How to measure "randomness" in a r.v.

Suppose that $X$ takes val's $a_1, \ldots, a_k$ with probability $p_1, \ldots, p_k$ resp. Shannon's entropy

$$H(X) = -\sum_i p_i \log p_i \geq 0.$$

Example: Fair coin $X$:

$$Prob(X = head) = \tfrac{1}{2}$$

$$Prob(X = tail) = \tfrac{1}{2}.$$

$$H(X) = -\tfrac{1}{2}\log_2 \tfrac{1}{2} - \tfrac{1}{2}\log_2 \tfrac{1}{2} = 1 \text{ bit.}$$

A memory cell that holds $1$ with probability $90\%$ and $0$ with probability $10\%$

is of memory size $\ll 1$ bit !!!

In practice, the distribution of $x$ is very hard to obtain; that is, the $p_1, \ldots, p_k$ are unknown. Then, how to estimate the randomness $H(x)$ of $x$?

many ways:

① Sampling. (read/take prob & statistics).

② Lompel-Ziv gives the answer.

// Note: Lompel-Ziv has a condition: The sampled
// sequence is "taken" or "generated" by an ergodic
// Markov Chain.

Step 1. Produce a "random" sequence of values
  $z_1, \ldots, z_n$,  each from dist. of $x$.

(Classic sampling.)

// Example: if you know $x$ is a coin with
// unknown dist'ns, then you toss it for $n$
// $n$ shall be larger — half million.
// thus, find obtain the outcomes $x_1, \ldots, z_n$.

Step 2. Run Lempel-Ziv (LZ) on $x_1, \dots, x_n$.

Step 3. Compression Ratio can be used to estimate $H(X)$.

---

Remark. ① LZ is the, for exple, LZMA in 7Z.

or gzip, ...
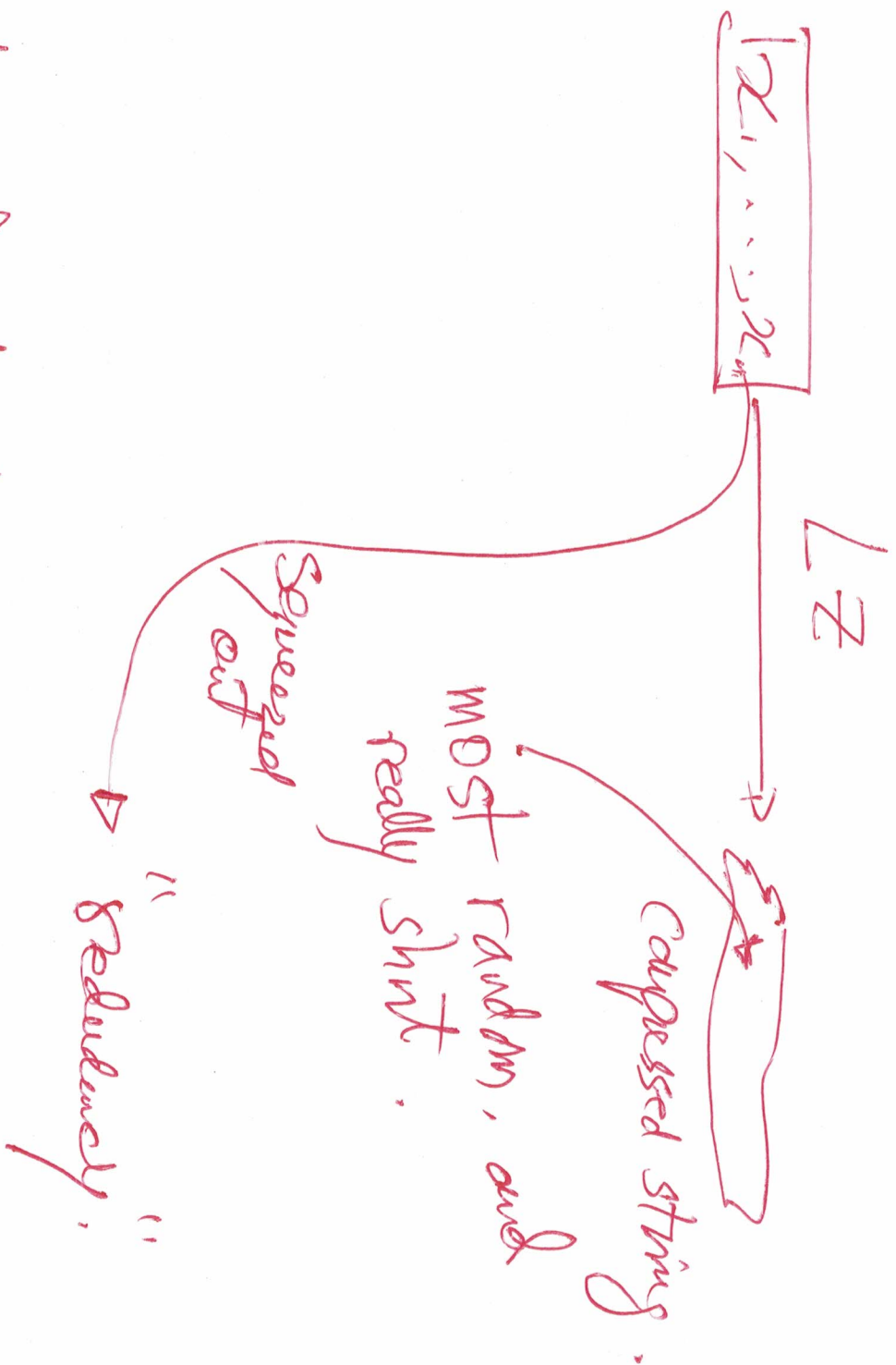
② LZ is optimal. (with side condition).

③ In a sequence $x_1, \dots, x_n$, randomness (entropy) is what is left when "redundancy" is reduced.

is squeezed out.

$x_1, \ldots, x_n$

LZ

compressed string.

Squeezed out
MOST random, and
really shirt.

"Redundancy".

④ LZ is lossless.

Read: "Lossless of covariate channels modeled by Transducers".

Hash, universal Hash, Locality-
Sensitive hash'g & Bloom Filter.

Big Picture.
a few things chosen from many things.
$\longrightarrow S$
$\nearrow U$

Wait: a short name for each thg that
is chosen.

_names are [M]_.

Math'cal set up.
$[M] = \{0, 1, 2, \dots, M-1\} \Leftarrow M$ indices
$U$: universe, a very large set of keys.
(e.g., $U = \Sigma^k$ th size $|\Sigma|^k$).

Ex: $\Sigma = \{0, 1\}$, $k = 100$. Then

$U$ is of all 100-bit strings. How many?

$$|U| = |\Sigma|^{100} = 2^{100}.$$

$S \subseteq U$, a subset of the Universe. (keys?)

Let $|S| = N.$  $(N \leq |U|).$

Consider $U$ as the set of all 4-bit

strings. So, $U = \{0000, 0001, \ldots, 1111\}$.
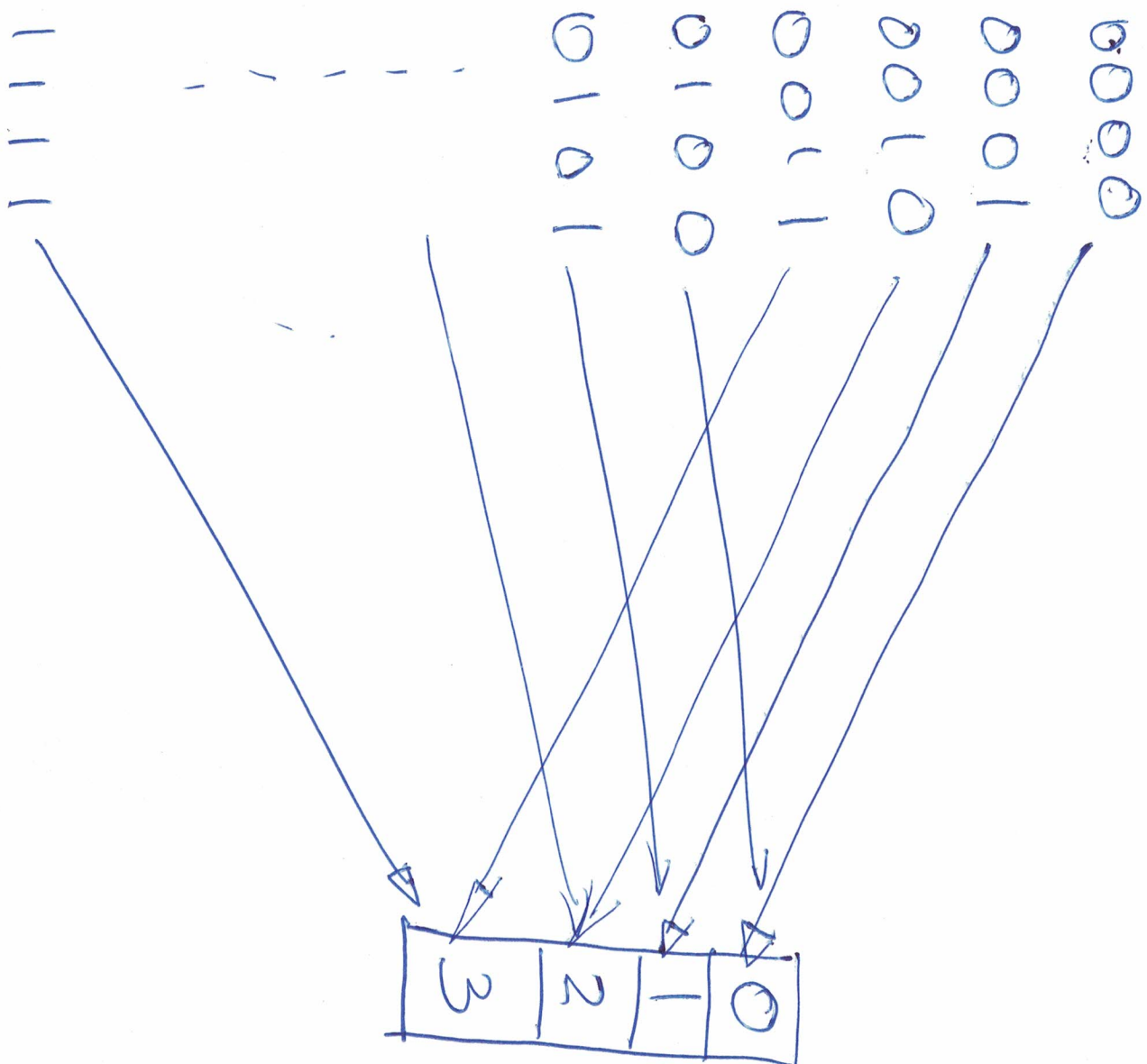
$$|U| = 2^4 = 16.$$

Let $M = 4$. ( 4 indices).

$h: U \longrightarrow [M] = \{0, 1, 2, 3\}$ defined as

$$h(a_1 a_2 a_3 a_4) = a_3 a_4 . \quad \text{(taking the last two bits).}$$

For example, $h(0010) = 10_{\text{in binary}} = 2$.

Universe $U$ $\xrightarrow{h}$ $[M]$

$1\,1\,1\,1$

$-\;-\;-\;-\;-\;-\;-$

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

| 3 | 2 | 1 | 0 |

Nightmare:

Take $S = \{0000, 0100, 1000\}$. Th

we have

$$N = |S| = 3.$$

we have $h(x) = 00$ for each $x \in S$.

That is, all elements in $S$ share the same $C$ keys.

A collision refers to this scenario:

$$\underbrace{h(x) = h(y)}_{\text{same index}} \text{ but } x = y.$$

Ideal case:

$$f(x) \neq h(y)$$

whenever $x \neq y$.

for all $x, y \in S$, and for all $S \subseteq U$

when $|S| = N \approx M \ll |U|$.

---

Intuitive Explanation:

In a party of 16 kids. I want to find
a way to assign short names to each kid
(e.g., each kid's name has only two bits)
such that any subset of 4 kids would have
distinct names. Possible? NO. kids in

$f$  ← $f$

← Universe.

$N = M$

$\leftarrow$ any subset

$= =$

Theorem. Let $|\mathcal{U}| > (N-1)M + 1$. Then
for each hash funct= $h$, There is a subset
$S$ with $|S| = N$ and
$$h(x) = h(y)$$
for all $x, y \in S$.

Proof. Pigion-hole. Read book.

Reason: we try to map $|S| = N$ keys to different locations/indices with minimal # of collisions possible. However, if $h$ is known, we can always pick a set of keys, $S$, to let $h$ fail terribly as is the theme.

Solution: to make $h$ unknown. How? But we have to know $h$. What to do? "$h$ is known But it's random."