

# Advanced Algorithms

## Lecture 16: Balls and bins

# Announcements

- **HW 4 out — due next Friday** (Graphs & shortest paths) .
- Mid-term poll: discussion
- Office hrs today — 1 PM → 2 PM .

# Randomized algorithms: analysis

- In some cases: trade-off between running time & probability of correctness (e.g., identity testing, ...)
- Las Vegas algorithms: always correct, but running time can sometimes be large (e.g., quick sort)

(running times are random variables)

# Expected running time

- Run time is a random variable — e.g., quick sort
  - Choose random number of the array as the “pivot”, divide array into two parts, sort recursively
  - As long as split is “roughly balanced”, problem size reduces significantly

- Can write a **recurrence** for the expected running time

$f(n)$ : expected runtime on array of size  $n$ .

- Law of “conditional expectation”

$$\mathbb{E}[X] = p(F) \cdot \mathbb{E}[X|F] + (1 - p(F)) \cdot \mathbb{E}[X|\overline{F}]$$

$\hookrightarrow \int x \operatorname{pr}[X=x] dx$

# Expectation good enough?

- Suppose that the expected running time is  $O(n \log n)$
- Can we upper bound probability that it is (say)  $n^2$  ?  $\rightarrow$
- What about  $n^{1.5}$  ?

$$\Pr\{\text{run time} > n^{1.5}\} < \frac{1}{100} ?$$

# General result

**Markov's inequality:** let  $X$  be a non-negative random variable with expectation  $C$ . Then  $\text{prob}[X > tC] \leq 1/t$ .  $\forall t \geq 1$

$$\text{Pr}[X > 10 \cdot C] \leq \frac{1}{10}$$

- Implication for quick sort?



$$C = 4n \log n$$

- “Boosting probability”



what if we want success prob. of 0.9999

$X$ : running time of quicksort on instance of size  $n$ .

$$\text{Pr}[X > 10 \cdot 4n \log n] \leq \frac{1}{10}$$

[w.p.  $\geq 0.5$ , running time  $\leq 2 \cdot C$ ]

What is a bound on  $\text{Prob}[X \geq \frac{n^2}{2}]$ .

$$\mathbb{E}[X] := C = n \log n$$



$$t = \frac{n}{2 \log n}$$

$$\Pr[X \geq t \cdot C] \leq \frac{1}{t} \implies \Pr[X \geq \frac{n^2}{2}] \leq \frac{2 \log n}{n}$$

What is a bound on  $\Pr[X > 2n]$ ?  $\rightarrow$  ~~a~~ no meaningful bound.

Goal: get much higher success prob.

# Amplification by repetition

"Basic" Markov:  $\text{Prob}[X > 2 \cdot n \log n] \leq \frac{1}{2}$ .

→ Impatient quicksort:

- try <sup>100</sup> times:

- run quicksort for  $2n \log n$  steps
- if it doesn't complete, terminate
- if it completes, ~~do~~ break;

prob.  $\leq \frac{1}{2}$

Prob that we fail all the time  $\leq \frac{1}{2^{100}}$

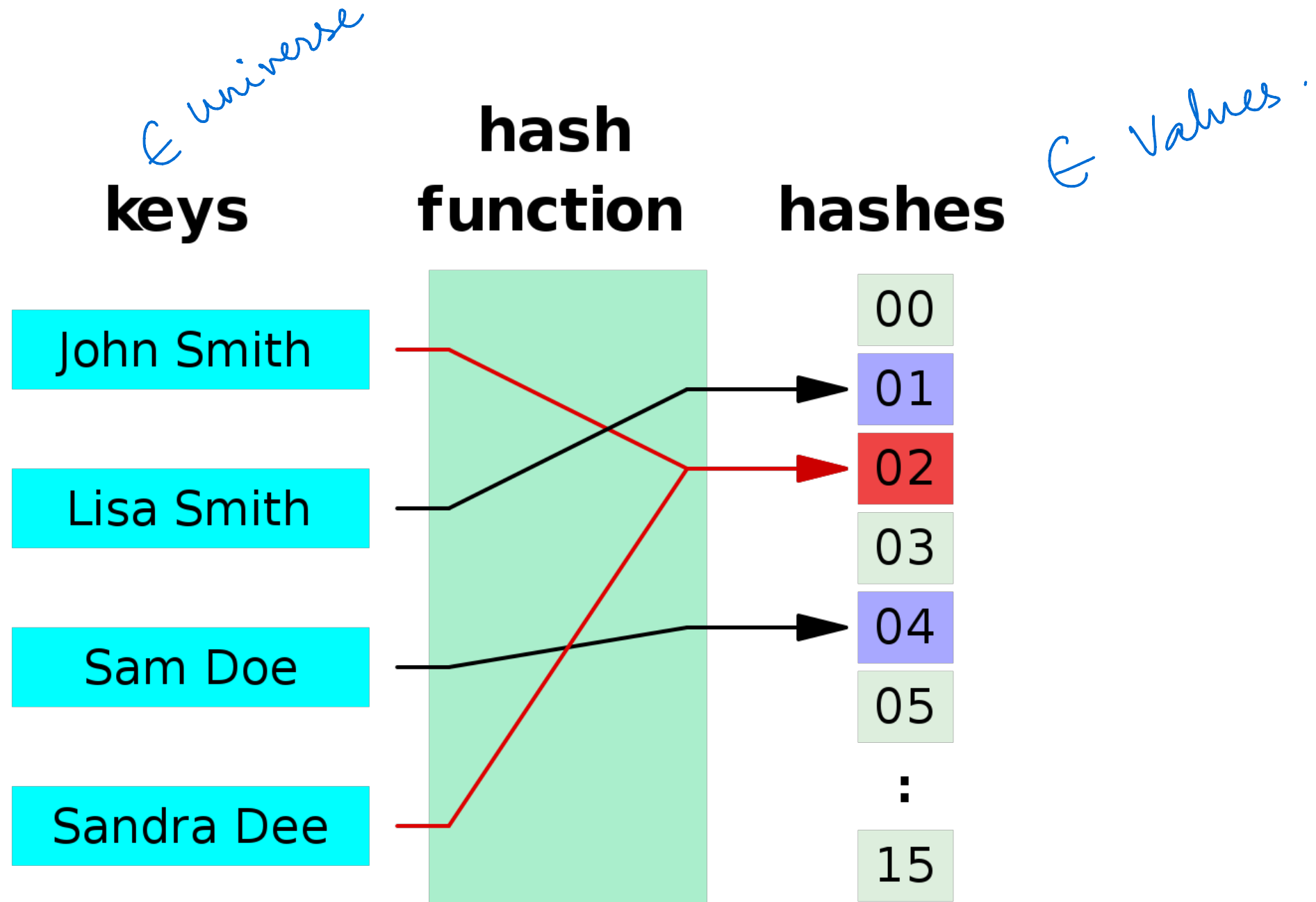
total running  
time =  
 $200 \cdot n \log n$



# Today's plan

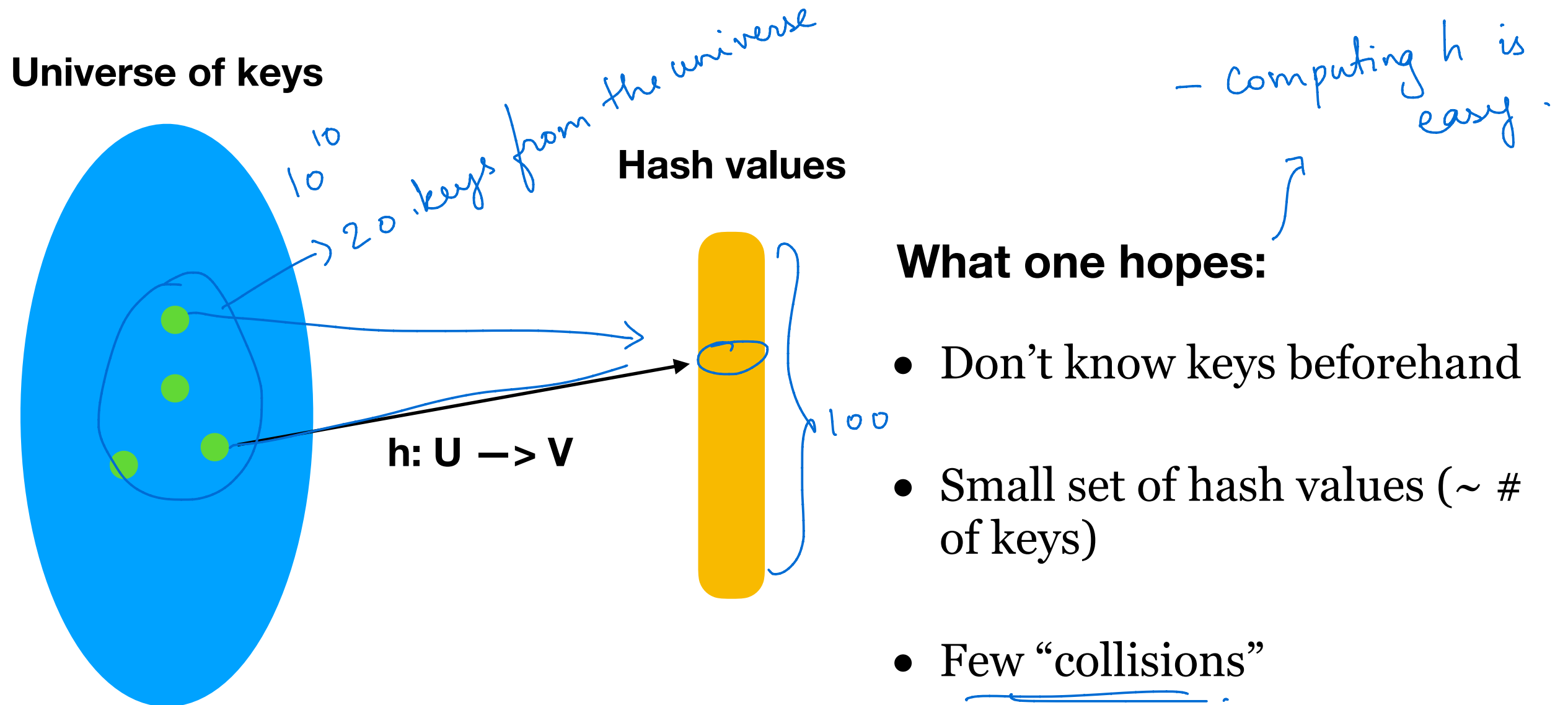
- Understanding hashing via “balls and bins”
- Linearity of expectation

# Hashing



(Src: wikipedia)

# Hashing



**Designing hash functions can be tricky...**

Balls and bins

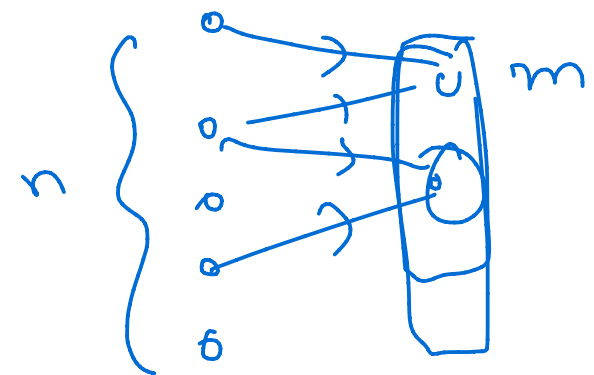
# Some questions

$$B_1, B_2, \dots, B_m$$

**Problem:** suppose we have  $n$  balls and  $m$  bins. Imagine throwing the balls into bins, independently and uniformly at random.

- What is the expected size of each bin?

$$n/m$$



- Suppose  $n = m$ ; What is the expected number of bins with exactly 4 balls?

$$c \cdot n$$

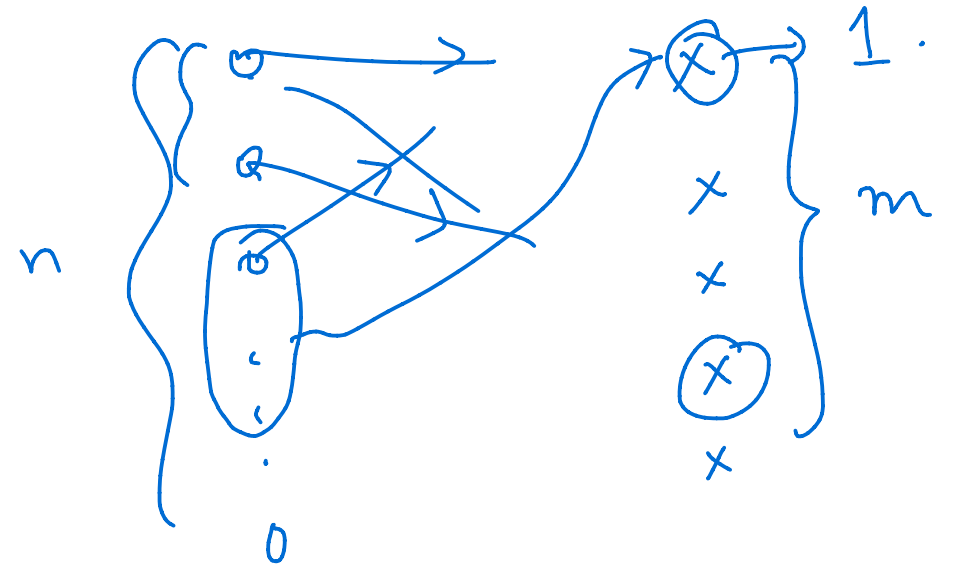
- Suppose  $n = m$ ; What is the probability that there exists a bin with  $(\log n)$  balls?

$$\sqrt{n}$$

# Expected size of ~~a~~ bin 1.

Definition of expectation:

$B$ : random variable which is  
the # balls landing in bin 1.



$$\mathbb{E}[B] = \sum_{k=1}^n k \cdot \Pr[B=k].$$

# Expected size of a bin

Definition of expectation:

**Moral from last week:** never compute expectations using the definition!

# Linearity of expectation

→ Let  $X, Y$  be any two random variables (~~if~~ they could be dependent). Then

$$\underline{\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y].}$$

$$\mathbb{E}[XY] \stackrel{?}{=} \mathbb{E}[X] \cdot \mathbb{E}[Y] ?$$

Main "use": We can often write random variables as sum of "simpler" r.v.s.



# Expected size of a bin

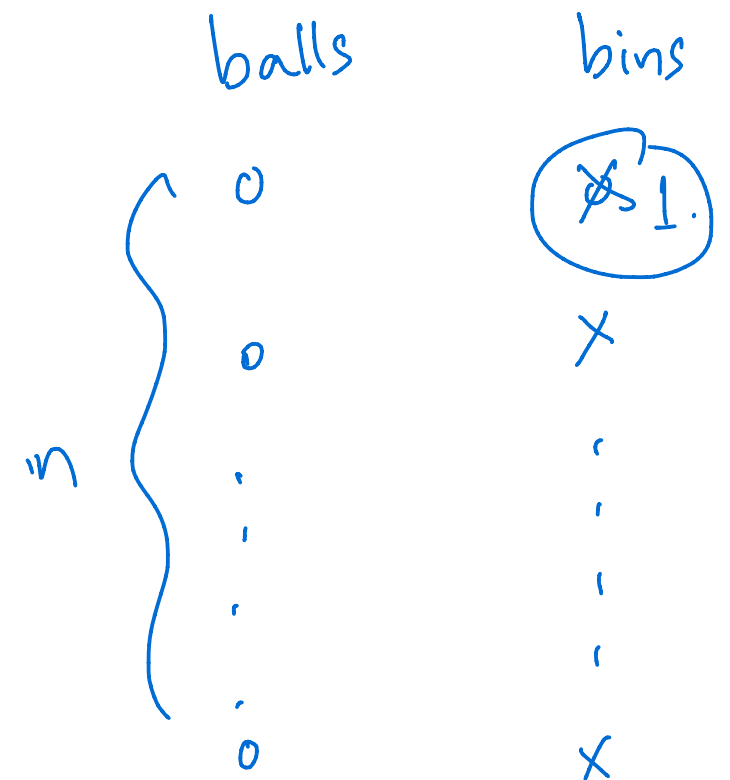
“Decomposing” into a sum of random variables...

$B_i \equiv$  random variable that is  
the # of balls in bin  $i$

$$B_i = n - \sum_{j \neq i} B_j$$

$$\mathbb{E}[B_i] = n - \sum_{j \neq i} \mathbb{E}[B_j]$$

$$\sum_j \mathbb{E}[B_j] = n \quad \leadsto \quad \text{using symmetry, all } \mathbb{E}[B_j] \text{ vals are equal; } \Rightarrow \mathbb{E}[B_j] = \frac{n}{m}.$$



- Take bin 1 ; let  $B$  be the r.v. denoting # of balls in it.

- define  $X_1 = \begin{cases} 1 & \text{if ball 1 went to bin 1} \\ 0 & \text{o/wise.} \end{cases}$

$$X_2 = \begin{cases} 1 & \text{if ball 2 } \dots \dots 1 \\ 0 & \text{o/wise.} \end{cases}$$

by defn,  $B = X_1 + X_2 + \dots + X_n \longrightarrow \star$ .

What is  $\mathbb{E}[X_i]$ ? =  $0 \cdot \Pr[X_i=0] + 1 \cdot \underbrace{\Pr[X_i=1]} = \frac{1}{n}$ .

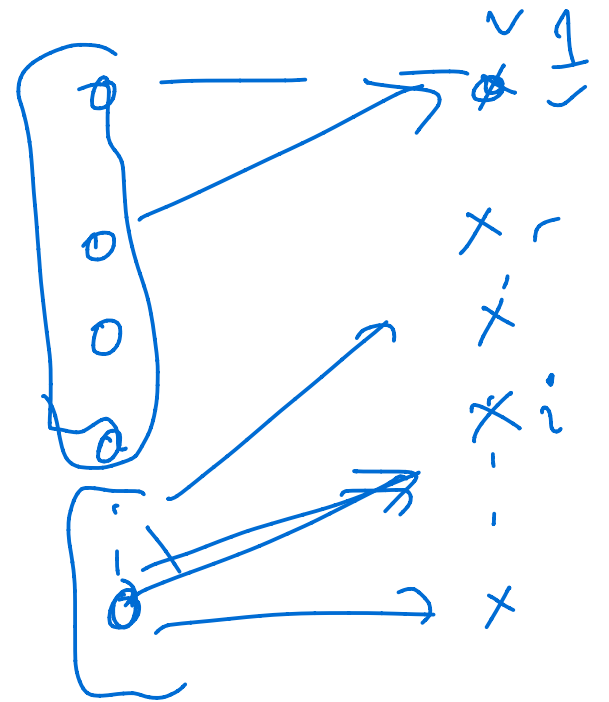
Using linearity of exp. at  ~~$\star$~~ , we get  $\mathbb{E}[B] = \frac{n}{n}$ .

$$(n=m).$$

# Expected # of bins with 4 balls

$Y$ : # of bins with precisely 4 balls

$$X_i = \begin{cases} 1 & \text{if bin } i \text{ receives precisely 4 balls} \\ 0 & \text{otherwise} \end{cases}$$



by definition,  $Y = X_1 + X_2 + \dots + X_n \Rightarrow E[Y] = \sum_{i=1}^n E[X_i]$ .

$$\text{What is } E[X_i]? = \Pr[X_i = 1] = \binom{n}{4} \cdot \frac{1}{n^4} \cdot \left(\frac{n-1}{n}\right)^{n-4}$$

$$E[Y] \approx \frac{n}{80} = \frac{n(n-1)(n-2)(n-3)}{4! \cdot n^4} \cdot \frac{(n-1)^{n-4}}{n^{n-4}} \approx \frac{1}{e} \cdot \frac{1}{24} \approx \frac{1}{80}.$$

Expected # of bins with 4 balls

# Some questions

**Problem:** suppose we have  $n$  balls and  $m$  bins. Imagine throwing the balls into bins, independently and uniformly at random.

- What is the expected size of each bin?
- Suppose  $n = m$ ; What is the expected number of bins with exactly 4 balls?
- Suppose  $n = m$ ; What is the probability that there exists a bin with  $(\log n)$  balls?

# The union bound

# The union bound