# Hash tables & probability

Ben Langmead

**Department of Computer Science**

# Hash Table

*"Hashing with chaining"* or *"chain hashing"*

# Hash Function

$U$

$N$

$n$

Assume accessing table slot is $O(1)$
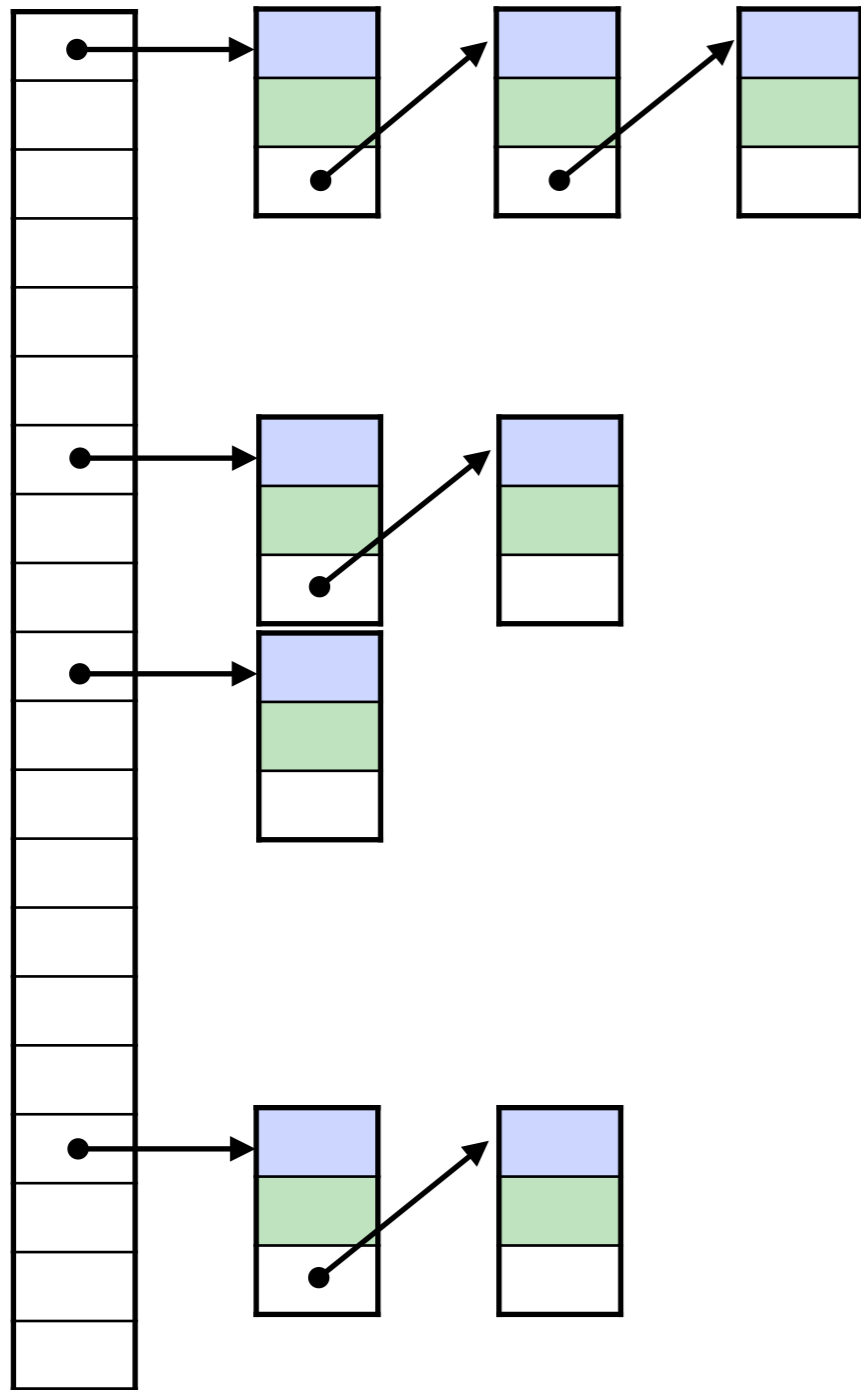
Assume hash function operates on any item from $U$ (integers, strings, etc) and is $O(1)$ time

# Hash Table



What "abstract data types" can we implement with this?

| map | set | counter |
|---|---|---|
| $\langle k_1, v_1\rangle$ | $\langle k_1\rangle$ | $\langle k_1, 7\rangle$ |
| $\langle k_2, v_2\rangle$ | $\langle k_2\rangle$ | $\langle k_2, 4\rangle$ |
| $\langle k_3, v_3\rangle$ | $\langle k_3\rangle$ | $\langle k_3, 8\rangle$ |
| $\langle k_4, v_4\rangle$ | $\langle k_4\rangle$ | $\langle k_4, 5\rangle$ |

# Hash Table

I add $m$ items to an $n$-bucket hash table

***Without*** probability, what can I say?

| Question | Assumption | Statement | Comment |
|---|---|---|---|
| **Does any bucket have more the one item?** | $m > n$ | Yes | Pigeonhole principle |
| **Is any bucket empty?** | $m < n$ | Yes | "Empty pigeonhole" principle |
| **What is the average bucket occupancy?** | - | $m/n$ | - |

Nothing profound here

# Hash Table

I have added $m$ items to a $n$-bucket hash table.  What "interesting questions" can I ask about the table's state?
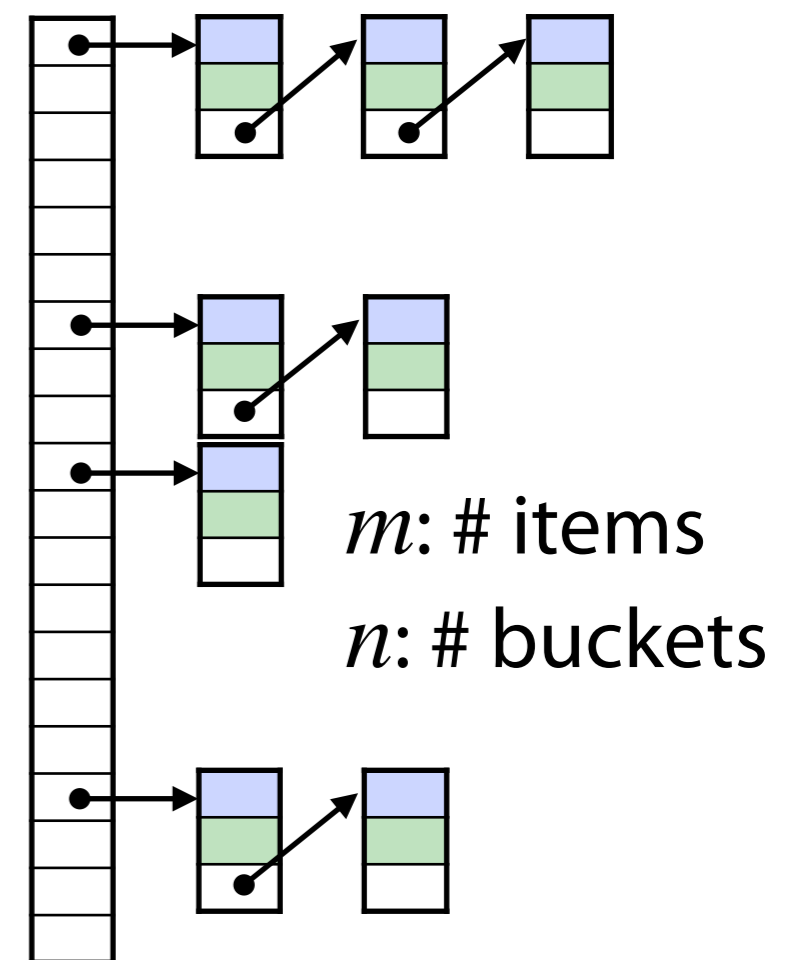
How many buckets are empty?

How many items are in the *median* bucket?

How many items are in the *average* bucket?

What's the chance all buckets are non-empty?
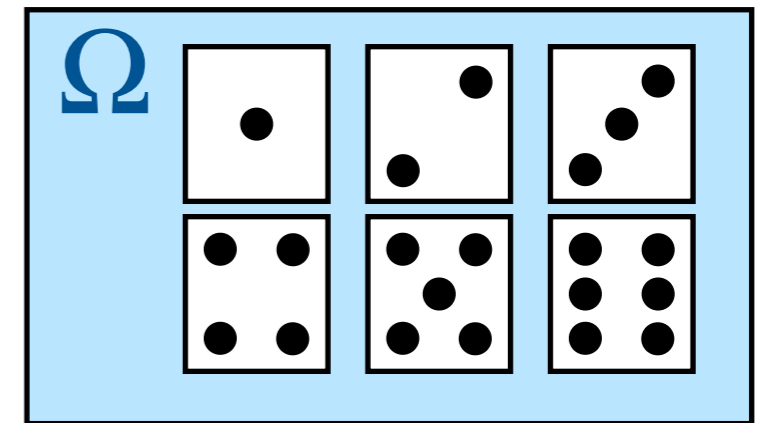
How many items are in the *fullest* bucket?

What's the chance no bucket has >1 item?

$m$: # items

$n$: # buckets

# Probability

*Sample space* ($\Omega$) is **set** of all possible outcomes

    E.g. $\Omega$ = { all possible rolls of 2 dice }



An *event* is a **subset** of $\Omega$

    $A$ = { rolls where $1^{st}$ die is odd }
    $B$ = { rolls where $2^{nd}$ die is even }

*When outcomes are equally likely*, can use "naive definition of probability"

$\Pr(A)$: fraction of outcomes that are in $A$

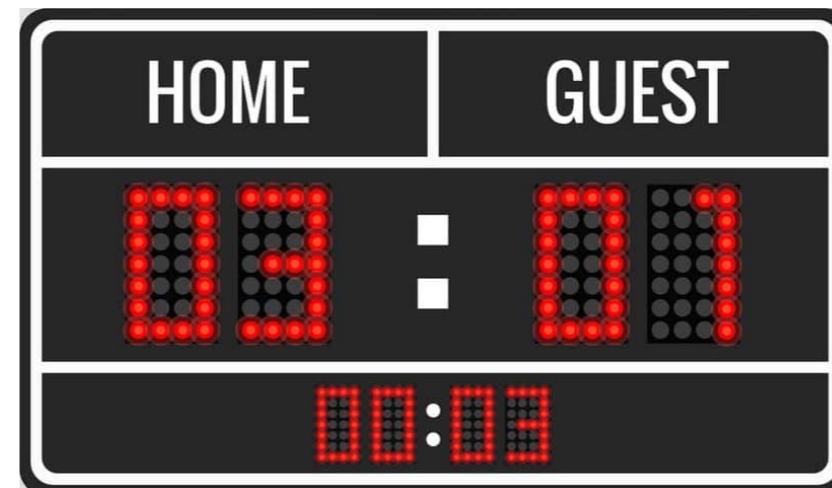$$\Pr(A) = |A| / |\Omega| = 18/36 = 0.5$$

# Probability

"Naive definition" of probability fails to apply when outcomes are not equally probable

Loaded coin

# goals scored in soccer game

## Probability function $\mathrm{Pr}$

$\mathrm{Pr} : \mathscr{P}(\Omega) \to \mathbb{R}$ , where $\mathscr{P}(\Omega)$ is "power set" (set of all subsets) of $\Omega$, satisfies conditions:

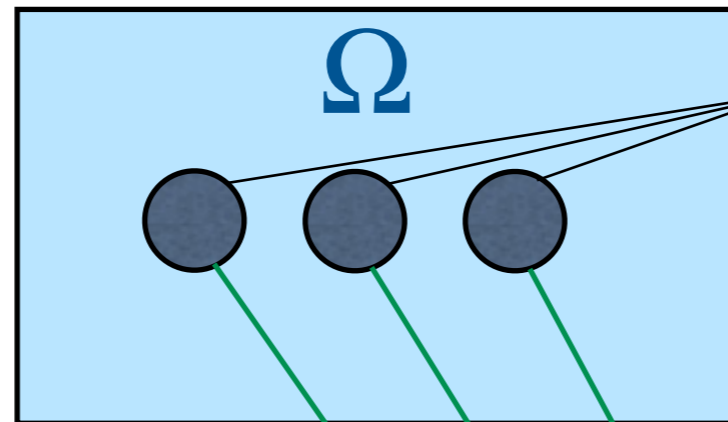1. For any event $E$, $\ 0 \le \mathrm{Pr}(E) \le 1$

2. $\mathrm{Pr}(\Omega) = 1$

3. Probabilities of disjoint events $E_1, E_2, \ldots$ add:

$$\mathrm{Pr}\left(\bigcup_{i \ge 1} E_i\right) = \sum_{i \ge 1} \mathrm{Pr}(E_i)$$
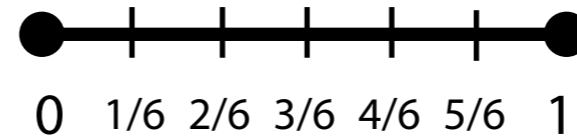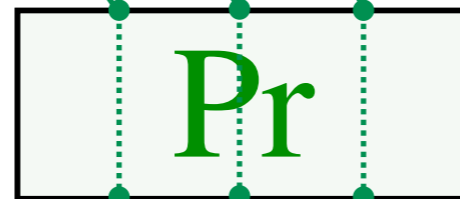
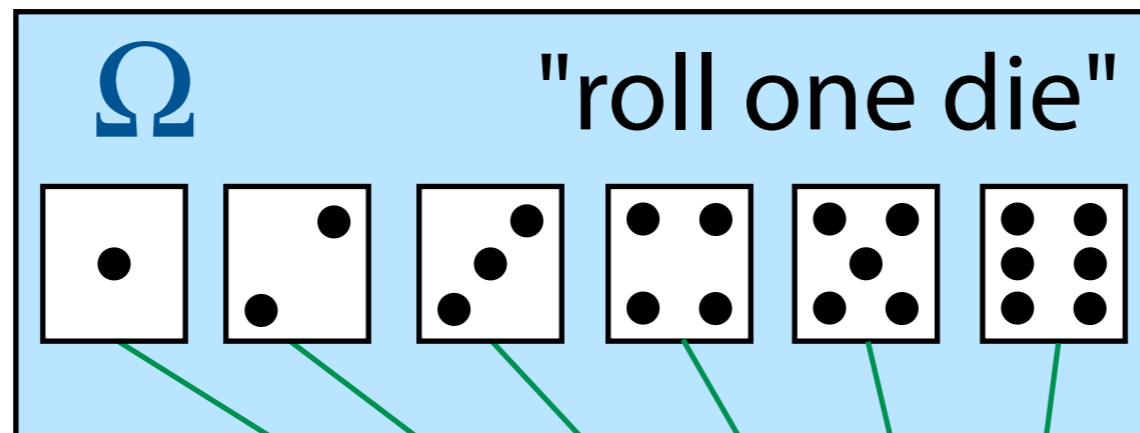Probability function Pr
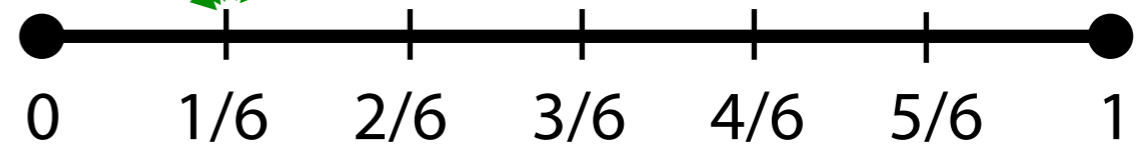
Sample space

Ω

outcomes

set

Pr

probability function

reals in $[0, 1]$

0   1/6  2/6  3/6  4/6  5/6   1

Probability function Pr

$\Omega$   "roll one die"

set

probability function

Probabilities of disjoint events add;

$\mathrm{Pr}(\{\boxed{\cdot}, \boxed{\vcenter{\hbox{∴}}}\}) = 1/3$

Pr

0   1/6   2/6   3/6   4/6   5/6   1

# Random variable

Random variables have two "natures"

$$X$$

**Function**, mapping *outcomes* from $\Omega$ to *numbers* (in $\mathbb{R}$)

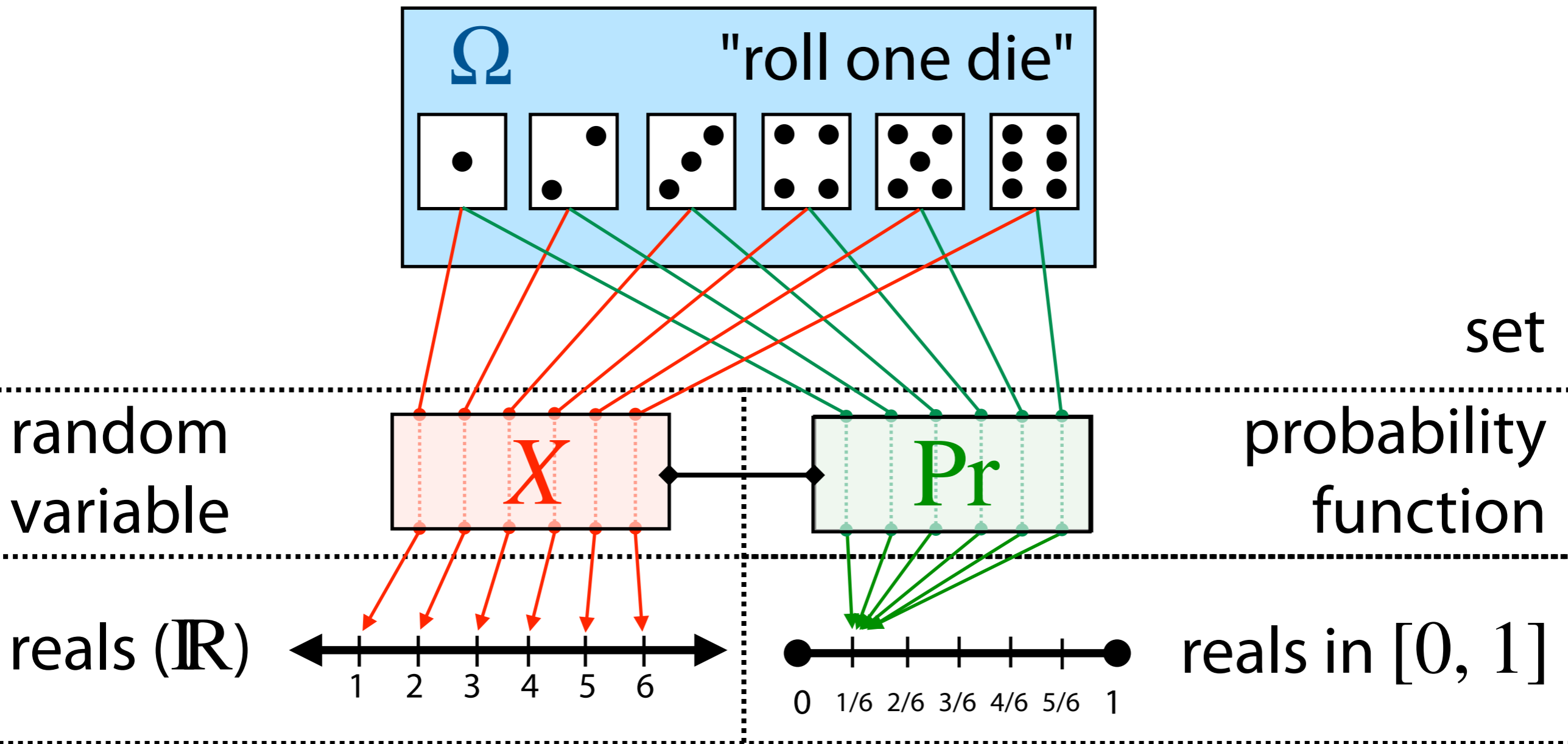$$X(\boxed{\,\vcenter{\hbox{:·:}}\,}) = 4$$

$$Y = 3.5 - X$$

**Potential experiment** with a *distribution* (a $\Pr$ for its $\Omega$) and numerical result

# Random variable

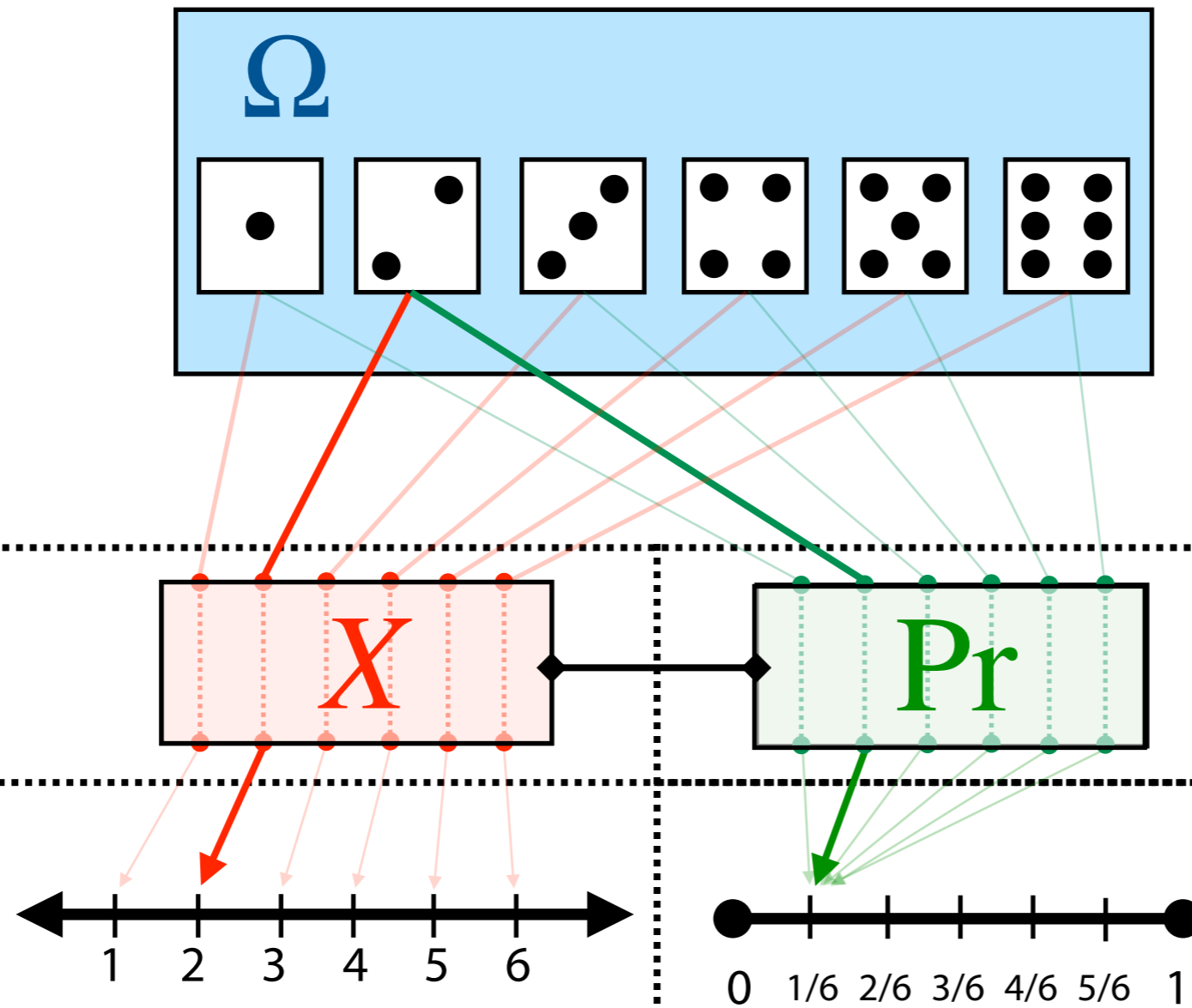We use capitals e.g. $X$, $Y$ to denote a random variable

Abbreviate with "r.v."

# Random variable & probability function

$\Omega$ — "roll one die"

set

random variable — $X$

probability function — Pr

reals ($\mathbb{R}$)

1  2  3  4  5  6

reals in [0, 1]

0  1/6  2/6  3/6  4/6  5/6  1

$$\text{Pr}(X = 2) = \text{Pr}(\boxdot) = 1/6$$

# Random variable & probability function



$$\mathrm{Pr}(X = 2) = \mathrm{Pr}(\boxed{\cdot\,\cdot}) = 1/6$$

# Random variable & probability function



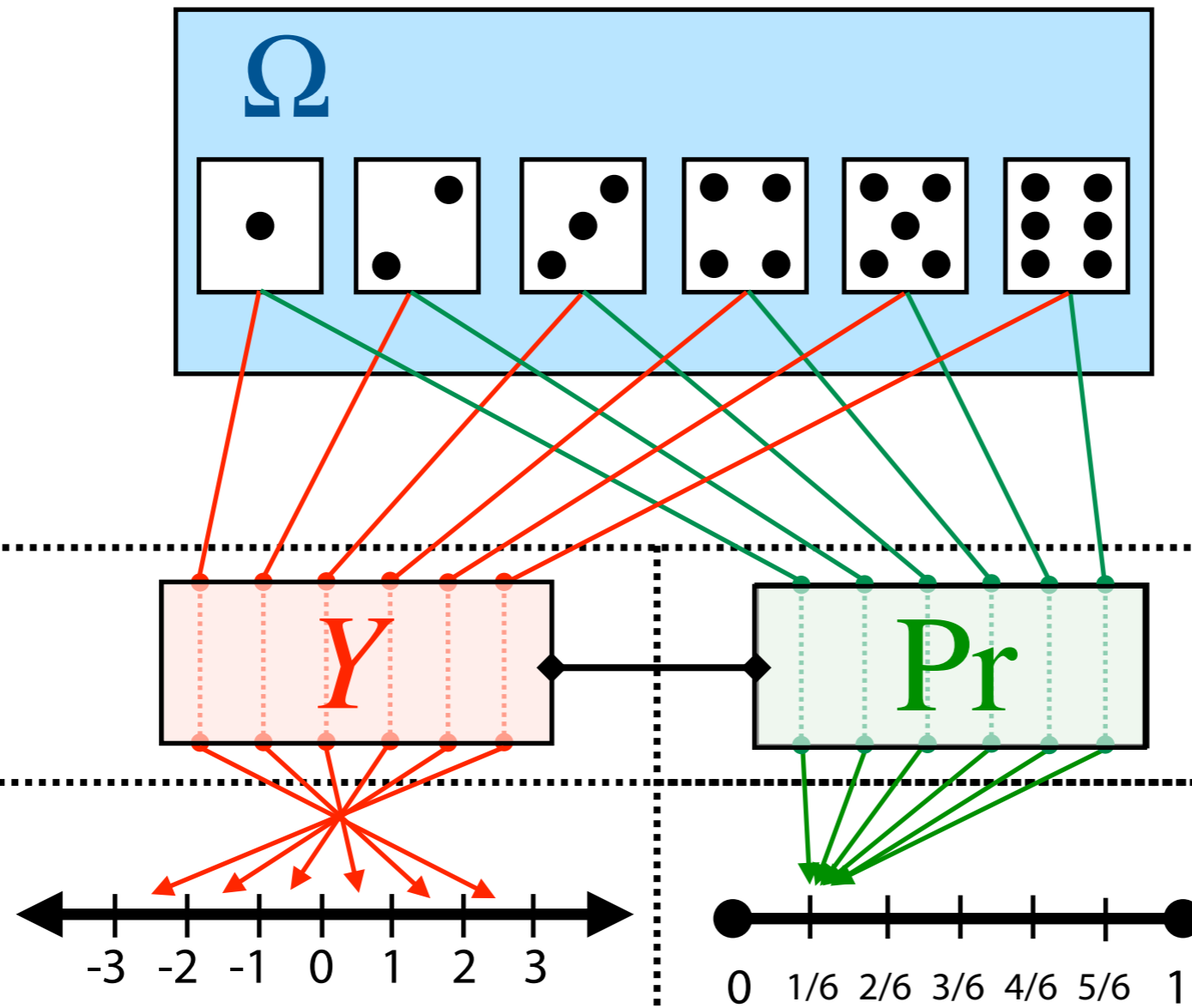$$\Pr(X \geq 4) = \Pr(\boxdot) + \Pr(\boxdot) + \Pr(\boxdot) = 1/2$$

# Random variable & probability function



$$Y = 3.5 - X$$

# Random variable & probability function



$$Y = 3.5 - X$$

# Expected value

Expectation ("expected value") of a discrete r.v. $X$, called $\mathbf{E}[X]$, is given by

$$\mathbf{E}[X] = \sum_x x \cdot \Pr(X = x)$$
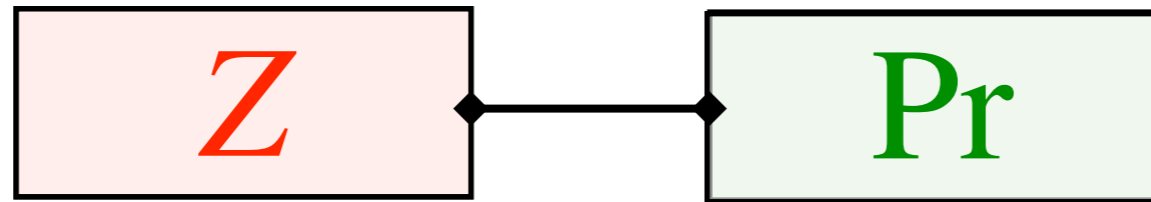
where summation is over values in range of $X$.

# Linearity of expectation

For discrete r.v.s $X_1, X_2, \ldots, X_n$

$$\mathbf{E}\left[\sum_{i=1}^{n} X_i\right] = \sum_{i=1}^{n} \mathbf{E}[X_i]$$

True whether or not $X_i$s are independent

# Expected value

$$Z \quad \quad Pr$$

$$Z = X + Y$$ where $X$ is fair die roll & $Y$ is fair coin flip

When $Z$ is a linear combination of other r.v.s, $\mathbf{E}[Z]$ can be easier to get than $Pr$

$\mathbf{E}[Z] = \mathbf{E}[X] + \mathbf{E}[Y]$ is simple $(3.5 + 0.5 = 4)$
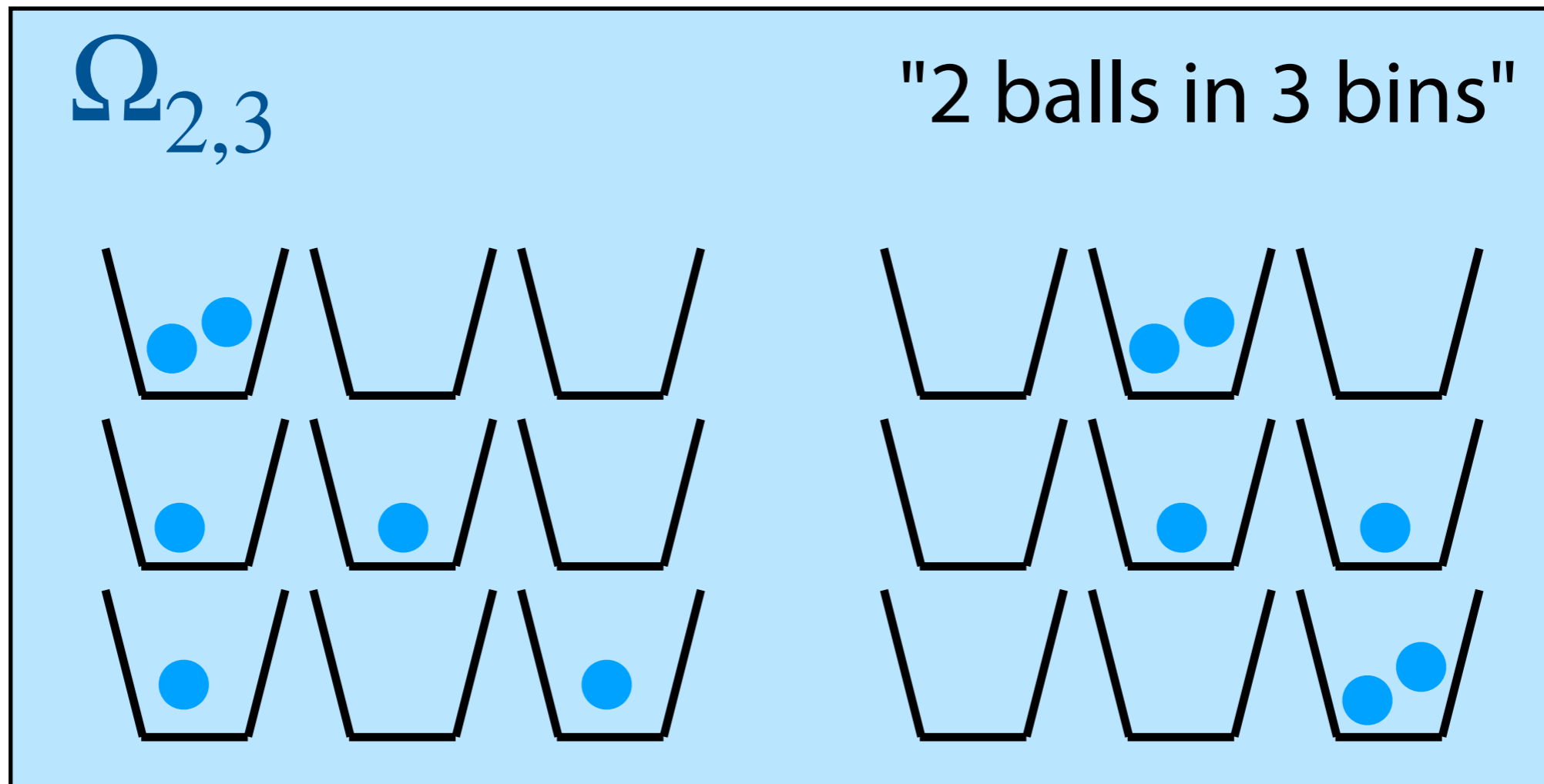
# Hash Table

I have added $m$ items to a $n$-bucket hash table

Besides this setup, what else do we need to define a random variable describing the table?

1. Sample space $\Omega$ ——— Possible allocation of items to buckets

2. Probability func. $\mathrm{Pr}$

3. Map $X$ from outcomes to reals

Depend on question asked, assumptions made about hash function

# Balls & bins

Throw $m$ balls into $n$ bins uniformly and independently

# Hash Table

I have added $m$ items to a $n$-bucket hash table. What "interesting questions" can I ask about the table's state?

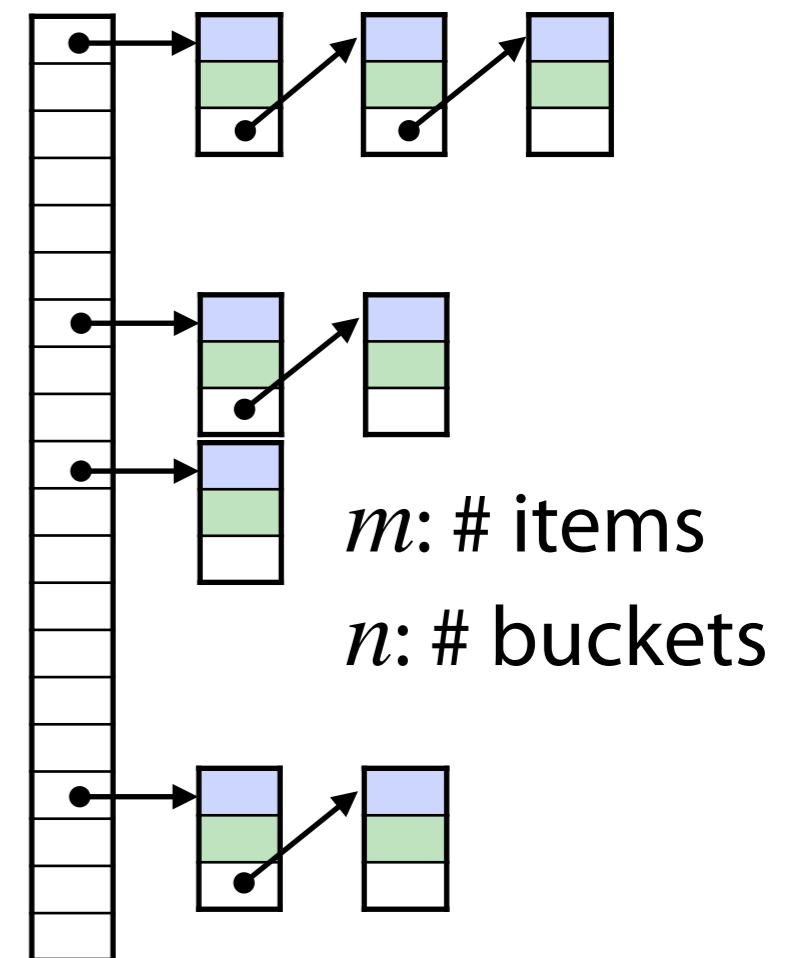How many buckets are empty?

How many items are in the *median* bucket?

How many items are in the *average* bucket?

What's the chance all buckets are non-empty?

How many items are in the *fullest* bucket?

What's the chance no bucket has >1 item?

$m$: # items

$n$: # buckets

# Balls & bins

I throw $m$ balls into $n$ bins uniformly and independently. What can I ask about the bins and their contents?

How many bins are empty?
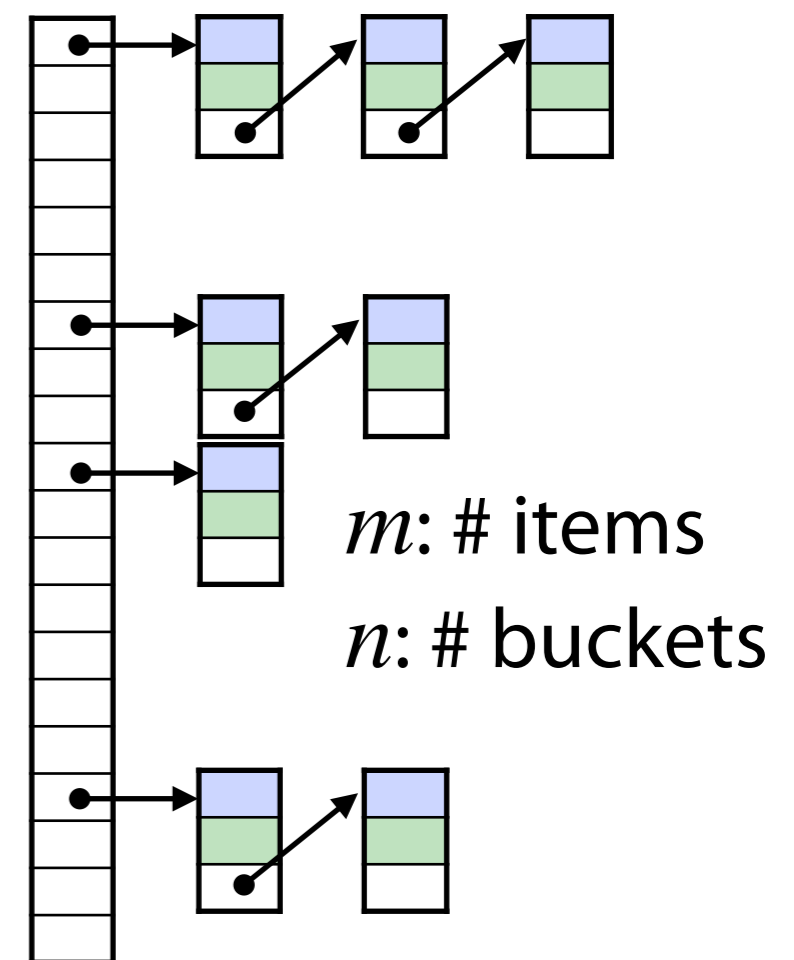
How many balls are in the *average* bin?

How many balls are in the *fullest* bin?

How many balls are in the *median* bin?

What's the chance all bins are non-empty?

What's the chance no bin has >1 item?

$m$: # items

$n$: # buckets

# Balls & bins

I throw $m$ balls into $n$ bins uniformly and independently. What can I ask?

| Category | Questions | | Approach |
|---|---|---|---|
| Empty/ non empty | How many buckets are empty? | What's the chance all buckets are non-empty? | **Coupon collector** |
| Collisions / no collisions | How many throws until there is a >0.5 chance of a collision? | What is the chance no bin has >1 item? | **Birthday problem** |
| Local (single bin) occupancy | What's the occupancy of a given bucket? | What is the chance a given bucket has >2 items? | **Binomial, Geometric, Poisson r.v.s** |
| Global occupancy | What is the *median* bucket occupancy? | What is the *maximum* bucket occupancy? | **Often hard** |