



# Introduction to Parallel & Distributed Computing

## Discrete Search & Termination Detection

Lecture 15, Spring 2014

Instructor: 罗国杰

[gluo@pku.edu.cn](mailto:gluo@pku.edu.cn)

# *In the Last Lecture ...*

---

- u **Parallel DFS**

- **distribute/request work in a stack**

- u **Parallel best-first search**

- **distribute/request work in a priority queue**

# *In the Last Lecture ...*

---

## u Parallel DFS

- **Work-splitting strategies**
  - **Send nodes near the bottom of the stack**
    - u suitable if the search space is uniform
  - **Send nodes near the cutoff depth**
    - u if a strong heuristic is available (to order successors so that goal nodes move to the left of the state-space tree)
  - **Send half the nodes between the bottom of the stack and the cutoff depth**
    - u suitable whenever the search space is uniform or highly irregular

# *In the Last Lecture ...*

---

## u Parallel DFS

### ■ Load-balancing schemes

- Asynchronous round robin
  - u work requests are generated independently
  - u 2+ processors may request work from the same donor processor
- Global round robin
  - u The *target* is store is a shared variable (shared address space) or at a designated processor (message passing)
  - u Work requests are evenly distributed, but there is contention for access to target
- Random polling
  - u Work requests are evenly distributed

# **Distributed Termination Detection**

- u Suppose we only want to print one solution
- u We want all processes to halt as soon as one process finds a solution
- u This means processes must periodically check for messages
  - Every process calls MPI\_Iprobe every time search reaches a particular level (such as the cutoff depth)
  - A process sends a message after it has found a solution

# MPI\_Iprobe

---

- u **Nonblocking test for a message**

```
int MPI_Iprobe(  
    int source,  
    int tag,  
    MPI_Comm comm,  
    int *flag,  
    MPI_Status *status);
```

- u **flag: true if a message with the specified source, tag, and communicator is available**
- u **It is not necessary to receive a message immediately after it has been probed for, and the same message may be probed for several times before it is received.**

# ***Simple (Incorrect) Algorithm***

- u A process halts after one of the following events has happened:
  - It has found a solution and sent a message to all of the other processes
  - It has received a message from another process
  - ~~It has completely searched its portion of the state space tree~~

# *Why Algorithm Fails*

---

- u If a process calls MPI\_Finalize before another active process attempts to send it a message, we get a run-time error
  - u How this could happen?
    - A process finds a solution after another process has finished searching its share of the subtrees
- OR
- A process finds a solution after another process has found a solution

# **Distributed Termination Problem**

---

- u **Distributed termination problem:** Ensuring that
  - all processes are inactive AND
  - no messages are en route
- u **A number of algorithms have been proposed.**
  - Token termination detection
    - [Dijkstra, Seijen, Gasteren, 1983]
  - Tree-based termination detection
    - [Rokusawa, Ichiyoshi, Chikayama, Nakashima, 1988]
  - ...

# **Token Termination Detection: A Simplified Scenario**

- u Assume that all processors are organized in a logical ring.
- u A simplified scenario
  - Once a processor goes idle, it never receives more work
  - Assume, for now that message passing can only happen from  $P_i$  to  $P_j$  if  $j > i$ .
- u Algorithm
  - Processor  $P_0$  initiates a token on the ring when it goes idle.
  - Each intermediate processor receives this token and forwards it when it becomes idle.
  - When the token reaches processor  $P_0$ , all processors are done.

# Token Termination Detection

u Assume that all processors are organized in a logical ring.

- Work transfers allowed in the system are from processor  $P_j$  to  $P_i$  such that  $i < j$

u Algorithm

- When  $P_0$  goes idle, it colors itself *white* and initiates a *white* token.
- If  $P_j$  sends work to  $P_i$  and  $j > i$  then processor  $P_j$  becomes *black*.
- If processor  $P_i$  has the token and  $P_i$  is idle, it passes the token to  $P_{i+1}$ .
  - If  $P_i$  is *black*, then the color of the token is set to *black* before it is sent to  $P_{i+1}$ .
  - If  $P_i$  is *white*, the token is passed unchanged.
- After  $P_i$  passes the token to  $P_{i+1}$ ,  $P_i$  becomes *white*.
- After the token passes  $P_0$ , the token becomes *white*.
- Terminates when processor  $P_0$  receives a *white* token and is itself idle.

# ***Token Termination Detection***

---

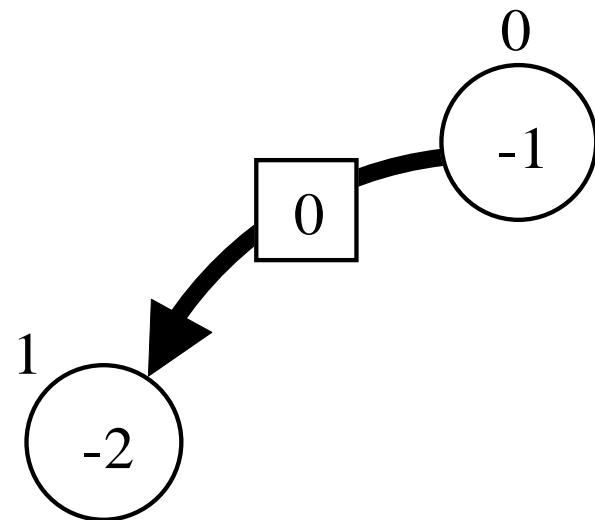
- u There are modified schemes when more information is passed
  - For example, to keep track of #messages in the communication channels
  - For example, to transfer the current best value searched

# **Dijkstra et al.'s Algorithm**

- u Each process has a color and a message count
  - Initial color is white
  - Initial message count is 0
- u A process that sends a message turns black and increments its message count
- u A process that receives a message turns black and decrements its message count
- u If all processes are white and sum of all their message counts are 0, there are no pending messages and we can terminate the processes

# *Dijkstra et al.'s Algorithm (cont.)*

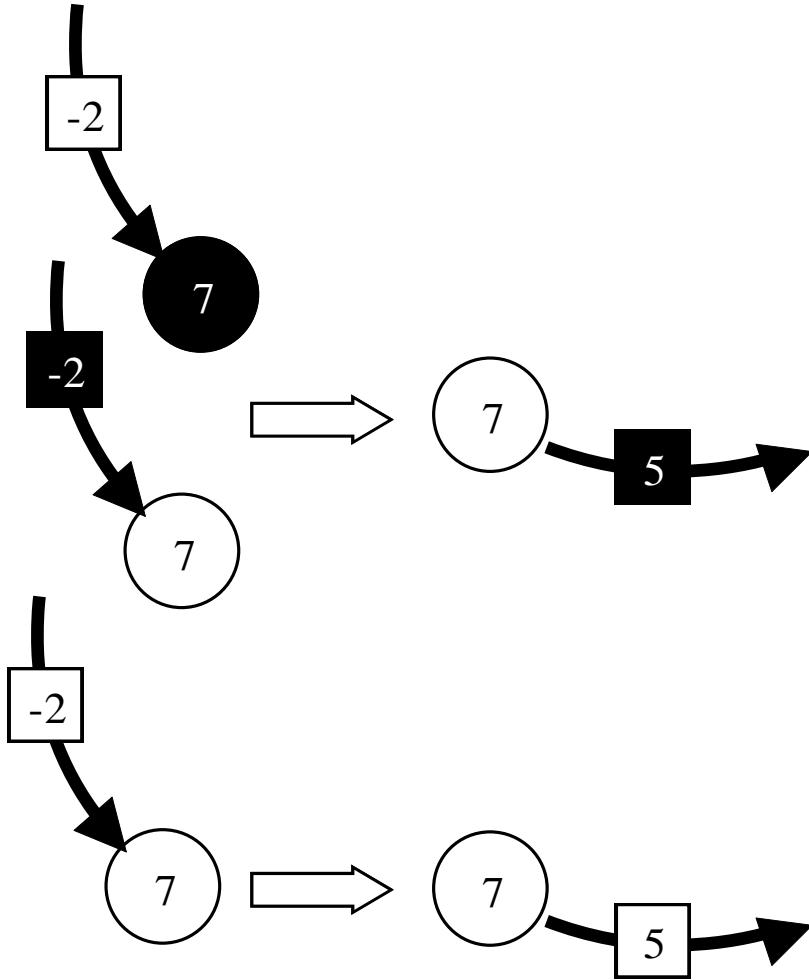
- u Organize processes into a logical ring
- u Process 0 passes a token around the ring
- u Token also has a color (initially white) and count (initially 0)



# *Dijkstra et al.'s Algorithm (cont.)*

- u A process receives the token
  - If process is black
    - Process changes token color to black
    - Process changes its color to white
  - Process adds its message count to token's message count
- u A process sends the token to its successor in the logical ring

# Dijkstra et al.'s Algorithm (cont.)

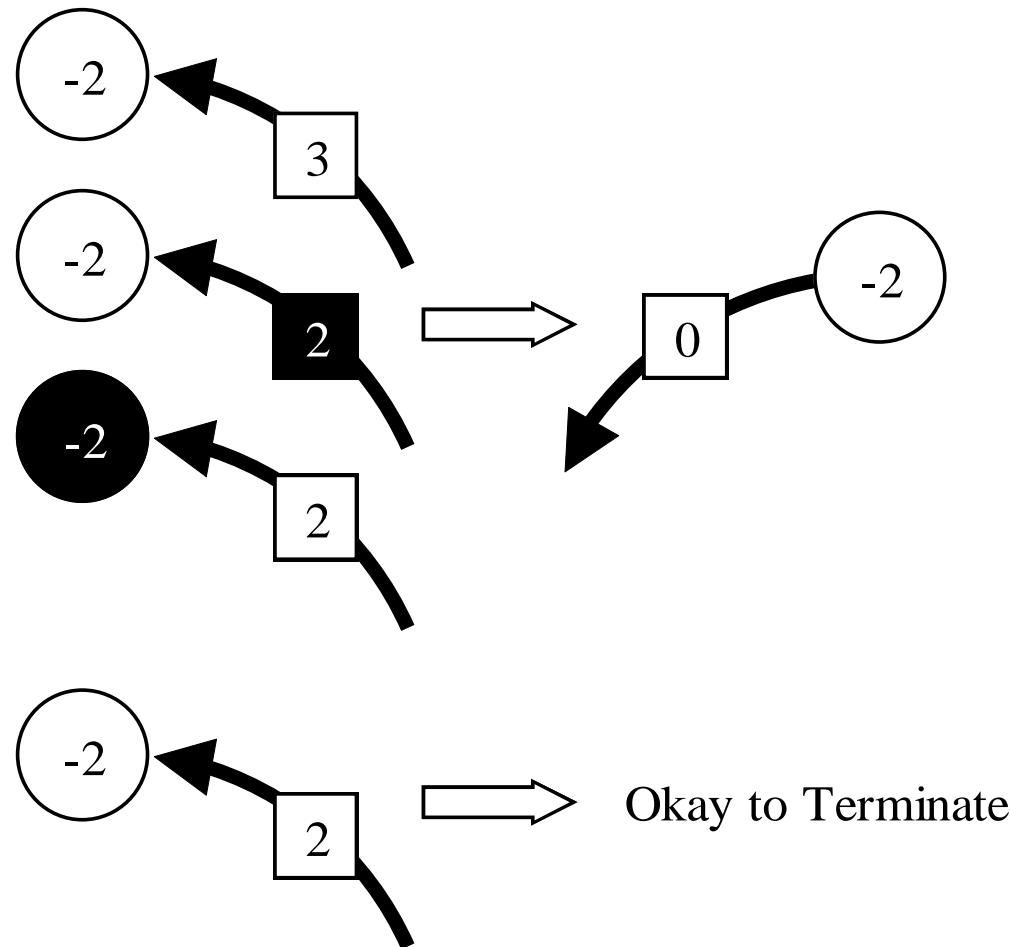


# *Dijkstra et al.'s Algorithm (cont.)*

u Process 0 receives the token

- Safe to terminate processes if
  - Token is white
  - Process 0 is white
  - Token count + process 0 message count = 0
- Otherwise, process 0 must probe ring of processes again

# Dijkstra et al.'s Algorithm (cont.)



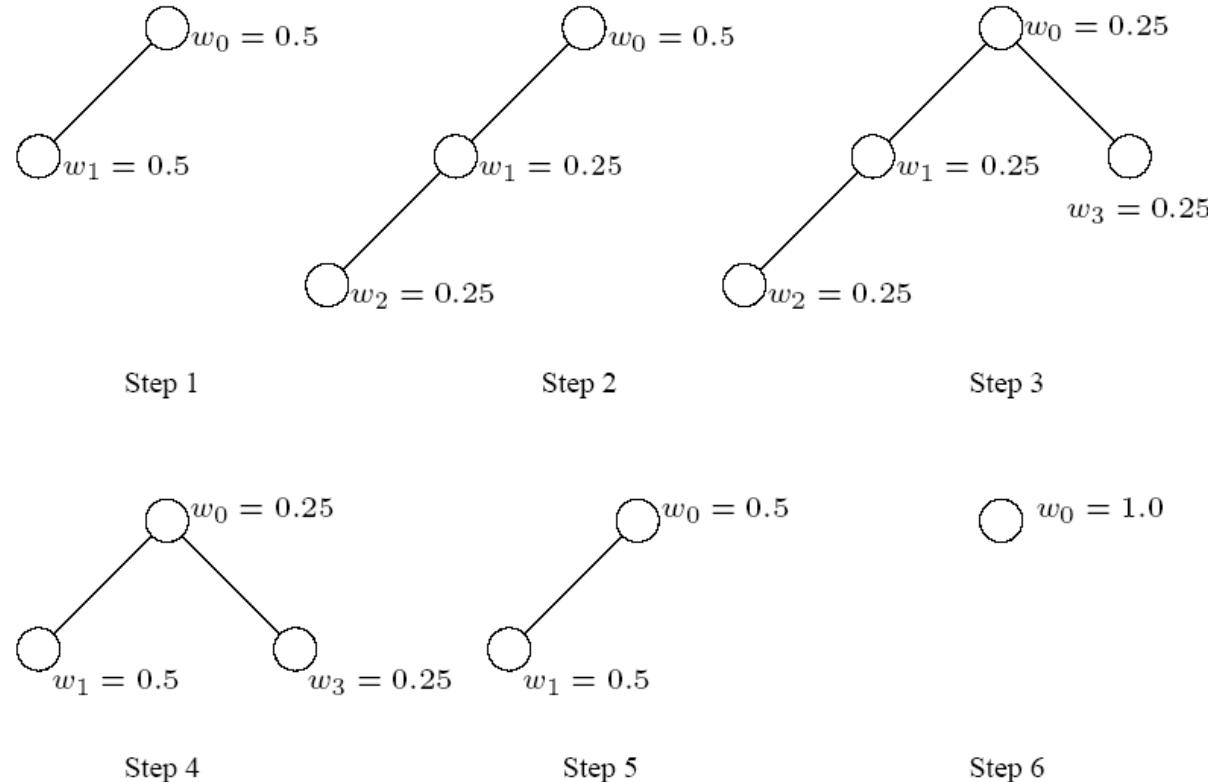
# **Tree-Based Termination Detection**

## **u Algorithm**

- Associate weights with individual workpieces. Initially, processor  $P_0$  has all the work and a weight of one.
- Whenever work is partitioned, the weight is split into half and sent with the work.
- When a processor gets done with its work, it sends its parent the weight back.
- Termination is signaled when the weight at processor  $P_0$  becomes one again.

## **u Note that underflow and finite precision are important factors associated with this scheme.**

# **Tree-Based Termination Detection**



**Steps 1-6 illustrate the weights  
at various processors after each work transfer**



# Introduction to Parallel & Distributed Computing

## Distributed Random Number Generation

Lecture 15, Spring 2014

Instructor: Guojie Luo

([gluo@pku.edu.cn](mailto:gluo@pku.edu.cn))

# **Outline**

---

- u **Monte Carlo method**
- u **Parallel (uniform) random number generators**

# **Monte Carlo Method**

---

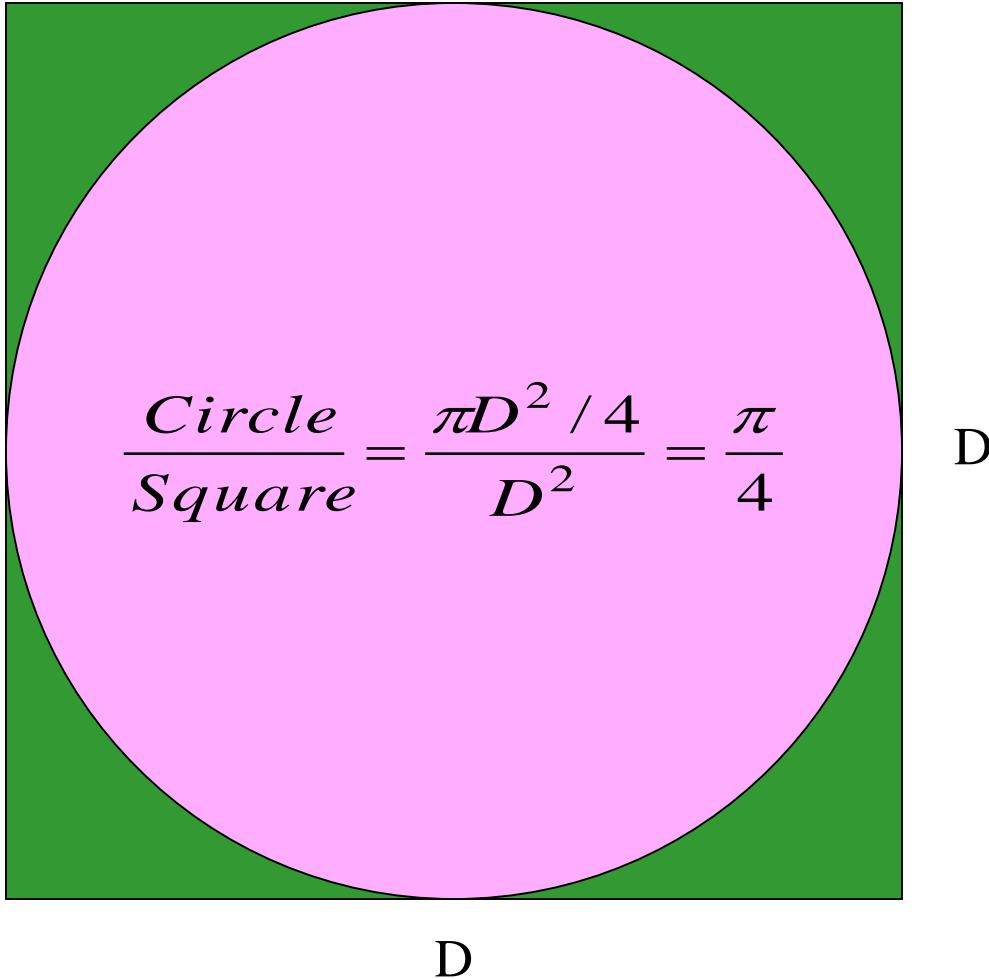
- u **Solve a problem using statistical sampling**
- u **Name comes from Monaco's gambling resort city**
- u **First important use in development of atomic bomb during World War II**

# ***Applications of Monte Carlo Method***

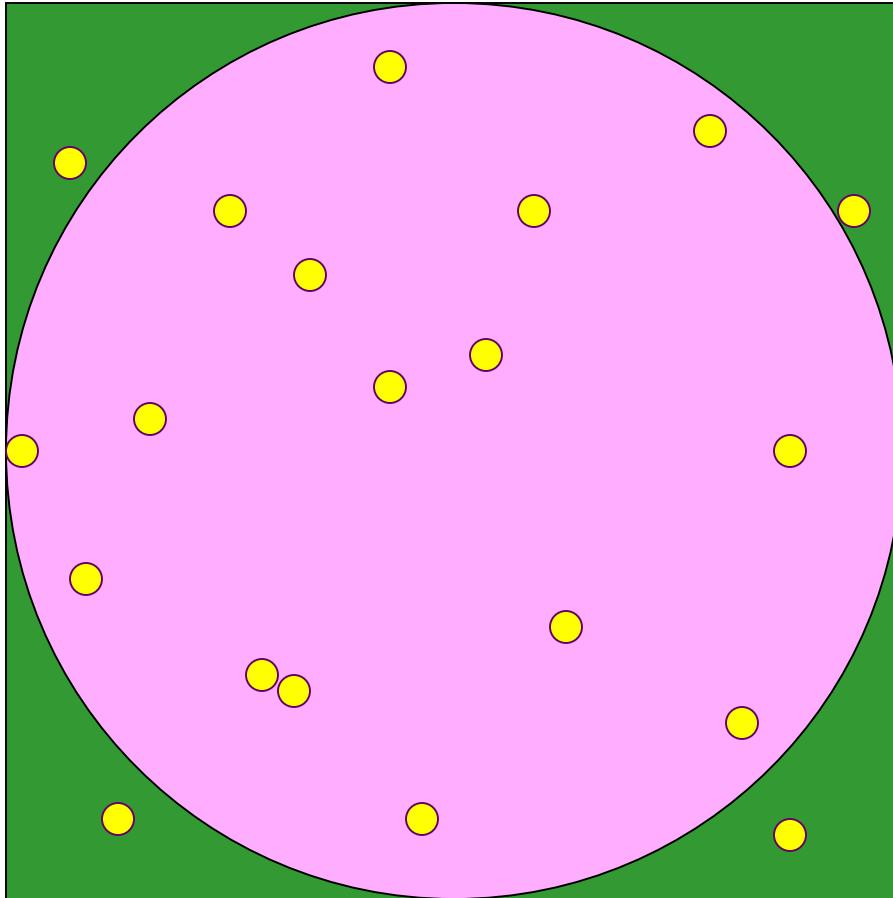
---

- u Evaluating integrals of arbitrary functions of 6+ dimensions
- u Predicting future values of stocks
- u Solving partial differential equations
- u Sharpening satellite images
- u Modeling cell populations
- u Finding approximate solutions to NP-hard problems

# *Example of Monte Carlo Method*



# *Example of Monte Carlo Method*



$$\frac{16}{20} \approx \frac{\pi}{4} \Rightarrow \pi \approx 3.2$$

# **Relative Error**

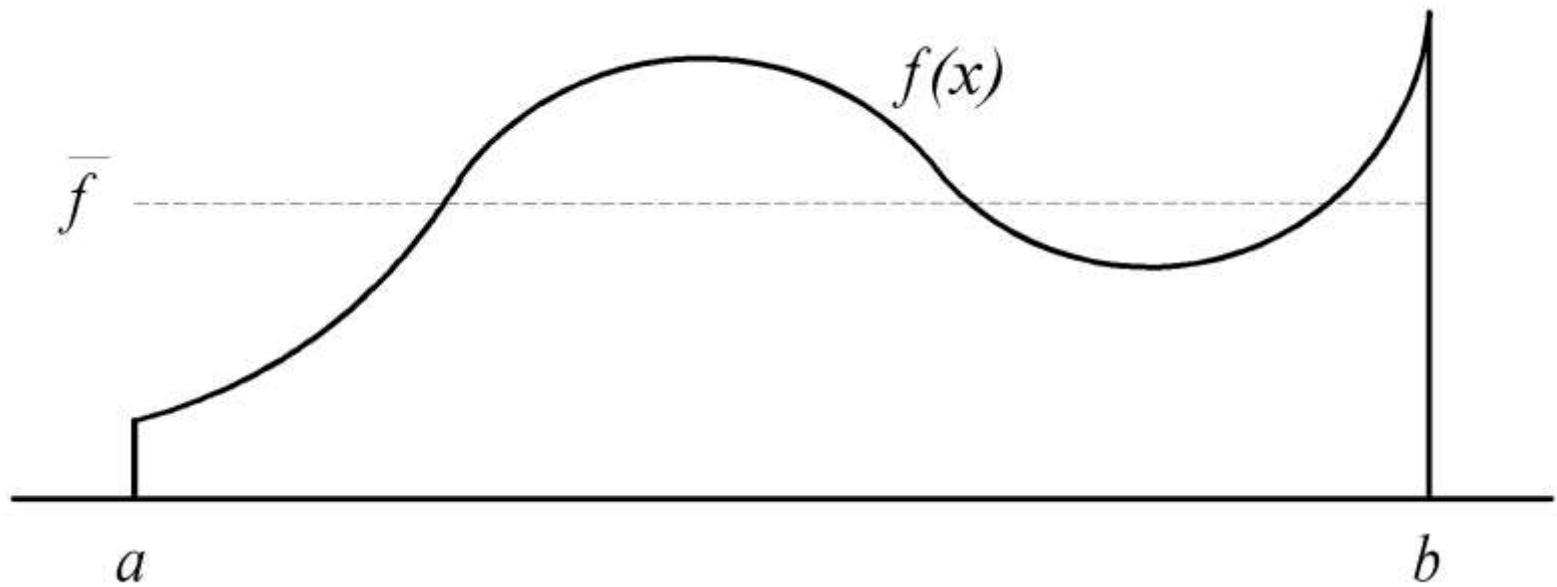
---

- u **Relative error is a way to measure the quality of an estimate**
- u **The smaller the error, the better the estimate**
- u **Relative error =  $|e-a|/a$** 
  - **$a$ : actual value**
  - **$e$ : estimated value**

# *Increasing Sample Size Reduces Error*

$n$	<i>Estimate</i>	<i>Error</i>	$1/(2n^{1/2})$
10	2.40000	0.23606	0.15811
100	3.36000	0.06952	0.05000
1,000	3.14400	0.00077	0.01581
10,000	3.13920	0.00076	0.00500
100,000	3.14132	0.00009	0.00158
1,000,000	3.14006	0.00049	0.00050
10,000,000	3.14136	0.00007	0.00016
100,000,000	3.14154	0.00002	0.00005
1,000,000,000	3.14155	0.00001	0.00002

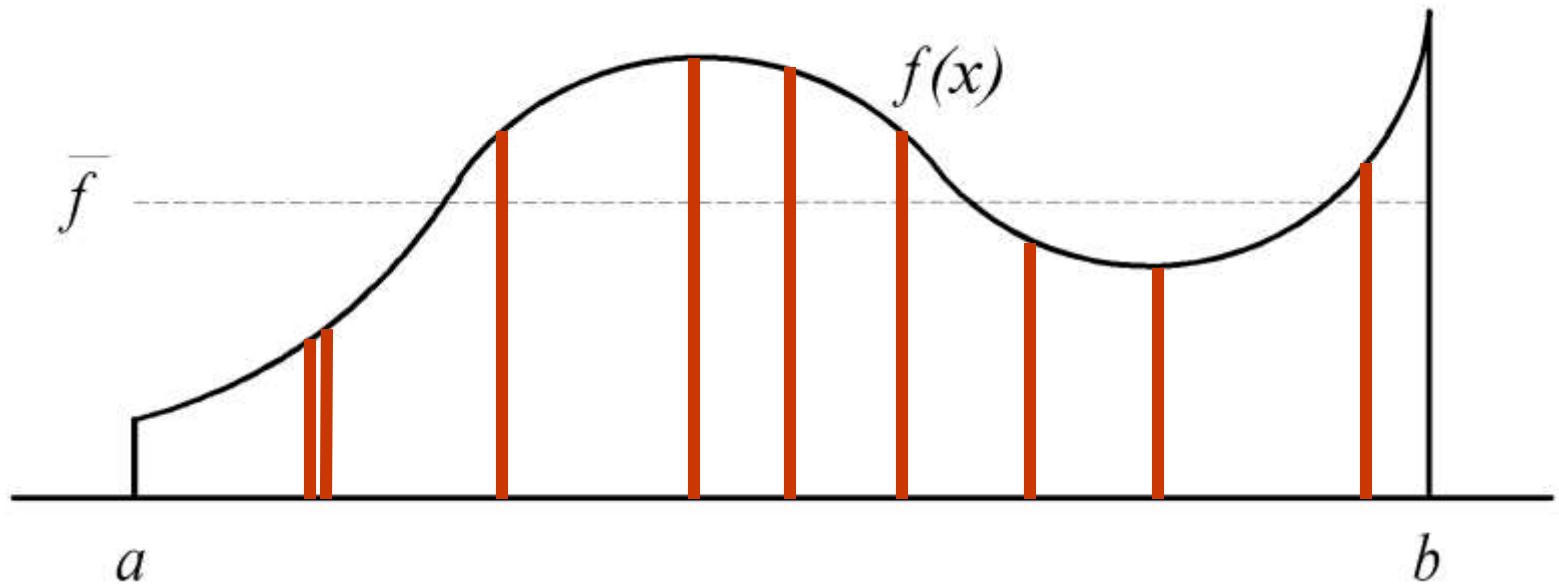
# Mean Value Theorem



$$\int_a^b f(x)dx = (b - a) \bar{f}$$

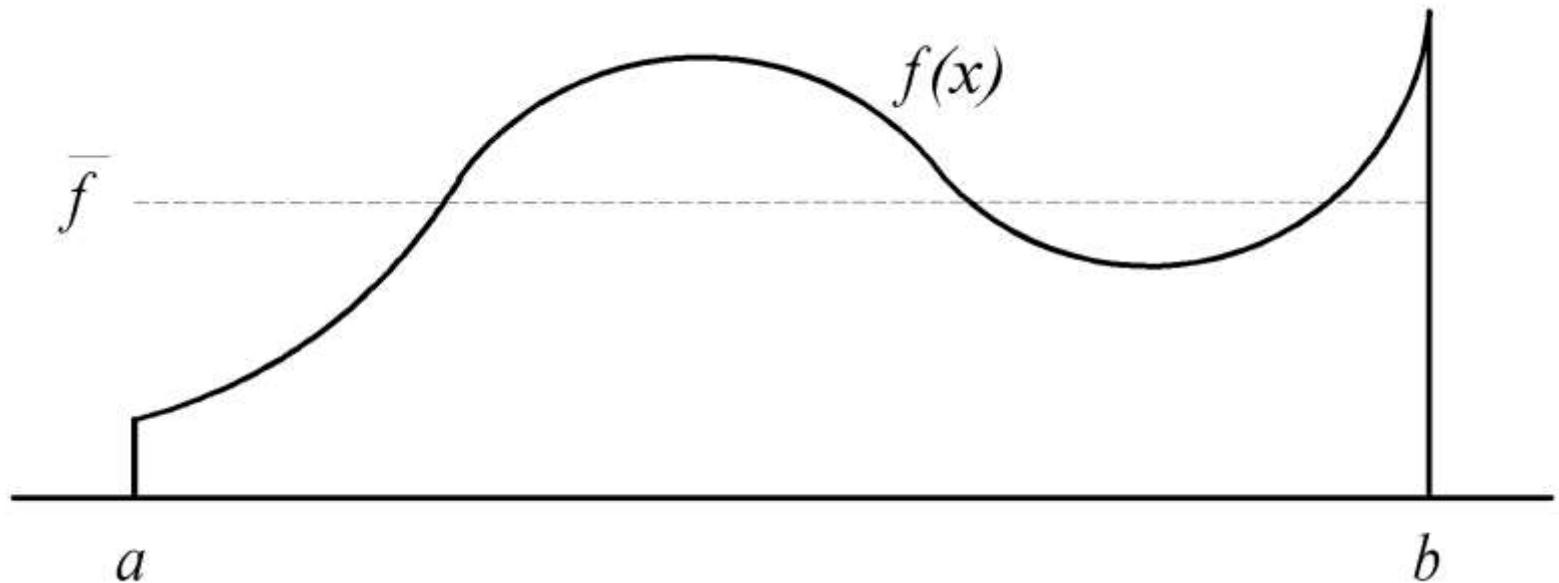
Michael J.Quinn, "Parallel Programming in C with MPI and OpenMP," 2003.

# Estimating Mean Value



The expected value of  $(1/n)(f(x_0) + \dots + f(x_{n-1}))$  is  $\bar{f}$

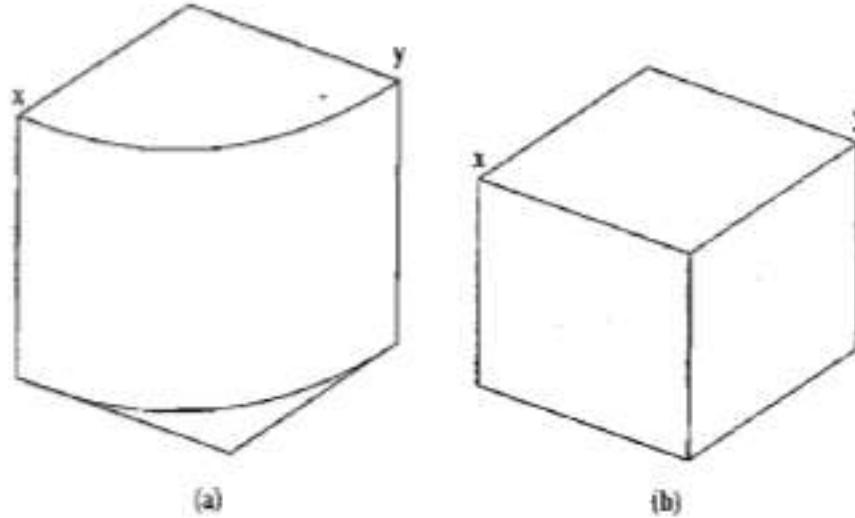
# Why Monte Carlo Works



$$\int_a^b f(x)dx = (b-a)\bar{f} \approx (b-a) \frac{1}{n} \sum_{i=0}^{n-1} f(x_i)$$

Michael J.Quinn, "Parallel Programming in C with MPI and OpenMP," 2003.

# Interpreting the $\pi$ Estimation Algorithm



use Monte Carlo to compute  $\text{volume}(a)$   
 $\text{volume}(a) / \text{volume}(b) = \pi/4$

# *Why Monte Carlo is Effective*

- u Error in Monte Carlo estimate decreases by the factor  $1/n^{1/2}$ 
  - (#sample:  $n \rightarrow 4n$ )  $\Rightarrow$  (error:  $\varepsilon \rightarrow \varepsilon/2$ )
- u Rate of convergence *independent* of integrand's dimension
- u Deterministic numerical integration methods do not share this property
- u Hence Monte Carlo superior when integrand has *6 or more* dimensions

# **Monte Carlo and Parallel Computing**

- u Monte Carlo methods often amenable to parallelism
  - Find an estimate about  $p$  times faster  
OR
  - Reduce error of estimate by  $p^{1/2}$
- u Based on the assumption that the random numbers are statistically independent
- u A principal challenge has been the development of good parallel random number generators
- u It is widely claimed that half of all super-computer cycles are dedicated to Monte Carlo calculations

# ***Random versus Pseudo-random***

---

- u Virtually all computers have “random number” generators
- u Their operation is deterministic
- u Sequences are predictable
- u More accurately called “pseudo-random number” generators
- u In this lecture “random” is shorthand for “pseudo-random”
- u “RNG” means “random number generator”

# *Properties of an Ideal RNG*

---

- u **Uniformly distributed**
- u **Uncorrelated**
- u **Never cycles**
- u **Satisfies any statistical test for randomness**
- u **Reproducible**
- u **Machine-independent**
- u **Changing “seed” value changes sequence**
- u **Easily split into independent subsequences**
- u **Fast**
- u **Limited memory requirements**

[Coddington, “Random number generators for parallel computers,” Technical Report, Syracuse University 1997.]

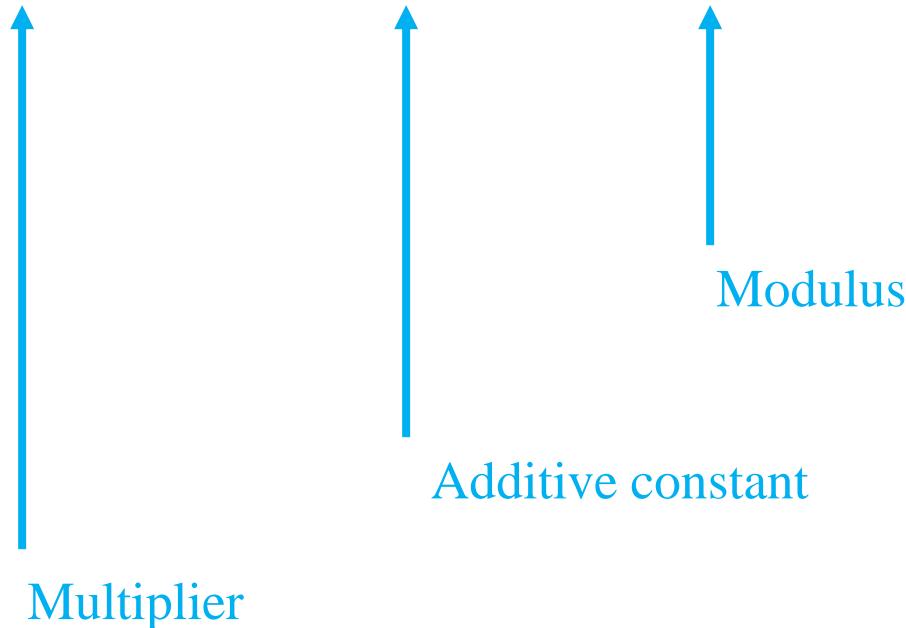
# No RNG is Ideal

---

- u Finite precision arithmetic  $\Rightarrow$  finite number of states  
 $\Rightarrow$  cycles
  - Period = length of cycle
  - If period > number of values needed, effectively acyclic
- u Reproducible  $\Rightarrow$  correlations
- u Often speed versus quality trade-offs
  - Since the time needed to generate a random number is typically a small part of the overall computation time of a program, speed is much less important than quality

# **Linear Congruential (线性同余) RNGs**

$$X_i = (a \times X_{i-1} + c) \bmod M$$



Sequence depends on choice of **seed**,  $X_0$

# *Period of Linear Congruential RNG*

- u Maximum period is  $M$
- u For 32-bit integers maximum period is  $2^{32}$ , or about 4 billion
- u This is too small for modern computers
- u Use a generator with at least 48 bits of precision

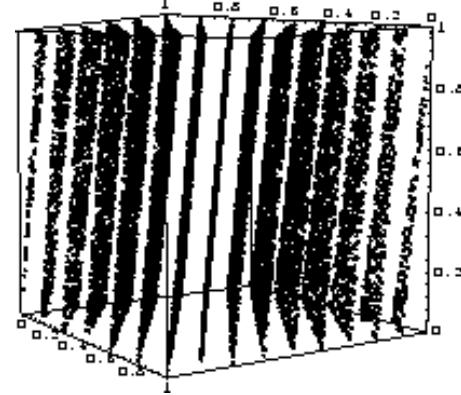
# *Producing Floating-Point Numbers*

---

- u  $X_i, a, c$ , and  $M$  are all integers
- u  $X_i$ 's range in value from 0 to  $M-1$
- u To produce floating-point numbers in range  $[0, 1)$ , divide  $X_i$  by  $M$

# Defects of Linear Congruential RNGs

- u Least significant bits correlated
  - Especially when  $M$  is a power of 2
- u  $k$ -tuples of random numbers form a lattice
  - Especially pronounced when  $k$  is large
  - Example: the infamous RANDU (IBM)
    - $M = 2^{31}$ ,  $a = 65539$ ,  $c = 0$ ,  $X_0 = 1$
- u Despite these flaws,
  - Linear congruential generators with 48 or more bits of precision and carefully chosen parameters “work very well for all known applications, at least on sequential computers”  
[Coddington 1997]



# *Lagged Fibonacci RNGs*

---

$$X_i = X_{i-p} * X_{i-q}$$

- u  $p$  and  $q$  are lags,  $p > q$
- u  $*$  is any binary arithmetic operation

- Addition modulo  $M$
- Subtraction modulo  $M$
- Multiplication (of odd integers) modulo  $M$
- Bitwise exclusive or

# *Properties of Lagged Fibonacci RNGs*

- u Require  $p$  seed values (more storage)
- u Careful selection of seed values,  $p$  and  $q$ , can result in very long periods and good randomness
  - [Coddington 1997] recommends setting  $(p,q)$  to at least (1279,1063)
- u Maximum period attainable with a  $b$ -bit  $M$ 
  - Exclusive or:  $2^p - 1$
  - Addition/subtraction:  $(2^p - 1)2^{b-1}$
  - Multiplication:  $(2^p - 1)2^{b-3}$

# **Ideal Parallel RNGs**

---

- u In addition to all properties of sequential RNGs,
- u No correlations among numbers in different sequences
- u Scalability
  - It should be possible to accommodate a large number of processes, each with its own stream(s)
- u Locality
  - A process should be able to spawn a new sequence of random numbers without interprocess communication

modified from Michael J.Quinn, "Parallel Programming in C with MPI and OpenMP," 2003.

# *Parallel RNG Designs*

---

- u Manager-worker
- u Leapfrog
- u Sequence splitting
- u Independent sequences

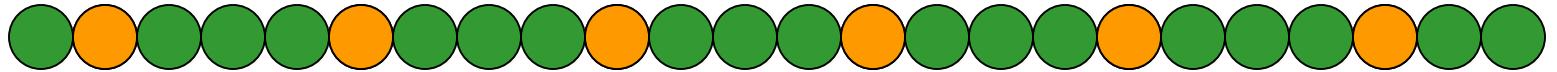
# **Manager-Worker Parallel RNG**

---

- u Manager process generates random numbers
- u Worker processes consume them
- u Disadvantages
  - Not scalable to an arbitrary number of processes
    - It may be difficult to balance the speed of the random number producer with the speed of the consumers of these numbers
  - Does not exhibit locality
    - It is communication-intensive

# *Leapfrog Method*

---



Process with rank 1 of 4 processes

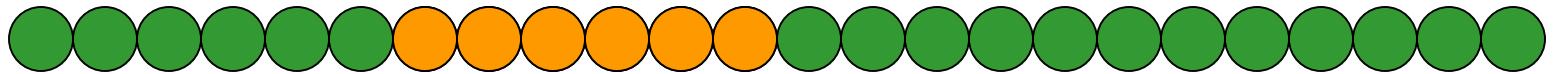
# *Properties of Leapfrog Method*

---

- u (Analog to cyclic decomposition)
- u Easy modify linear congruential RNG to support jumping by  $p$
- u Can allow parallel program to generate same tuples as sequential program
  - For example, to generate the 2-tuples  $(X_{2r}, X_{2r+1})$
- u Does not support dynamic creation of new random number streams

# *Sequence Splitting*

---



Process with rank 1 of 4 processes

# **Properties of Sequence Splitting**

---

- u (Analog to block decomposition)
- u Forces each process to move ahead to its starting point
  - It may take a long time
  - This only needs to be done at the initialization
- u Can be modified to support dynamic creation of new sequences
  - For example, a process creating a new stream could give up half of it section to the new stream

# **Independent Sequences**

---

- u Run sequential RNG on each process
- u Start each with different seed(s) or other parameters
- u Example: linear congruential RNGs
  - Works well for up to 100 streams with different additive constants [Percus & Kalos, 1989]
- u Example: lagged Fibonacci RNGs
  - SPRNG library provides around  $2^{1008}$  distinct streams
- u Supports goals of locality and scalability