

# Probability in Computing

## LECTURE 5: MORE APPLICATIONS WITH PROBABILISTIC ANALYSIS, BINS AND BALLS

# Agenda

- ◆ Review: Coupon Collector's problem and Packet Sampling
- ◆ Analysis of Quick-Sort
- ◆ Birthday Paradox and applications
- ◆ The Bins and Balls Model

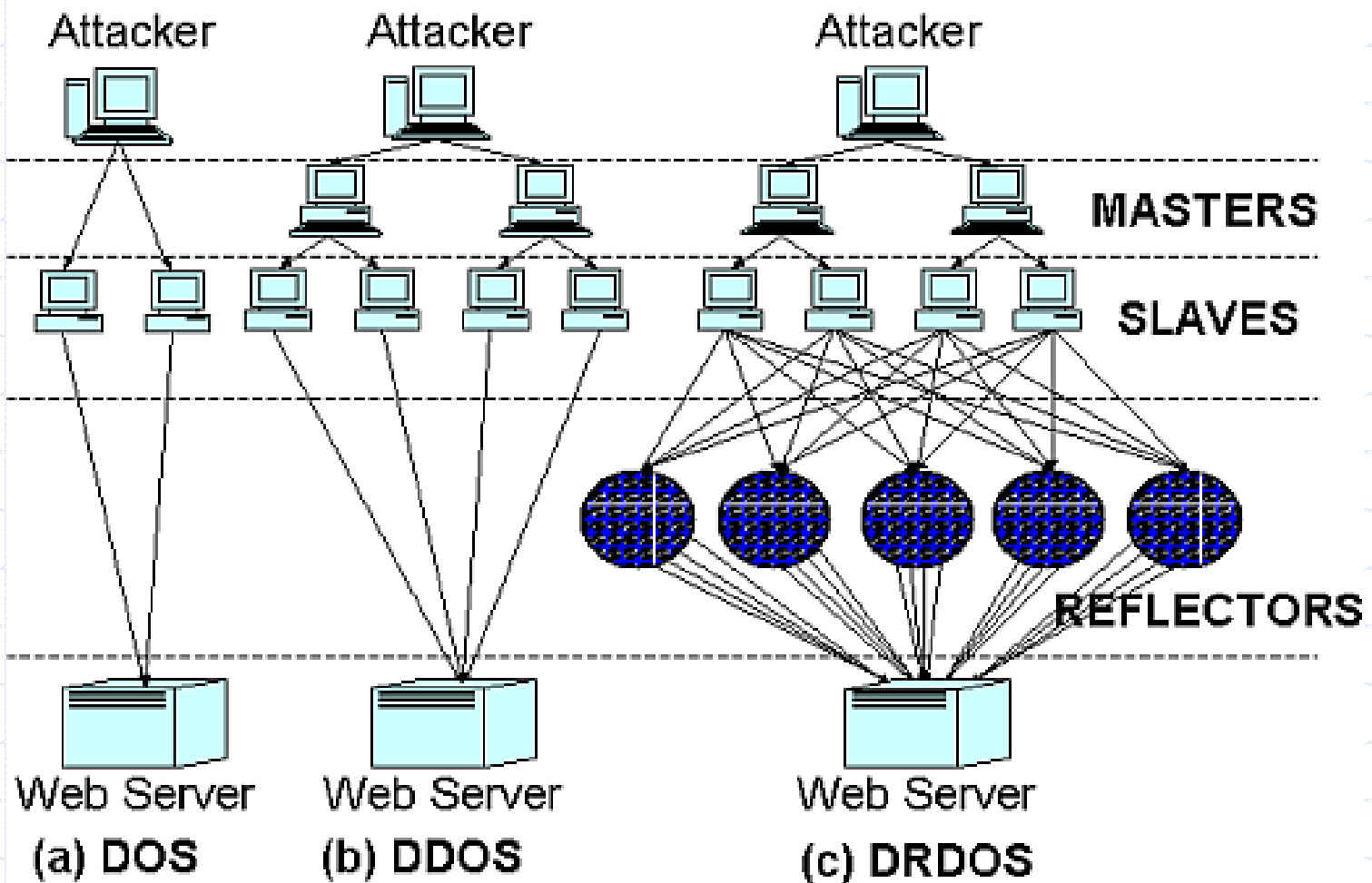
# Coupon Collector Problem

- ◆ Problem: Suppose that each box of cereal contains one of  $n$  different coupons. Once you obtain one of every type of coupon, you can send in for a prize.
- ◆ Question: How many boxes of cereal must you buy before obtaining at least one of every type of coupon.
- ◆ Let  $X$  be the number of boxes bought until at least one of every type of coupon is obtained.
- ◆  $E[X] = nH(n) = n \ln n$

# Application: Packet Sampling

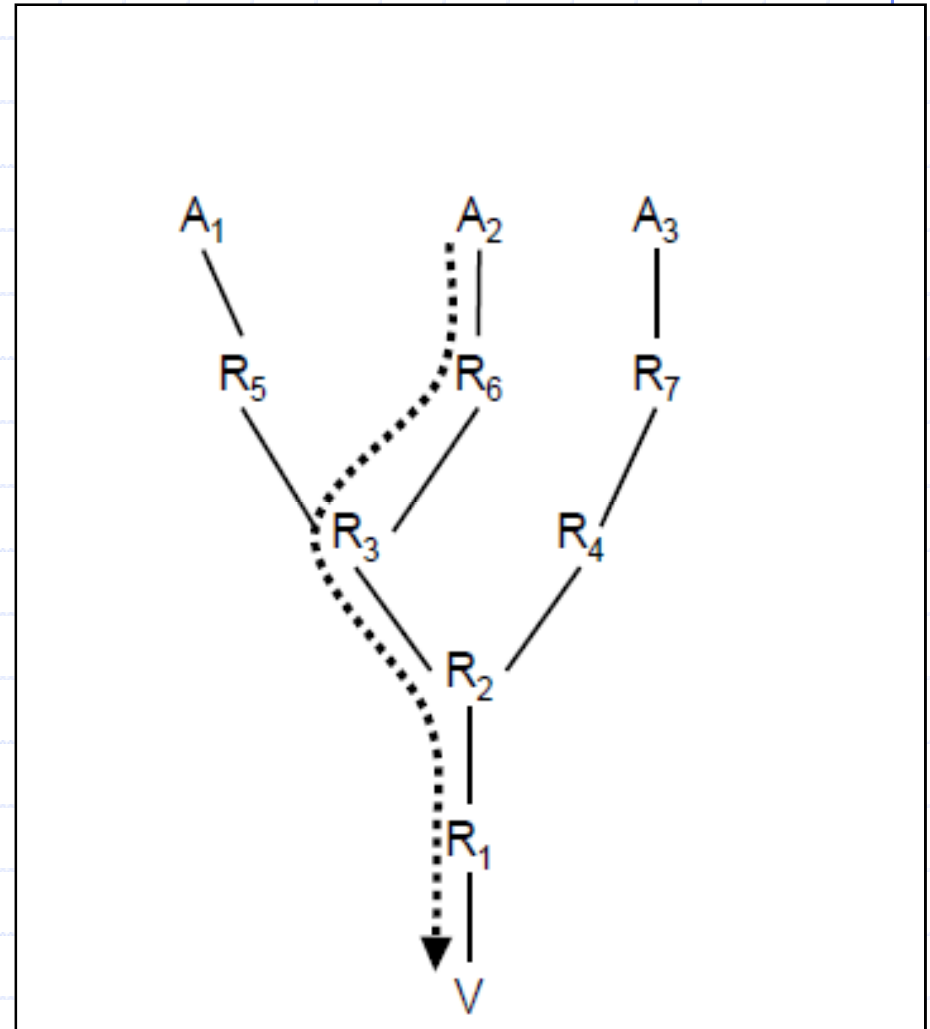
- ◆ Sampling packets on a router with probability  $p$ 
  - The number of packets transmitted after the last sampled packet until and including the next sampled packet is geometrically distributed.
- ◆ From the point of destination host, determining all the routers on the path is like a coupon collector's problem.
- ◆ If there's  $n$  routers, then the expected number of packets arrived before destination host knows all of the routers on the path  $= n \ln(n)$ .

# DoS attack

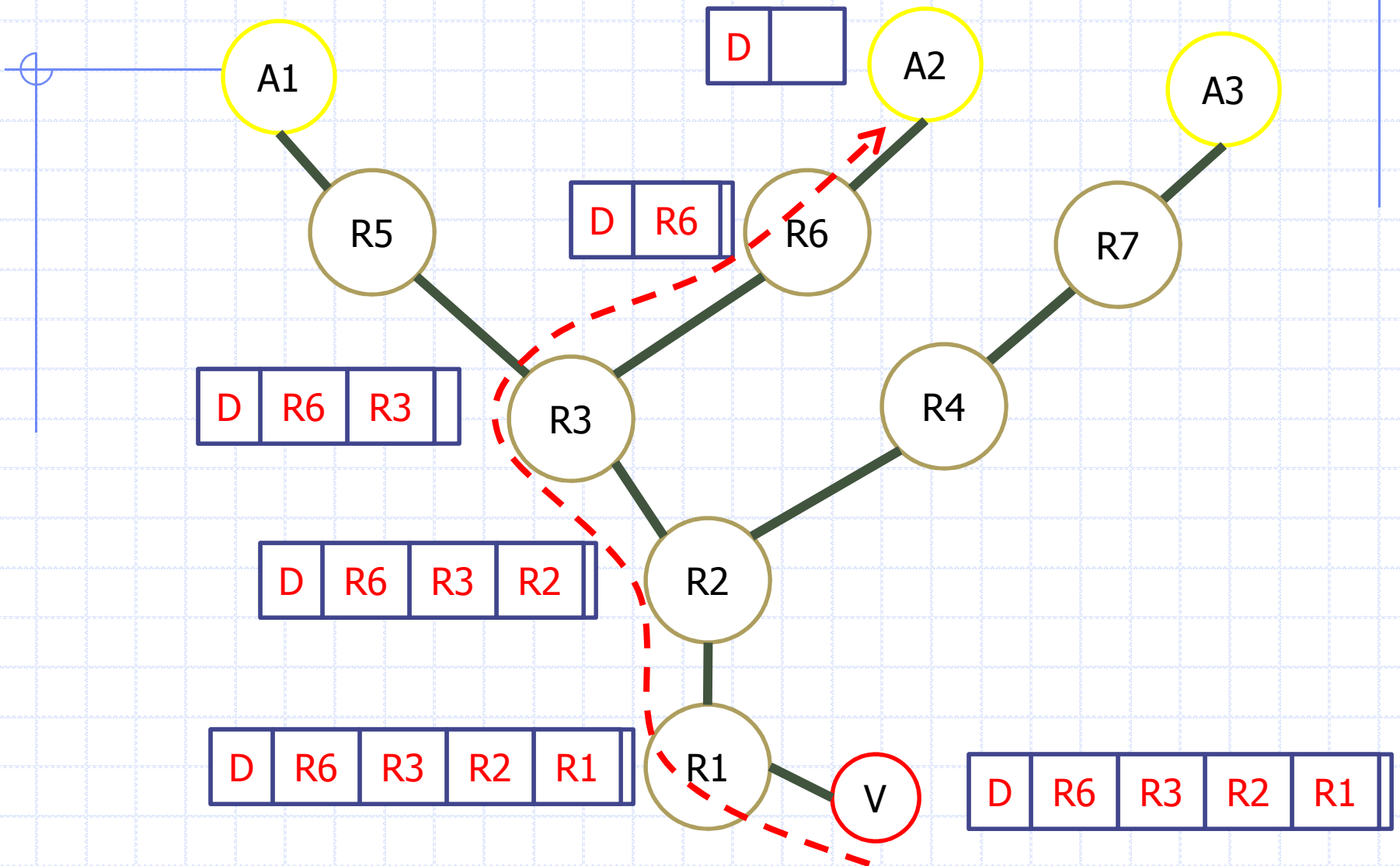


# IP traceback

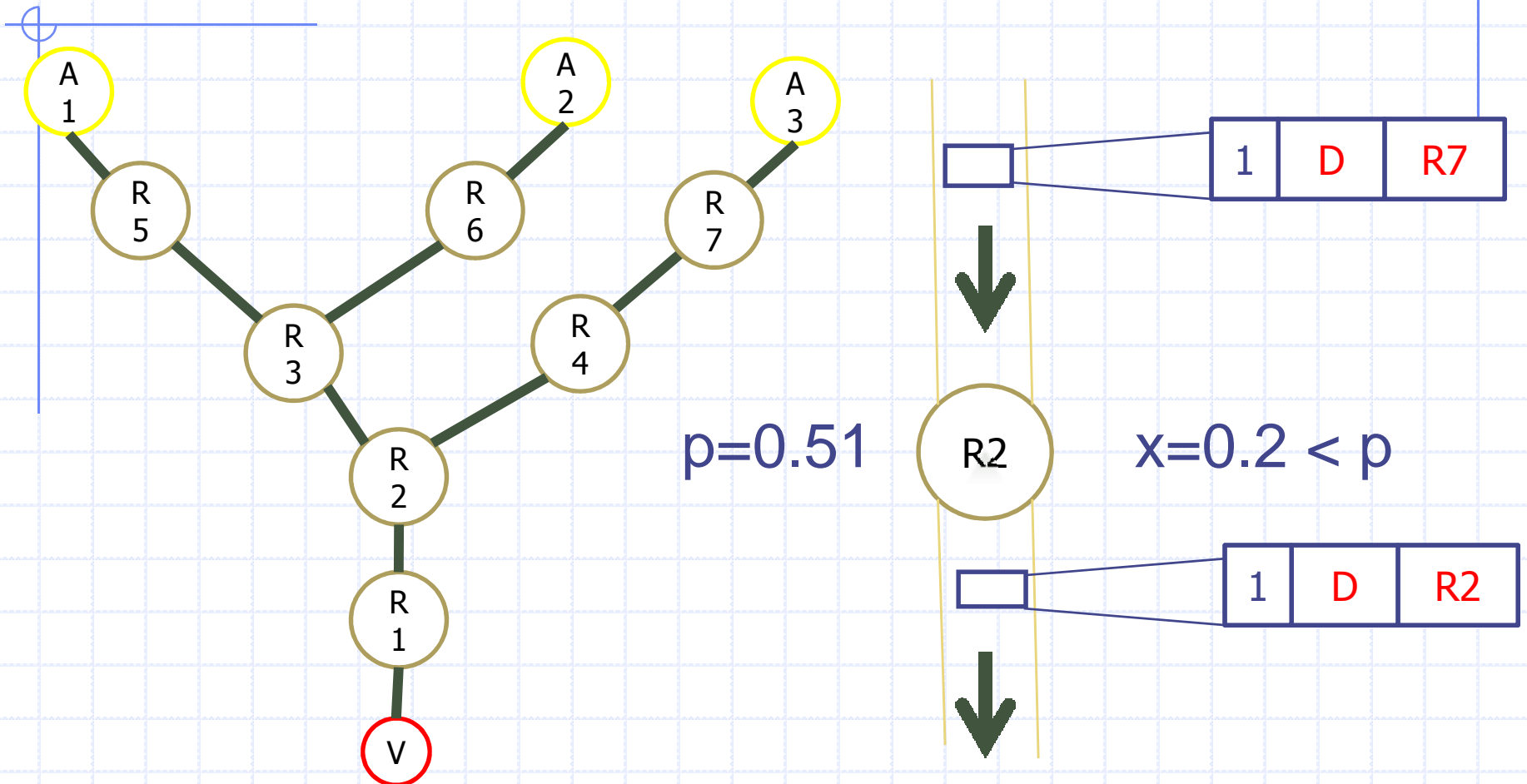
- ◆ Marking and Reconstruction
  - Node append vs. node sampling



# Node append



# Node Sampling





# Expected Run-Time of QuickSort

## Quicksort Algorithm:

**Input:** A list  $S = \{x_1, \dots, x_n\}$  of  $n$  distinct elements over a totally ordered universe.

**Output:** The elements of  $S$  in sorted order.

1. If  $S$  has one or zero elements, return  $S$ . Otherwise continue.
2. Choose an element of  $S$  as a pivot; call it  $x$ .
3. Compare every other element of  $S$  to  $x$  in order to divide the other elements into two sublists:
  - (a)  $S_1$  has all the elements of  $S$  that are less than  $x$ ;
  - (b)  $S_2$  has all those that are greater than  $x$ .
4. Use Quicksort to sort  $S_1$  and  $S_2$ .
5. Return the list  $S_1, x, S_2$ .

# Analysis

- ◆ Worst-case:  $n^2$ .
- ◆ Depends on how we choose the pivot.
- ◆ Good pivot (divide the list in two nearly equal length sub-lists) vs. Bad pivot.
- ◆ In case of good pivot  $\rightarrow n \lg(n)$ . [by solving recurrence]
- ◆ If we choose pivot point randomly, we will have a randomized version of QuickSort.

# Analysis

◆  $X_{ij}$  be a random variable that

- Takes value 1 if  $y_i$  and  $y_j$  are compared with each other
- 0 if they are not compared.

◆  $E[X] = \sum \sum E[X_{ij}]$

◆  $E[X_{ij}] = 2 / (j - i + 1)$  (when we choose either  $i$  or  $j$  from the set of  $Y_{ij}$  pivots  $\{y_i, y_{i+1}, \dots, y_j\}$ )

◆ Using  $k = j - i + 1$ , we can compute  $E[X] = 2n \ln(n)$

# Detail analysis

$$\begin{aligned}\mathbf{E}[X] &= \sum_{i=1}^{n-1} \sum_{j=i+1}^n \frac{2}{j-i+1} \\&= \sum_{i=1}^{n-1} \sum_{k=2}^{n-i+1} \frac{2}{k} \\&= \sum_{k=2}^n \sum_{i=1}^{n+1-k} \frac{2}{k} \\&= \sum_{k=2}^n (n+1-k) \frac{2}{k} \\&= \left( (n+1) \sum_{k=2}^n \frac{2}{k} \right) - 2(n-1) \\&= (2n+2) \sum_{k=1}^n \frac{1}{k} - 4n.\end{aligned}$$

# Birthday “Paradox”

What is the probability that two persons in a room of 30 have the same birthday?

# Birthday Paradox

◆ Ways to assign  $k$  different birthdays without duplicates:

$$\begin{aligned} N &= 365 * 364 * \dots * (365 - k + 1) \\ &= 365! / (365 - k)! \end{aligned}$$

◆ Ways to assign  $k$  different birthdays with possible duplicates:

$$D = 365 * 365 * \dots * 365 = 365^k$$

# Birthday “Paradox”

Assuming real birthdays assigned randomly:

$N/D$  = probability there are no duplicates

$1 - N/D$  = probability there is a duplicate

$$= 1 - 365! / ((365 - k)!(365)^k)$$

# Generalizing Birthdays

$$P(n, k) = 1 - n! / (n-k)! n^k$$

Given  $k$  random selections from  $n$  possible values,  $P(n, k)$  gives the probability that there is at least 1 duplicate.



# Birthday Probabilities

$$P(\text{no two match}) = 1 - P(\text{all are different})$$

$$P(2 \text{ chosen from } N \text{ are different})$$

$$= 1 - 1/N$$

$$P(3 \text{ are all different})$$

$$= (1 - 1/N)(1 - 2/N)$$

$$P(n \text{ trials are all different})$$

$$= (1 - 1/N)(1 - 2/N) \dots (1 - (n - 1)/N)$$

$$\ln(P)$$

$$= \ln(1 - 1/N) + \ln(1 - 2/N) + \dots \ln(1 - (k - 1)/N)$$

# Happy Birthday Bob!

$$\ln(P) = \ln(1 - 1/N) + \dots + \ln(1 - (k-1)/N)$$

$$\text{For } 0 < x < 1: \ln(1 - x) \leq -x$$

$$\ln(P) \leq -(1/N + 2/N + \dots + (k-1)/N)$$

Gauss says:

$$1 + 2 + 3 + 4 + \dots + (k-1) + k = \frac{1}{2} k(k+1)$$

So,

$$\ln(P) \leq -\frac{1}{2} (k-1) k / N$$

$$P \leq e^{-\frac{1}{2} (k-1) k / N}$$

$$\text{Probability of match} \geq 1 - e^{-\frac{1}{2} (k-1) k / N}$$

# Applying Birthdays

$$P(n, k) > 1 - e^{-k*(k-1)/2n}$$

◆ For  $n = 365$ ,  $k = 20$ :

$$P(365, 20) > 1 - e^{-20*(19)/2*365}$$

$$P(365, 20) > .4058$$

◆ For  $n = 2^{64}$ ,  $k = 2^{32}$ :  $P(2^{64}, 2^{32}) > .39$

◆ For  $n = 2^{64}$ ,  $k = 2^{33}$ :  $P(2^{64}, 2^{33}) > .86$

◆ For  $n = 2^{64}$ ,  $k = 2^{34}$ :  $P(2^{64}, 2^{34}) > .9996$

◆ Application: Digital Signatures

# Balls into Bins

- ◆ We have  $m$  balls that are thrown into  $n$  bins, with the location of each ball chosen independently and uniformly at random from  $n$  possibilities.
- ◆ What does the distribution of the balls into the bins look like
  - “Birthday paradox” question: is there a bin with at least 2 balls
  - How many of the bins are empty?
  - How many balls are in the fullest bin?

Answers to these questions give solutions to many problems in the design and analysis of algorithms

# The maximum load

- ◆ When  $n$  balls are thrown independently and uniformly at random into  $n$  bins, the probability that the maximum load is more than  $3 \ln n / \ln \ln n$  is at most  $1/n$  for  $n$  sufficiently large.

- By Union bound,  $\Pr [\text{bin 1 receives } \geq M \text{ balls}] \leq \binom{n}{M} \left(\frac{1}{n}\right)^M$ .
- Note that:

$$\binom{n}{M} \left(\frac{1}{n}\right)^M \leq \frac{1}{M!} \leq \left(\frac{e}{M}\right)^M.$$

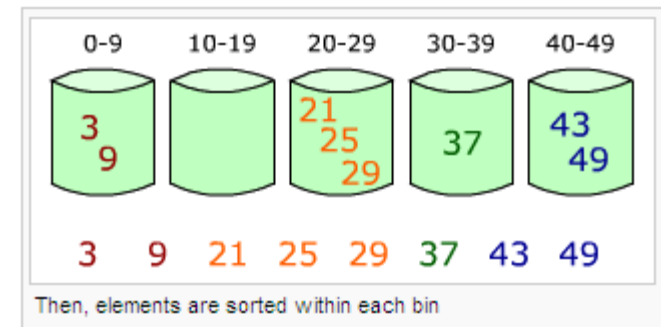
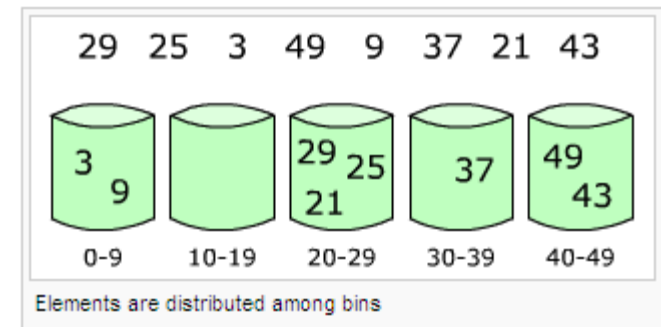
- Now, using Union bound again,  $\Pr [\text{any bin receives } \geq M \text{ balls}]$  is at most

$$n \left(\frac{e}{M}\right)^M \leq n \left(\frac{e \ln \ln n}{3 \ln n}\right)^{3 \ln n / \ln \ln n}$$

which is  $\leq 1/n$

# Application: Bucket Sort

- ◆ A sorting algorithm that breaks the  $\Omega(n \log n)$  lower bound under certain input assumption
- ◆ Bucket sort works as follows:
  - Set up an array of initially empty "buckets."
  - Scatter: Go over the original array, putting each object in its bucket.
  - Sort each non-empty bucket.
  - Gather: Visit the buckets in order and put all elements back into the original array.



- ◆ A set of  $n = 2^m$  integers, randomly chosen from  $[0, 2^k)$ ,  $k \geq m$ , can be sorted in expected time  $O(n)$

- Why: will analyze later!

# The Poisson Distribution

## ◆ Consider $m$ balls, $n$ bins

- $\Pr[\text{a given bin is empty}] = \left(1 - \frac{1}{n}\right)^m \approx e^{-m/n}$ ;
- Let  $X_j$  is a indicator r.v. that os 1 if bin  $j$  empty, 0 otherwise
- Let  $X$  be a r.v. that represents # empty bins

$$\mathbf{E}[X] = \mathbf{E}\left[\sum_{i=1}^n X_i\right] = \sum_{i=1}^n \mathbf{E}[X_i] = n\left(1 - \frac{1}{n}\right)^m \approx ne^{-m/n}$$

- Generalizing this argument,  $\Pr[\text{a given bin has } r \text{ balls}] =$

$$\binom{m}{r} \left(\frac{1}{n}\right)^r \left(1 - \frac{1}{n}\right)^{m-r} = \frac{1}{r!} \frac{m(m-1) \cdots (m-r+1)}{n^r} \left(1 - \frac{1}{n}\right)^{m-r}$$

- Approximately,  $p_r \approx \frac{e^{-m/n} (m/n)^r}{r!}$

- So: **Definition 5.1:** A discrete Poisson random variable  $X$  with parameter  $\mu$  is given by the following probability distribution on  $j = 0, 1, 2, \dots$ :

$$\Pr(X = j) = \frac{e^{-\mu} \mu^j}{j!}.$$

# Limit of the Binomial Distribution

We have shown that, when throwing  $m$  balls randomly into  $b$  bins, the probability  $p_r$  that a bin has  $r$  balls is approximately the Poisson distribution with mean  $m/b$ . In general, the Poisson distribution is the limit distribution of the binomial distribution with parameters  $n$  and  $p$ , when  $n$  is large and  $p$  is small. More precisely, we have the following limit result.

**Theorem 5.5:** *Let  $X_n$  be a binomial random variable with parameters  $n$  and  $p$ , where  $p$  is a function of  $n$  and  $\lim_{n \rightarrow \infty} np = \lambda$  is a constant that is independent of  $n$ . Then, for any fixed  $k$ ,*

$$\lim_{n \rightarrow \infty} \Pr(X_n = k) = \frac{e^{-\lambda} \lambda^k}{k!}.$$

This theorem directly applies to the balls-and-bins scenario. Consider the situation where there are  $m$  balls and  $b$  bins, where  $m$  is a function of  $b$  and  $\lim_{n \rightarrow \infty} m/b = \lambda$ . Let  $X_n$  be the number of balls in a specific bin. Then  $X_n$  is a binomial random variable with parameters  $m$  and  $1/b$ . Theorem 5.5 thus applies and says that

$$\lim_{n \rightarrow \infty} \Pr(X_n = r) = \frac{e^{-m/n} (m/n)^r}{r!},$$