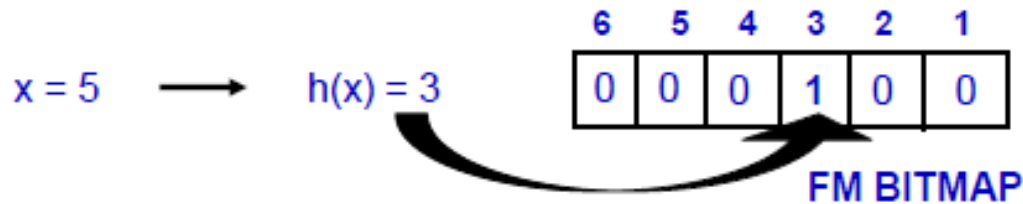# CS6931 Database Seminar

Lecture 5: Streaming Model (continued)

# The Distinct Counting Problem
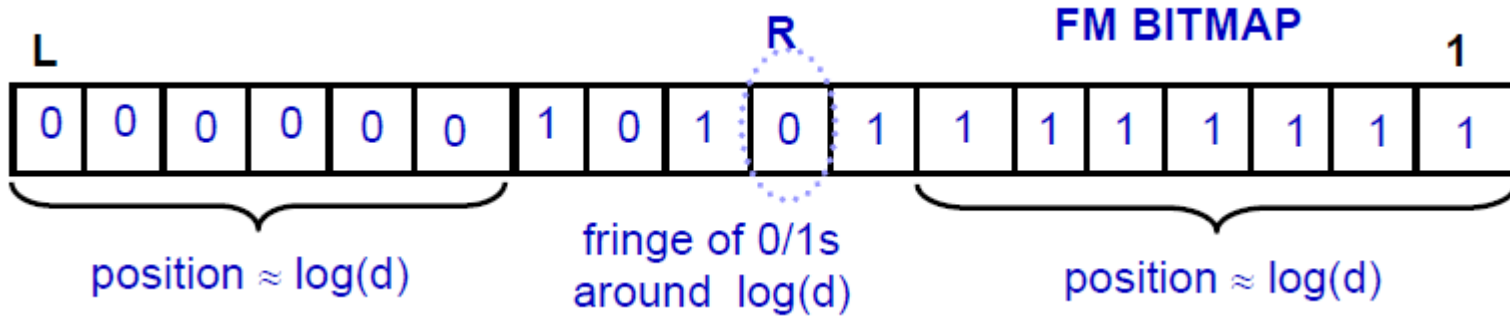
# FM Sketch

- ‚Estimates number of distinct inputs (count distinct)
- ‚Uses hash function mapping input items to i with prob $2^{-i}$
  - i.e. $Pr[h(x) = 1] = ½$, $Pr[h(x) = 2] = ¼$, $Pr[h(x)=3] = 1/8$ …
  - Easy to construct h() from a uniform hash function by counting trailing zeros
- ‚Maintain FM Sketch = bitmap array of L = log U bits
  - Initialize bitmap to all 0s
  - For each incoming value x, set $FM[h(x)] = 1$

|  | 6 | 5 | 4 | 3 | 2 | 1 |
|---|---|---|---|---|---|---|

x = 5 ⟶ h(x) = 3

| 0 | 0 | 0 | 1 | 0 | 0 |
|---|---|---|---|---|---|

FM BITMAP

# FM Sketch

- If d distinct values, expect d/2 map to FM[1], d/4 to FM[2]…
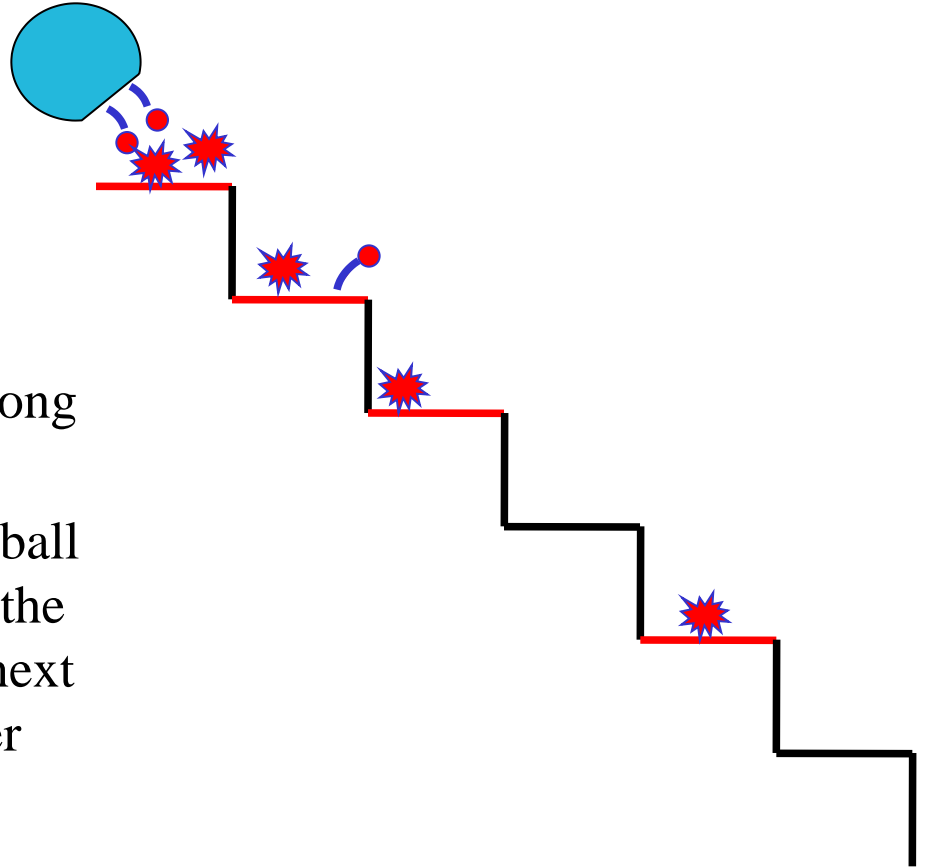


- Let R = position of rightmost zero in FM, indicator of log(d)
- Basic estimate d = c2R for scaling constant c ≈ 1.3
- Average many copies (different hash fns) improves accuracy

- With $O(1/\varepsilon 2 \log 1/\delta)$ copies, get $(\varepsilon,\delta)$ approximation
  - 10 copies gets ≈ 30% error, 100 copies < 10% error

# COUNT Sketches

- Problem: Estimate the number of distinct item IDs in a data set with only one pass.

- Constraints:
  - Small space relative to stream size.
  - Small per item processing overhead.
  - Union operator on sketch results.

- Exact COUNT is impossible without linear space.

- First approximate COUNT sketch in [FM'85].
  - $O(\log N)$ space, $O(1)$ processing time per item.

# **Counting Paintballs**

- Imagine the following scenario:
  - A bag of $n$ paintballs is emptied at the top of a long stair-case.
  - At each step, each paintball either bursts and marks the step, or bounces to the next step. 50/50 chance either way.
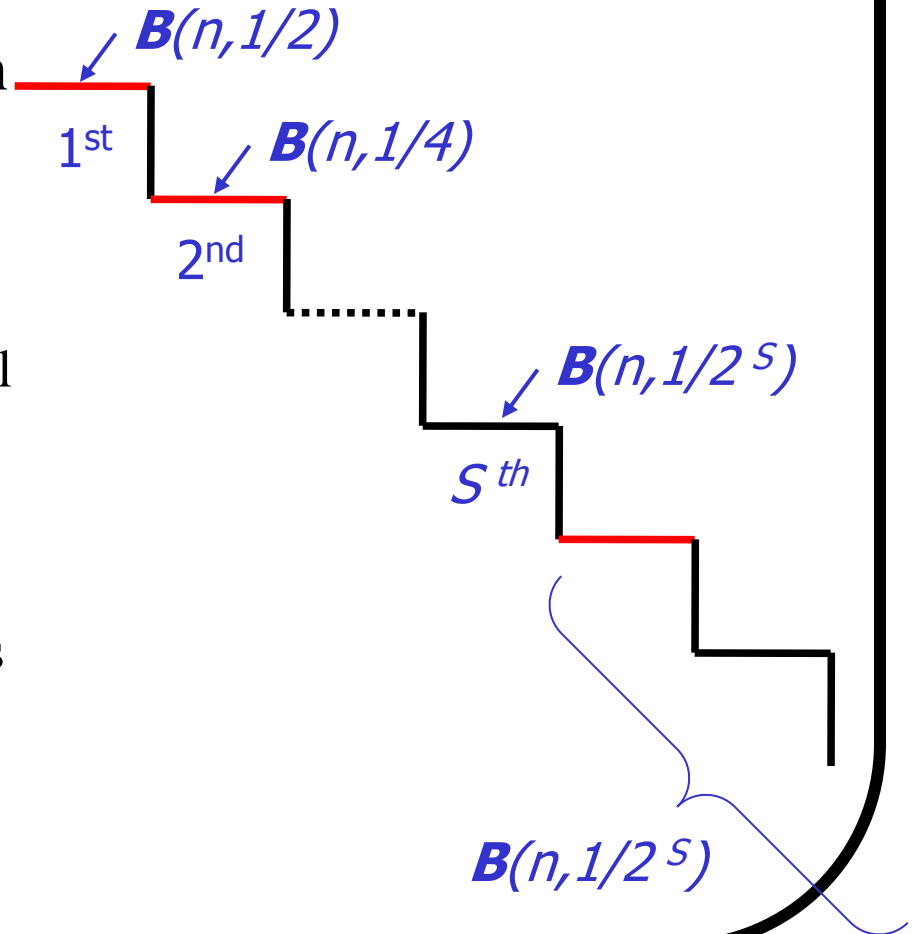
Looking only at the pattern of marked steps, what was $n$?

# Counting Paintballs (cont)

- What does the distribution of paintball bursts look like?

  – The number of bursts at each step follows a binomial distribution.

  – The expected number of bursts drops geometrically.

  – Few bursts after $log_2\, n$ steps

$B(n,1/2)$

1st

$B(n,1/4)$

2nd

$B(n,1/2^S)$

$S^{th}$

$B(n,1/2^S)$

# Counting Paintballs (cont)

- Many different estimator ideas [FM'85,AMS'96,GGR'03,DF'03,...]
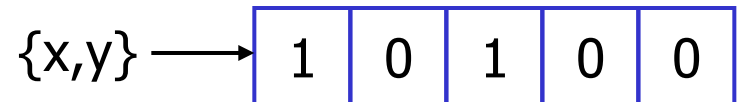- Example: Let *pos* denote the position of the highest unmarked stair,

$$E(pos) \approx log_2(0.775351\ n)$$

$$\sigma^2(pos) \approx 1.12127$$

- Standard variance reduction methods apply
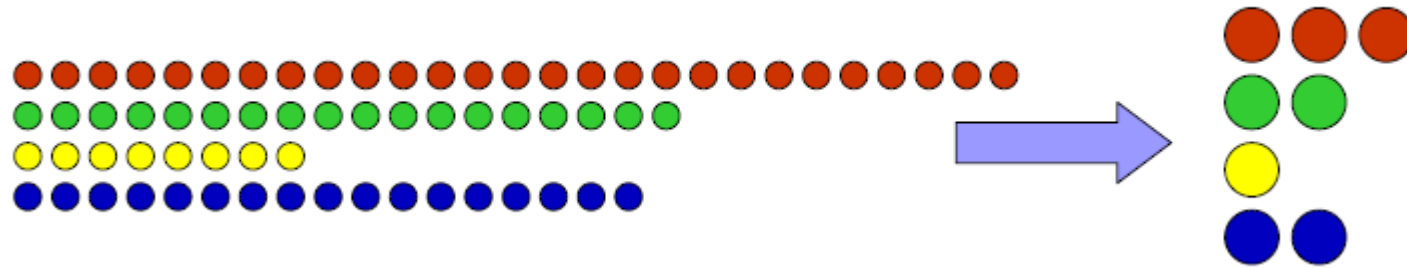- Either O(log n) or O(log log n) space

# Back to COUNT Sketches

- The COUNT sketches of [FM'85] are equivalent to the paintball process.
  - Start with a bit-vector of all zeros.
  - For each item,
    * Use its ID and a hash function for coin flips.
    * Pick a bit-vector entry.
    * Set that bit to one.
- These sketches are duplicate-insensitive:

$$\forall A, B \ (Sketch(A) \mathbin{/\!/} Sketch(B)) = Sketch(A \cup B)$$

{x} $\longrightarrow$ | 1 | 0 | 0 | 0 | 0 |

{y} $\longrightarrow$ | 0 | 0 | 1 | 0 | 0 |

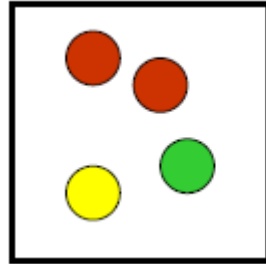{x,y} $\longrightarrow$ | 1 | 0 | 1 | 0 | 0 |

The Sampling Problem

# Sampling From a Stream



- Fundamental prob: sample m items uniformly from stream
  - Useful: approximate costly computation on small sample
- Challenge: don't know how long stream is
  - So when/how often to sample?
- Two solutions, apply to different situations:
  - Reservoir sampling (dates from 1980s?)
  - Min-wise sampling (dates from 1990s?)

# **Reservoir Sampling**



- Sample first m items
- Choose to sample the i'th item (i>m) with probability m/i
- If sampled, randomly replace a previously sampled item
- Optimization: when i gets large, compute which item will
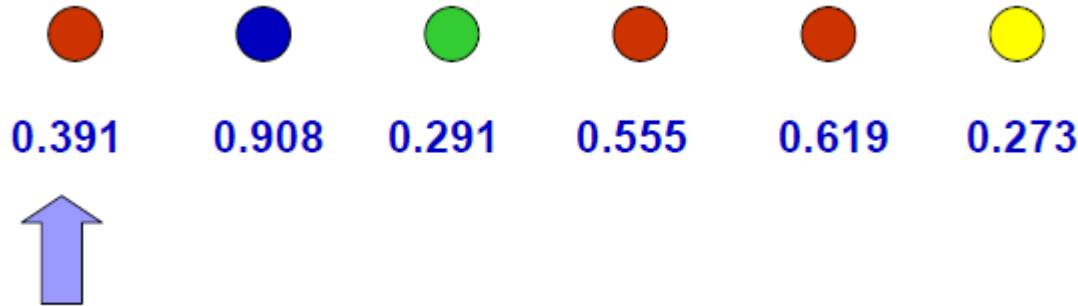  be sampled next, skip over intervening items. [Vitter 85]

# Reservoir Sampling

- Analyze simple case: sample size m = 1
- Probability i'th item is the sample from stream length n:
  - Prob. i is sampled on arrival $\times$ prob. i survives to end

$$\frac{1}{i} \times \frac{i}{i+1} \times \frac{i+1}{i+2} \cdots \frac{n-2}{n-1} \times \frac{n-1}{n}$$

$$= 1/n$$

- Case for m > 1 is similar, easy to show uniform probability
- Drawbacks of reservoir sampling: hard to parallelize

# Min-wise Sampling

- For each item, pick a random fraction between 0 and 1
- Store item(s) with the smallest random tag [Nath et al.'04]



| 0.391 | 0.908 | 0.291 | 0.555 | 0.619 | 0.273 |

- Each item has same chance of least tag, so uniform
- Can run on multiple streams separately, then merge