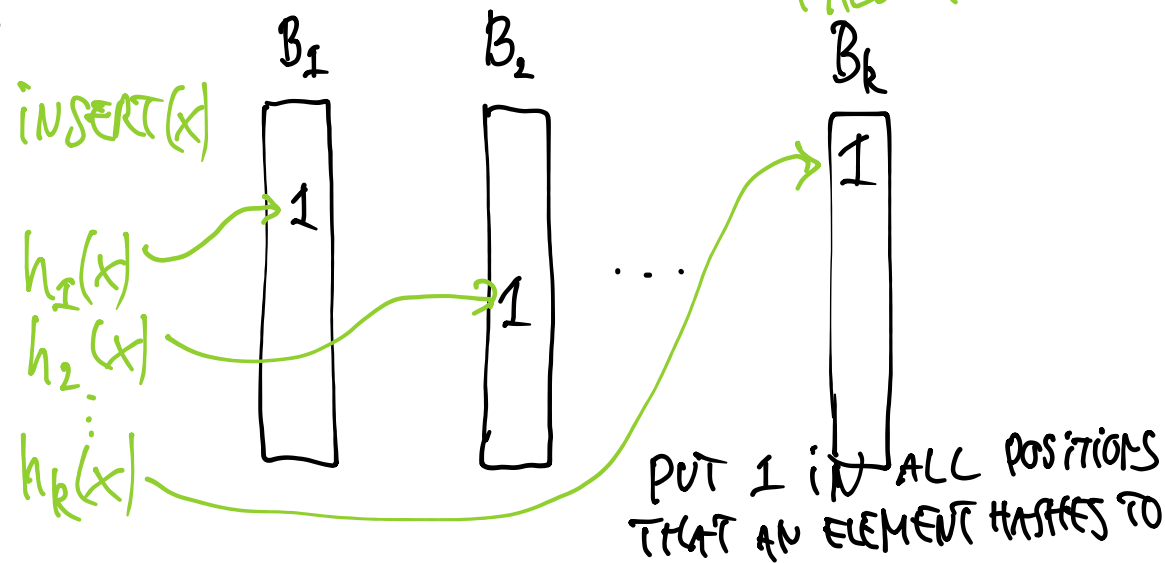
- APPROXIMATE VERSION OF A HASHPSET, PARAMETER $\epsilon > 0$
- CAN INSERT ELEMENTS INTO SET S
- CAN ASK ABOUT MEMBERSHIP QUERY $Q(x) = \begin{cases} 1 & \text{if } x \in S \\ 0 & \text{if } x \notin S \end{cases}$
- SPACE USAGE, $|S|=n$: $m = O(n \log(1/\epsilon))$ BITS

OFTEN MUCH LESS THAN
REQUIRED TO STORE ELEMENTS

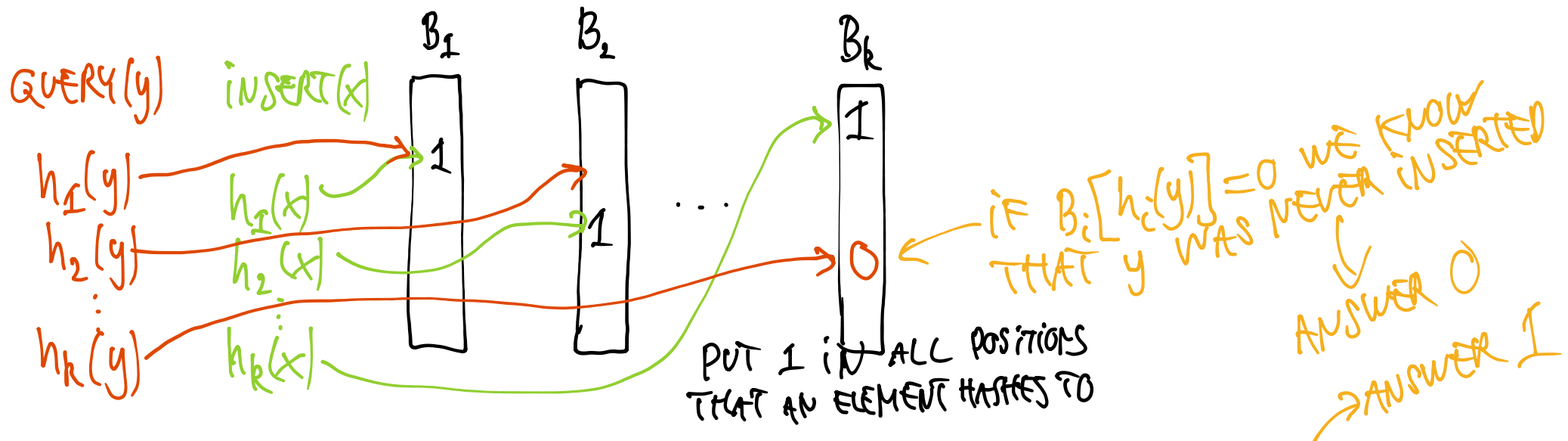
WITH PROBABILITY $\geq 1 - \epsilon$
1 IF $x \notin S$
WITH PROB. $\leq \epsilon$

• HOW IT WORKS. PARAM.

- ~ HASH FUNCTIONS h_1, \dots, h_k
- MAPPING ELEMENTS TO $\{1, \dots, m/k\}$
- BIT VECTORS B_1, \dots, B_k , EACH OF SIZE m/k



BLOOM FILTERS, CONTINUED



IF $B_i[h_i(y)] = 1$ FOR $i=1, \dots, k$, THEN y MAY HAVE BEEN INSERTED
 (NON-OPTIMAL)

PARAMETER CHOICE: LET $k = \log_2(1/\epsilon)$, $m = 2kn = 2n \log_2(1/\epsilon)$.

FOR y NEVER INSERTED:

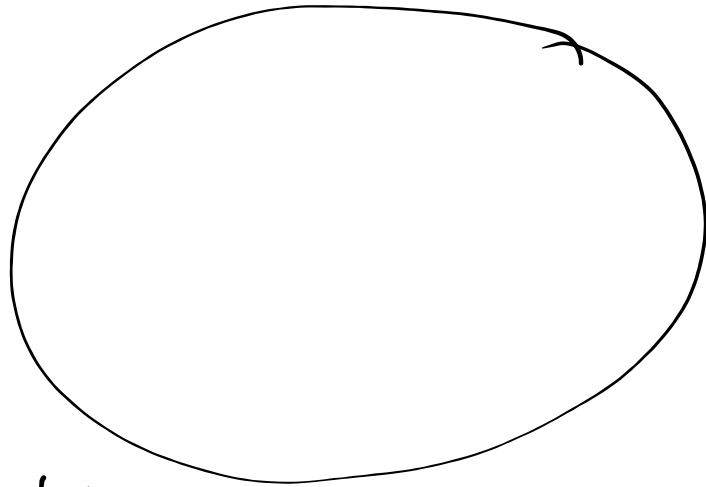
$$\Pr[B_i[h_i(y)] = 1 \text{ FOR } i=1, \dots, k] = \prod_{i=1}^k \Pr[B_i[h_i(y)] = 1] \leq \prod_{i=1}^k \frac{1}{2} = \epsilon.$$

↑
 AT MOST n OUT OF $2n$ ENTRIES ARE 1

NEED TO KNOW (BOUND ON) SET SIZE IN ADVANCE

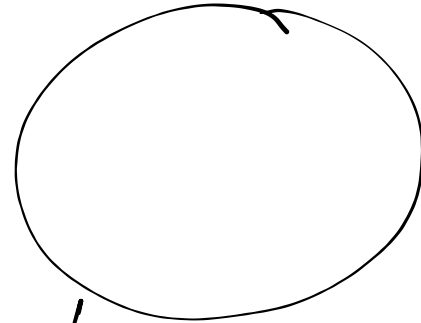
BLOOM FILTER APPLICATION

ESTIMATING THE NUMBER OF COMMON ELEMENTS

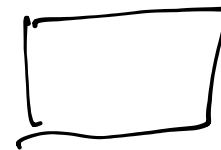


LARGE SET A

HOW LARGE
IS $A \cap B$?



SMALLER SET B, BUT
DOES NOT FIT IN MEMORY



BLOOM FILTER FOR
B, FITS IN MEMORY

ALGORITHM

FOR EACH $x \in A$, LOOK UP x
IN BLOOM FILTER FOR B , t POSITIVES

ESTIMATE $|A \cap B| \approx t - \epsilon |A|$

↑
EXPECTED NUMBER OF
FALSE POSITIVES