Perfect Hash.   ⊖——  no collision.

$|S| = N$

Pigeon

U

$O(N)$ holes.

$N=2N$

$N=1$

Each bin $j$ contains $n_j$ balls.

$|S| = N$

Throw N balls to N bins. It's known:

$$Exp\left(\sum_{j=0}^{N-1} n_j^2\right) \leq 2N.$$

We also have:

Prob$\left\{\sum_{j=0}^{N-1} n_j^2 > 4N\right\} < \frac{1}{2}.$

Two steps.

(1). $H$: universal hash family —
$h$ is randomly from $H$.

$$|s| = N$$



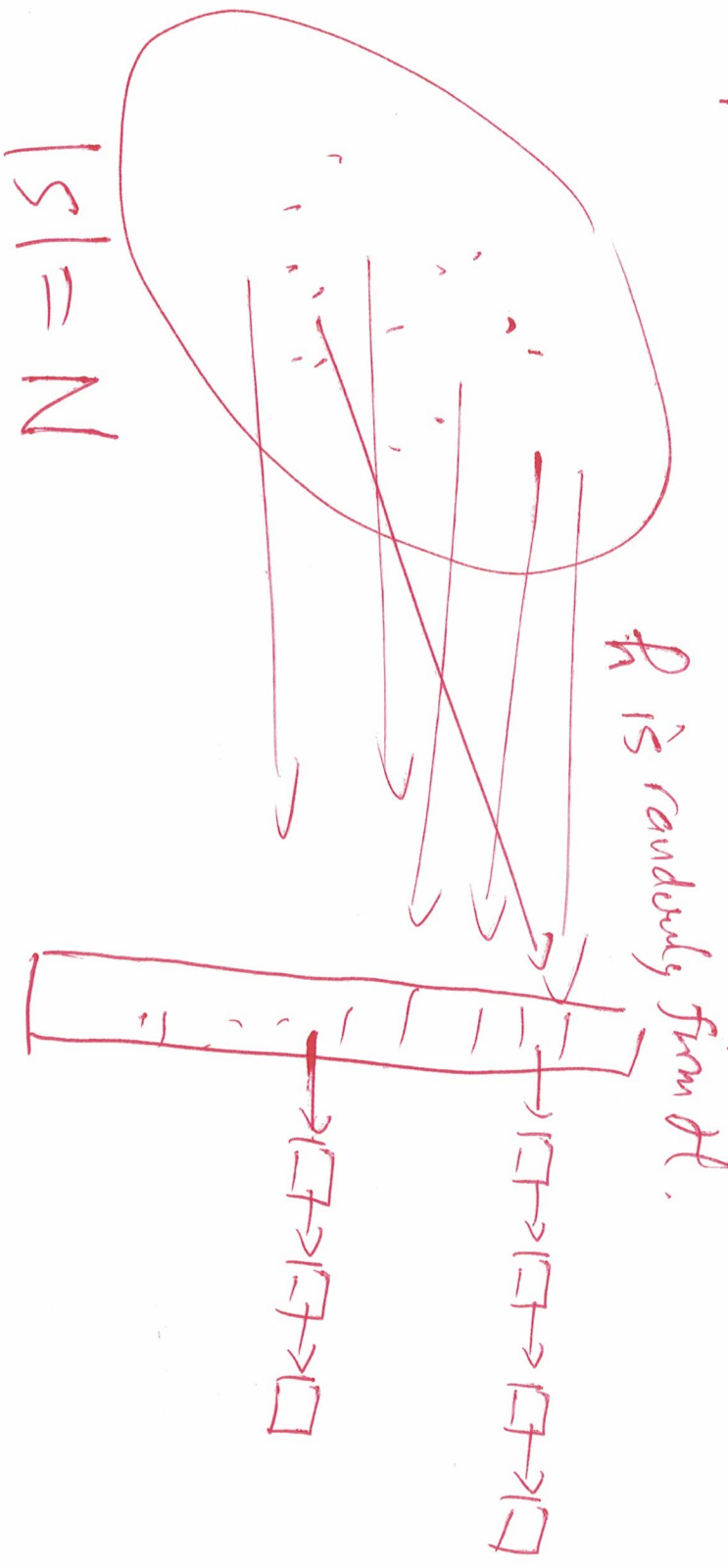$$M = O(N) \neq N.$$

(2). hash the collisions again (using the $O(N)^2$
no collision method).

# Perfect Hashing Alg:

(1). Repeat

Select $h \in \mathcal{H}$
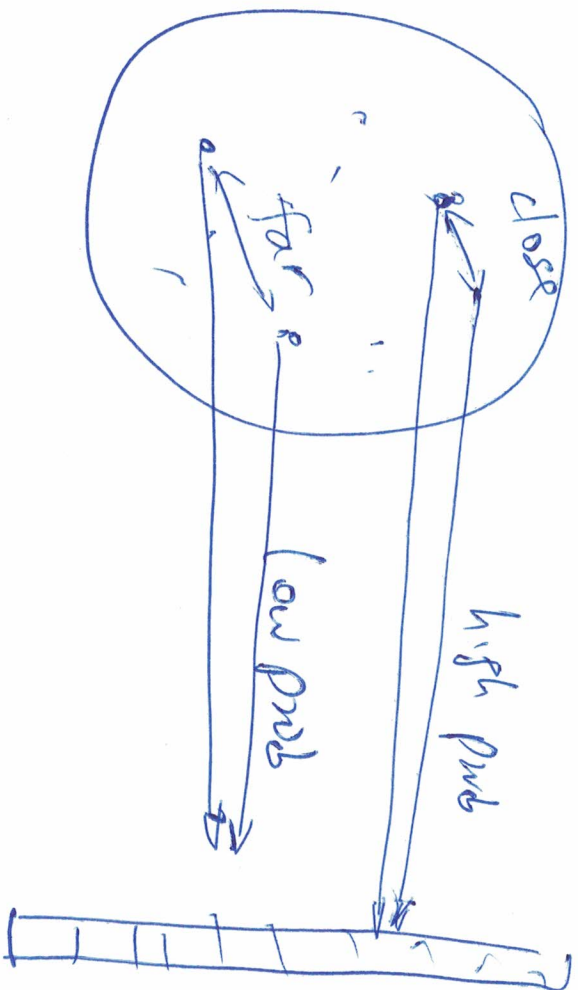
Until $\sum_{j=0}^{N-1} n_j^2 \leq 4N$.

(2). For each j-th slot with $n_j$ keys.

$(n_j > 1)$. // collision case

Select $h_j$ such that the $n_j$ keys

are hashed into $n_j^2$ slots with no

collision.

total space: $\leq 5N$

# LOCality Sensitive Hashing.



close

far

high prob

low prob

$H$: hash family $U \longrightarrow N$. where $U$ is a metric
space. Let $h$ be random on $H$ and
$H$

(1). if $d(p_1, p_2) \leq r_1$ Then

$$\text{Prob}\{h(p_1) = h(p_2)\} \geq P_1,$$
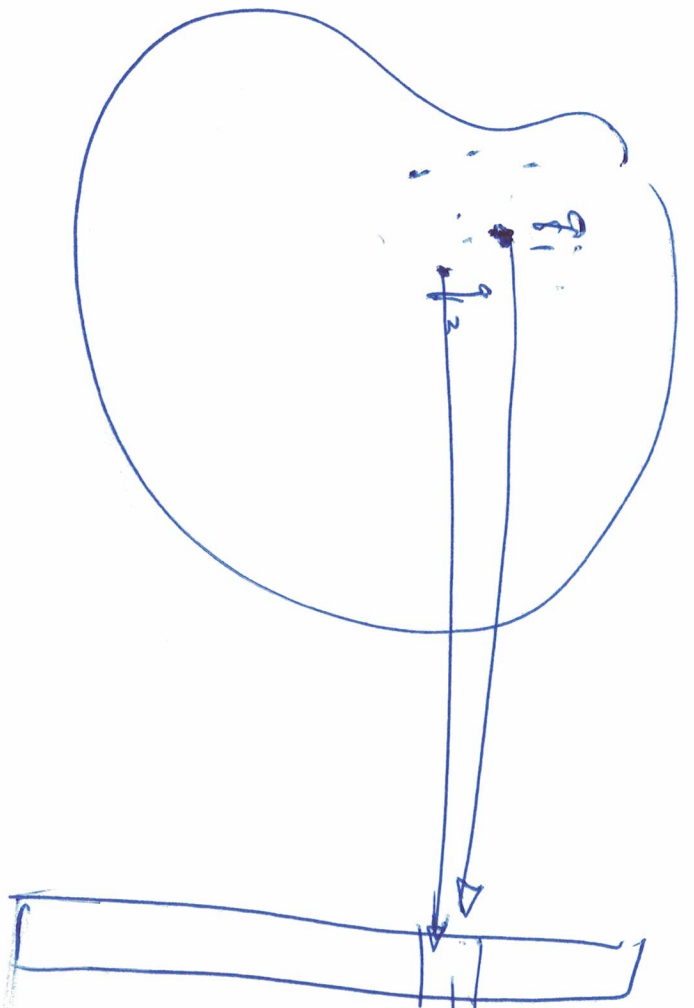
(2). if $d(p_1, p_2) \geq r_2$ then

$$\text{Prob}\{h(p_1) = h(p_2)\} \leq P_2,$$

where $p_1, p_2 \in \mathcal{U}$, $P_1 > P_2$, and $r_1 < r_2$.

Then, $H$ is called $(r_1, r_2, P_1, P_2)-\text{sensitive}$.

Given $p_1, \dots, p_n$ where $n$ is large find those $(p_i, p_j)$ pairs with $d(p_i, p_j)$ Small. How?

$q_i$

$q_2$

collisons with light
(the points with prob. stay close)

Intuition.

$$h_{a,b}(v_1) = h_{a,b}(v_2) \text{ with } v_1, v_2 \text{ fixed}$$
and $a, b$ random.



$$a \cdot v_1 + b r \qquad a \cdot v_2 + b r$$

if

(1).
$$\frac{a \cdot (v_1 - v_2)}{r} \in [0, 1) \qquad // \ D \in [0, r) .$$
$$\leftarrow r \rightarrow \leftarrow r \rightarrow$$

(2).
$$b r \in [0, r - D].$$

Read SCG'04 with The poof) derived
$h_a$ (1) and (2).

(Datar etc., SCG '04).

Let $U = \mathbb{R}^d$. $H$ : The family of all

$$f_{a,b}(\cdot) \ , \quad \text{where}$$

① $a \in \mathbb{R}^d$ with each component $a_i$ is chosen randomly from normal distribution (with $\mu = 0, \sigma = 1$).

② $b$ is a random number is $[0,1]$.

③ $r$ is some positive number.

④ $$f_{a,b}(v) = \left\lfloor \frac{a \cdot v}{r} + b \right\rfloor$$

$\forall \ v \in \mathbb{R}^d$, a vector

# Bloom filter.

( a data structure to rep. a set ).

Define

$$h(a, x) = \left( \text{rotate } x \text{ to the} \atop \text{left by } a \text{ bits} \right) \mod 16$$

Exp]o.

$$x = 1001100010010$$
$$a = 3.$$
$$1100010010\underbrace{0100} \mod 16$$
$$h(a, x) = 0100 = 4.$$

We are given:

$$h_1 = h(2, x)$$

$$h_2 = h(4, x)$$

$$h_3 = h(3, x).$$

Consider a set $\{x_1, x_2\} = \{11000100,$
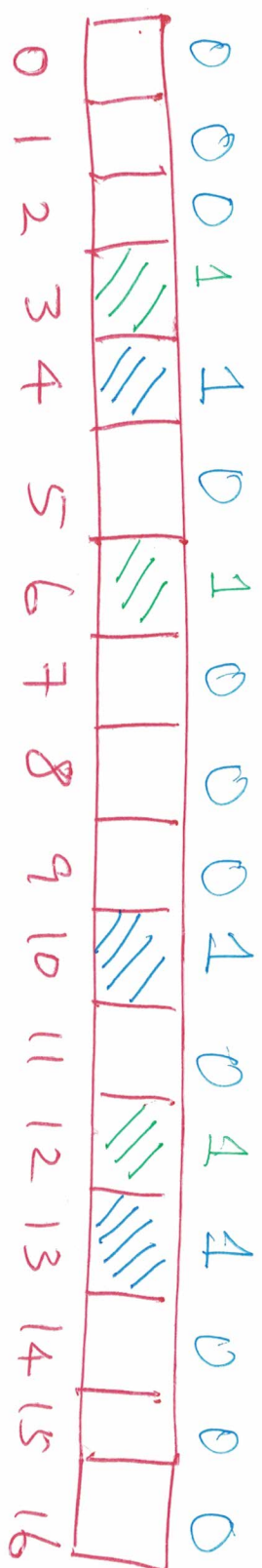$01000100011\}$.

| | |
|---|---|
| $h_1(x_1) = 3$ | $h_1(x_2) = 13$ |
| $h_2(x_1) = 12$ | $h_2(x_2) = 4$ |
| $h_3(x_1) = 6.$ | $h_3(x_2) = 10$. |

I add $x_1$ to Bloom filter.

$$0\ 0\ 0\ 1\ 0\ 0\ 0\ 1\ 0\ 0\ 0\ 1\ 0\ 1\ 0\ 0\ 0$$

I add $x_2$ to Bloom filter:

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|

Given: $x$

Query: $x \in$ the set?

Compute $h_1(x)$, $h_2(x)$, $h_3(x)$.

and see if the $h_1(x)$-bit, $h_2(x)$-bit and $h_3(x)$-bit are all set to 1. in the Bloom filter. If yes, $x \in$ the set not $x \notin$ the set.

False Positive: when query $x$ for $h$

$h_1(x), \ldots, h_k(x)$, all the bit $= 1$

at positions $h_1(x), \ldots, h_k(x)$, the probability

of false positive:

$$\left(1 - \left(1 - \frac{1}{m}\right)^{kn}\right)^k \approx \left(1 - e^{-\frac{kn}{m}}\right)^k.$$

In practice, how to choose the $k$?

$$k \approx \frac{m}{n} \cdot \ln 2.$$