

Hashing & Streaming Algorithms

$$\underline{h}: [N] \longrightarrow [M]$$

$$h \in \mathcal{H}$$

$$N \gg M$$

\mathcal{H} Universal family of hash functions:

$$\forall x \neq y \quad P_{h \in \mathcal{H}}[h(x) = h(y)] = \frac{1}{M}$$

e.g. $x = x_1 x_2 x_3 x_4$

$$x_i \in [256]$$

Instead $p = 257$

Pick $a = (a_1, a_2, a_3, a_4)$ random mod 257

$$\underline{h_a}(x) = a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 \pmod{p} \quad \boxed{\quad}$$

$$|\mathcal{H}| = 257^4$$

pairwise indep? No.

Pairwise independent family of hash functions:

$$\mathcal{H} \quad \forall x \neq y, u, v \quad P[h(x)=u \text{ and } h(y)=v] = \frac{1}{M^2}$$

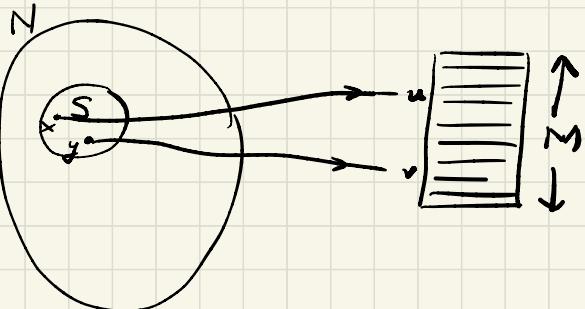
pick $a = a_1 a_2 a_3 a_4 \pmod{p}$

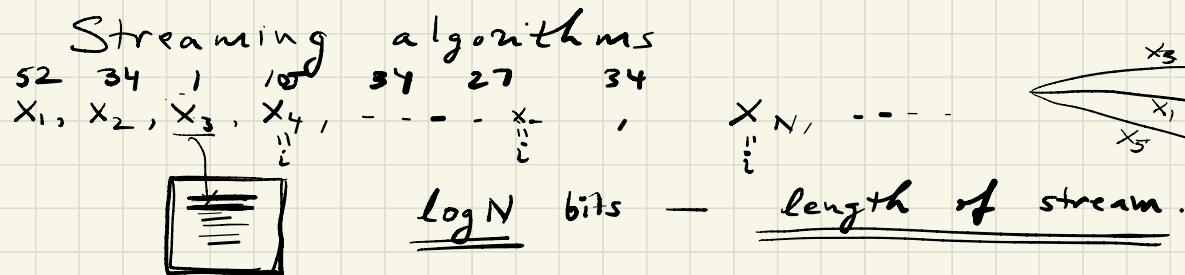
$$\mathcal{H}' \quad h'_{a,b}(x) = a_1 x_1 + a_2 x_2 + a_3 x_3 + a_4 x_4 + b \pmod{p}$$

$$|\mathcal{H}'| = 257^5$$

pairwise indep \Rightarrow universal

$$[N] = \{0, 1, \dots, N-1\}$$





Frequencies: $i = m_i$ times

$$m_{34} = 3$$

Frequency moments: $F_K = \sum_i m_i^K$

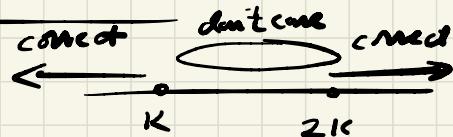
$F_0 = \sum_i m_i^0 = \# \text{ distinct elements in stream.}$

$F_1 = \sum_i m_i = \text{length of stream}$

$F_2 = \sum_i m_i^2 \approx \text{"variance"}$

distinct elt ... thin factor of 2.

$K:$ # dist elts $\leq K$ — Small



dist elts $\geq 2K$ — large

$K < \# \text{ dist elts} < 2K$ — arbitrary answer.

Solution: Pick h random universal hash for:

Check $\underline{h(x_i)} = 0?$

$h'(x_i) = 0?$

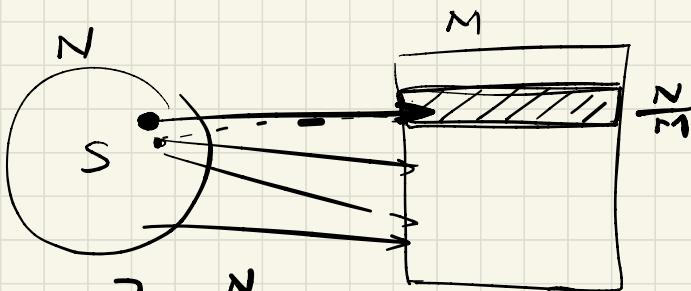
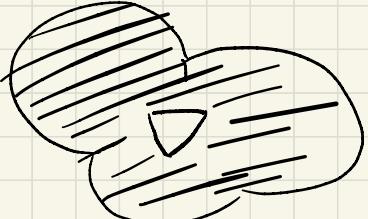
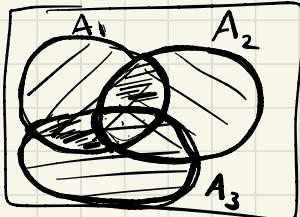
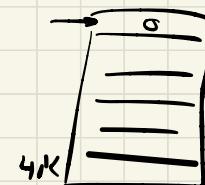
Assume $\geq 2K$ distinct elts

$P[\text{one of them hashes to 0}] \geq \frac{3}{8}$.

$i = 1, \dots, 2K$

$A_i = \text{event } i \text{ hashes to 0.}$

$$\begin{aligned}
 P[A_i] &\geq \sum_i P[A_i] - \sum_{i,j} P[A_i \cap A_j] \\
 &= 2K \times \frac{1}{4K} - \binom{2K}{2} \frac{1}{(4K)^2} \\
 &= \frac{1}{2} - \frac{\binom{2K}{2}}{2} \frac{1}{(4K)^2} \\
 &= \frac{1}{2} - \frac{1}{8} = \frac{3}{8}.
 \end{aligned}$$



Variance calculation is valid for pairwise independence.

$$E[\# \text{ in bin } j] = \frac{N}{M}$$

$$\text{Var}(\# \text{ in bin } j) = \left(\frac{1}{M}\right)\left(1 - \frac{1}{M}\right)N$$

Chebyshev