

## cluster Programming

**Message-Passing Model** for **High Performance Cluster** (more data to computing)

- All process on same program

difference subtasks → have to communicate → network

Message Passing Interface Process (MPI Process)  $\sim$  1 thread  $\rightarrow$  1 function

1) mpirun [programname] -np 3 on cmd  
number of process = node

2) no job process

3) no job process  $\rightarrow$  MPI\_Init()  $\rightarrow$  MPI\_Finalize()

4) no job process of Rank (MPI\_Rank())  $\sim$  Process ID

### Communication

One-to-One / Point-to-Point	Collective
<pre>MPI-Send (void *buf, int count, MPI_Datatype dtype            dest, int tag, MPI_Comm comm)            destination rank, tag, label            group of process</pre>	1) One-to-Many • Broadcast vs Scatter • Gather vs Reduction
<pre>MPI-Receive (void *buf, int count, MPI_Datatype dtype               source, int tag, MPI_Comm comm,               from who               (rank))</pre>	2) All-to-All • All gather

**Map Reduce Programming Model** for **Big Data Cluster** (more computing to data)

- Google for Webpage indexing

→ Application: Inverted index

- High level

→ User friendly interface for distributed system

- Distributed File System

→ I/O intensive R/W on disk

- locality of data → minimized communication overhead

- Hadoop = de facto standard

**process** Map : Applied map function to the local data (in one node)

Reduce : sum Key value pairs (in one node)  $\sim$  sort, grouping → Disk R/W

Shuffle : Process each group of data (execute)

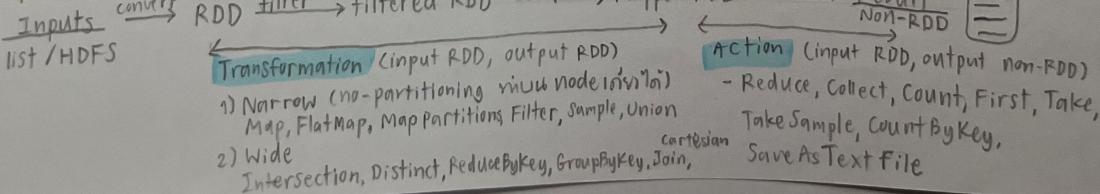
Make it Fault Tolerance	Task Crash Retry on other nodes	Node Crash Relaunch on other nodes	Slow tasks move to other node
-			

## Apache Spark

- Resilient Distributed Dataset (RDD) like Hadoop  $\rightarrow$  operation on RDD just like in HDFS  
- all operation  $\sim$  in memory  $\rightarrow$  more flexible

- lazy evaluation  $\rightarrow$  evaluate when needed to optimize cost

### Process



## Introduction to Distributed System

- not centralized
- autonomous computing elements (not just computer)
- connected by communication network (not just local)
- appears as a single coherent system
- passing messages
- components (software) located at networked computer

### Example

The internet, Email, Ebanking,  
Multiplayer online game, video streaming  
Social Network, Cloud storage, Sensor Network,  
Blockchain

### Trends

Ubiquitous Computing អំពីរស្តីក្នុងបន្ទាន់  
Cloud computing អំពីទិន្នន័យ (Computational service server)

## Characteristics of Distributed System

- Resource Sharing
  - Communication
  - Flexibility to build system
  - Scalability
  - Reliability
- ពាក្យសម្រាប់បង្កើតបច្ចុប្បន្ន / ពិនិត្យលក្ខណៈរបៀប

### Consequences of Distributed System

- Concurrent Program Execution
- Local information ព័ត៌មានផ្លូវការណាមួយនៅក្នុងបន្ទាន់
- No global clock → គ្មានបំបាត់នៃការងារ
- Independent failure

## Challenge (Parallel - performance, Distributed - Failure)

### Heterogeneity = និរនោគការងារ

Key : Standard and middleware

- Network - use internet protocol
- Computer Hardware - Big Endian  
Little Endian } converter
- OS - set standard
- Programming Language - Character Representation  
Data Structure
- Implementation - Developer (Coding Style)

### Transparency និងក្នុង single system

Key: user រួមចុះឱ្យការការងារដែលមានការងារ

និងការងារដែលមានការងារ

Access (How to access resource)

Location (Path / Link)

Concurrency (ការងារបានរំលែកបង្ហាញ)

Replication (ការកំណត់បញ្ជីការងារ)

Failure (ការងារបានរំលែកបង្ហាញ)

Mobility (ការកំណត់បញ្ជីការងារ)

performance

Scaling (ក្រឡាយការងារ)

### Security

Confidentiality

និងការអនុវត្តន៍

Integrity

ឯកតារក្នុងបញ្ជីការងារ

Availability

ការងារបានរំលែកបង្ហាញ

### Scalability

Extend physical Resource

Extend software Resource

Redesign (to avoid bottleneck)

### Failure Handling

Detect = ស្នើសុំ ex. Checksum, Message ID

Masking = Union

Retransmission រាយការ

Replication និងការអនុវត្តន៍

Drop ឯកតារក្នុងបញ្ជីការងារ

Tolerating = និងការអនុវត្តន៍

Recovery = rolled back

Redundancy = ឈឺការ

Concurrency និងការងារបានរំលែកបង្ហាញ

### Failures (ការងារបានរំលែកបង្ហាញ)

1) Network is reliable

2) Latency = 0 → Discard Timeout

3) Bandwidth =  $\infty$

4) Network is Secure

5) Topology doesn't change

6) There is one administrator

7) Transport cost = 0

8) Network is homogeneous

{Ring  
Star  
Tree}

## Distributed System Architecture

### Architectural Model ຈະນວຍດຸວຍ

- 1) Pattern of Placement of Components (Task)
- 2) Communication Paradigms (Interprocess / Remote procedure)
- 3) Roles and Responsibility (Client/Server/ Peer)

### Middleware

- interface abstraction
- interface between application and network

**Overlay Network** logical network over physical network

- 1) Structure - well define ex, tree, ring, star
- 2) Unstructure - randomly connection

## Client-Server Architecture

- servers may be clients of other servers
- can have more servers than clients → improve performance + increase reliability by replication

ex. Web

HTTP protocol

Network File System (NFS) Remote procedure call (RPC) ສອກຫາ ແລະ ປັບປຸງ ບໍລິສັດ

(ຈຸ່ວຍທີ່ມາ)

\* proxy servers and cache

ຂອງລົງທະບຽນ ລາຍໄລ ເພື່ອກຳນົດໃຫຍ່ ທັງໝາຍເຊີນຕົວໆ ຖ້າມີ proxy servers

↳ ລາຍໄລ Web Server

## Multitier Architecture

- layers of servers

- n-tier

## Peer-to-Peer Architecture

- All processes have similar roles = both provider and receiver

ex. P2P File Sharing

- 1) Napster P2P File Sharing ↗ Master index of centralized file → centralized, loaded server
- 2) P2P with circular (ring) ມີຄວາມຄວບຄົວ
- 3) P2P with Hierarchical (star) ↗ superpeer nodes (Ultra peer nodes) ຢູ່ກົກປູດເຄີຍໄວ້

## Hierarchical Architecture

ex. Network Time Protocol Servers

- Root nodes and leaf nodes

## Communication

### Fundamentals of Communication

(Any communication in the world)

Purposes, Media, Message, Protocols, Addressing, Transmission system

-Format

-How to deliver

-Content

### Basic Networking Concept

#### packet-switching technique

message

↓ divided, encapsulated into several packets

packet

↓ transmit via switches

receiver

### Network protocols : rules and formats

- sequence of messages
- addressing of endpoints
- format of data in messages

### Network Layers

- each layer has its own protocols to communicate at the same level
- provides services to upper layer
- uses services from lower layer
- propagated in sequence

#### Advantage

- 1) Reduce complexity  
Easy to understand and implement
- 2) Flexibility  
Implement and modify do not effect other layers

### Middleware Layer (between OS - Application)

- ① Lower Middleware Layer
- Message Passing
    - Sync (同步消息, 有序消息)
    - Async (异步消息, data push buffer window)
    - Blocking receive = សម្រេចទុកដាក់ឡើង
    - Non-blocking receive = សម្រេចទុកដាក់ឡើងនៅពេលបាន (polling)

- ② Upper Middleware Layer
- Message-Passing Interface (MPI) - group collective
  - Indirect Communication - through intermediary
    - space uncoupling, time uncoupling
  - Remote Procedure Call (RPC) - hidden communication
    - client call → wait → return
    - server request → reply
    - local procedure in server

### Encapsulation

- Each layer build a protocol data unit (PDU)

by adding header (and trailer)

→ result in multiple layers of headers and trailers

### ISO Open Systems Interconnection (OSI) protocol model

App

Presentation

Session

Transport

Network

Data Link

Physical

Upper  
Middle  
Lower

### TCP/IP protocol model

ex. HTTP, FTP, DNS, NTP, SSH

App → Message streams  
UDP or TCP

Transport → UDP or TCP

Internet → IP

Network Interface

Underlying Network

↳ circuits and Hardware that transmits data (Fiber / Cable)

→ transmit packets between nodes directly connected (Ethernet protocol)

→ Transfer packets between computer through routers (Internet protocol : IP)  
Switch in Datalink layer  
Router in Network layer

→ Lowest Level of Message (packet)

Messages are addressed to communication port

User datagram protocol (vs) Transmission control protocol (TCP)

- reduce overhead
- realtime application
- network management
- Failure Handling: checksum
- fail → drop full buffer → drop
- out-of-order message

### connection-oriented

→ connect routers

### reliable

### stream-oriented

Failure Handling → overhead  
Checksum, Sequence Number

Time out + Retry

- retransmission if necessary

- acknowledge to sender

• Message Queue manage by queue manager

- write allowing (1+ process)

- async

- put append to specific queue

- get block until nonempty and remove 1st message

- poll get but not block

- notify

• Publish-Subscribe System (Pub-Sub)

Distributed event-based system

Publisher → subscriber

publish, advertise, subscribe, Notify

# Web Architecture

Internet = Inter-networking infrastructure based on IP and TCP/UDP  
 World-Wide Web = Internet Application based on HTTP (hypertext transfer protocol)

static web pages > retrieving and displaying

dynamic web page > run scripts (the same as run local application)

## Web Enabling Technologies

Web server	- Resource provider
Web browser	- User Interface
HTTP	- App layer protocol over TCP/IP
HTML	- Data format for web page
URL (Uniform Resource Locator)	
MIME (Multipurpose Internet Mail Extensions)	/rfc2045/rfc2046/rfc2047 e-mail
Server-side Script	- CGI, PHP, JSP, Node.js
Client-side script	= JavaScript, VBScript

## HTTP protocol

- Text-based protocol — simple
- HTTP request to web server
- HTTP response from web server as HTML
- Non-SSL: HTTP (HTTP on a secured connection SSL)
- Port: (HTTP 80, port 443)

## Working Flow

- static web :
- 1) Do request to index.html
  - 2) Get HTML file from web server
  - 3) Display on screen
- dynamic web :
- 1) Do request to script file to web server
  - 2) web server execute via App server etc.
  - 3) Get HTML output to browser from web server

## Architecture

### 1) Single-Tier Architecture

- Monolithic application
- easy to manage
- scale-up only
- dumb remote terminal

### 2) Multi-Tier Architecture $\rightarrow$ presentation + logic + data

#### 2.1 Two-Tier Architecture (Thin Client)

- data sharing
- need fast clients
- app deployment

- (Presentation) and (Logic+Data)
- migration between client and server

#### 2.2 Three-Tier architecture

- UI/UX
- Scale-out, Flexibility
- More complex + Failure

#### 2.3 Four-Tier Architecture

- Browser + Webserver + App server + Database
- migration between client and server

## Internet Service Application / Service Oriented Architecture

Domain Application : E-mail, search, news, e-commerce, data storage

Hardware : Cluster

Software : Customized, multiter

Software Release Frequency : Days

Key metrics : Availability, functionality, scalability, manageability, throughput

## Performance Metrics

- 1) Concurrent users (Session) = Currently execute user
  - 2) Concurrent Connections = (1 user = n connections)
  - 3) Throughput = complete request / time
  - 4) Service Time / Response Time
- more resource  $\rightarrow$  larger, system upgrade  $\rightarrow$  faster

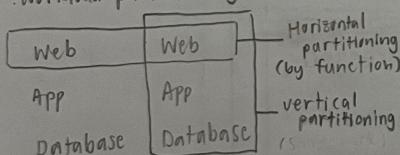
$$5) \text{Availability} = \frac{\text{Uptime}}{\text{Uptime} + \text{downtime}} \quad \text{four nines} = 99.99\%$$

downtime  $\rightarrow$  Failure  
 $\rightarrow$  Maintenance

## Scaling Web Application

### 1) Web Server/App Server optimization

: Workload partitioning (distributed workload)



### 2) Data layer optimization : Server load balancing

DNS Round-robin ; not high availability, DNS Caching

Reverse Proxy : User  $\rightarrow$  proxy  $\rightarrow$  back-end servers  
 pre-config proxy

Content Switching :根据不同 Content type/  
 Request type / User type

### 3) Network optimization : Memory / Storage Hierarchy

#### web caching

+ reduce latency, save bandwidth  
 Ø Popularity = everyone wants to access

Ø Locality = more recent content

Client : Proxy cache + Local cache

Server : Caching Load Balancer

#### Closeness = less latency

#### Content Delivery Network (CDN)

near content  $\rightarrow$  near end-users

high availability  $\rightarrow$  high performance

## Types of Large-Scale Internet Services // Usecase

### Types

Online service / Internet Portal

Content hosting service

Read Mostly — read > write

### Example

Hotmail

File Sharing

Wikipedia

### Front

functional partitioning

all the same

all the same

### Back

Single File, Single Database

data partitioned, in servers to n storages

full replication, all to all

## Time and Event Ordering

Event - the occurrence of a single action

- ex. sent / receive message, change of memory
- Some need timestamp accurately ex. auditing

**Real Time Clock (RTC)** = electronic oscillator circuit

- using the vibration of quartz crystal to control the frequency
- of I<sub>C</sub> microcontroller clock generator

→ not very accurate  
due to variation of quartz  
and temperature

**Clock skew** = Time A - Time B (in duration)

**Clock Drift** = Frequency A - Frequency B (in speed)

want **Monotonicity** time always increases

**Continuity** time doesn't jump

→ ດ້ວຍເລີກໄວ້ໃຫຍ່ ປົບປັດ ອານຸຍາວ ດູວ່າ Rate (ອະນຸຍາວ)

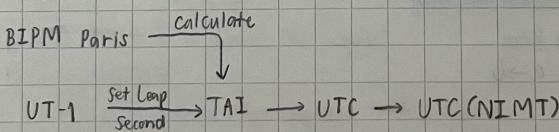
Sync with highly accurate  
external clock

→ **Atomic clock** (Cesium)

## Standard Time

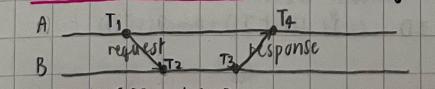
### Coordinated Universal Time (UTC)

- International Atomic Time (TAI) = average of 400 atomic clocks worldwide
  - Universal Time (UT-1) = based on rotation of Earth and around the Sun
  - Adjusted by adding leap seconds (same concept as leap year)  
= adding 1 second to some of 30 June and 31 Dec
  - broadcast through short-wave radio and GPS satellites
- computers sync to this



## Clock Synchronization and Reference Clock

with reference clock **Christian's Algorithm**



with pair (offset, delay)

$$\text{Offset} = \frac{1}{2}((T_2 - T_1) + (T_3 - T_4))$$

$$\text{Delay} = \frac{1}{2}((T_4 - T_1) - (T_3 - T_2))$$

in n round trip delay window

if offset > 0, set time = currTime + offset

else, slow down clock

without reference clock **Berkeley Algorithm**

- 1) leader asks times from others
- 2) average
- 3) sendback the offset

## Network Time Protocol (NTP)

- to sync clocks over the internet, use **UDP** (user datagram protocol) = less delay
- many sync time servers → hierarchy level = stratum
- reliable and scalability

## Logical Time

- physical time has delay.
- relatively order, not absolutely time

### Happened-Before Relation by Lamport

- potential causal ordering
- denote by  $\rightarrow$
- HB1 (intraprocess)** : If  $\exists$  process  $p_i$  ( $a \rightarrow_i b$ ) then ( $a \rightarrow b$ )
- HB2 (interprocess)** :  $m_i \text{ send}(m) \rightarrow m_j \text{ receive}(m)$
- HB3 (transitive)** : If ( $a \rightarrow b$  and  $b \rightarrow c$ ) then ( $a \rightarrow c$ )
- If they are not related by  $\rightarrow$  (cannot tell who is before)  
then they are said to be **concurrent**, denote by  $\parallel$   
(if no  $a \rightarrow b$  then all)

### Logical Clock / Lamport Clock

- monotonically increasing software **counter**
- use concept of HB relation
- logical timestamp for each process  $P_i$  with logical clock  $L_i$  (1 process / 1 LC)
- LC1** :  $L_i$  is incremented by 1 **before** each event at process  $P_i$
- LC2** : (a) process  $P_i$  **send**( $m$ ) to process  $P_j$   $t = L_i$  **before** receive( $m$ )  
(b) process  $P_j$  **receive**( $m, t$ ) **set**  $L_j = \max(L_j, t)$ ,  
**in LC1**, **before** receive( $m$ )

so  $a \rightarrow b$  imply  $L(a) < L(b)$   
but  $L(a) < L(b)$  **not** imply  $a \rightarrow b$

- **Timestamp** of an event by process  $a$  is **Logical Timestamp** ( $L$  բարձրացնելու)
- $L$  բարձրացնելու **delay** կամ  $d$
- **totally ordered**  $L$  →  $m$  **Totally ordered logical clocks**
- $L$  = **Logical Timestamp + process ID** ( $L = Host ID + Process ID$ )
- \* There is **no way** to have perfect knowledge on ordering of events

## coordination

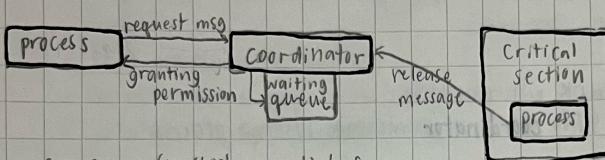
- value and action agreement

**Mutual exclusion:** exclusive access to some resource, critical section

Requirement	Mutual exclusion	
	Safety	Liveness
Safety	- one process run critical section at a time	process performing
Liveness	- Request always get response	No deadlock / starvation
Fairness	- fairness among processes	processes
Permission-based	- avoid deadlock by monitoring process	
Token-based	- avoid token circulate access resource	by token

## Centralized Permission-based Algorithm

- One process is elected as the **coordinator**



- 2 messages (request + grant) before 1 message (release) after
- Disadvantage: delay on enter(), sync delay, single point of failure and performance bottleneck

## Distributed Permission-based Algorithm : Ricart and Agrawala Algorithm

- sends message (resource name, process id, current logical time)

to all processes and waits until all OK

- မျှတော်မူနဲ့ 3 ပုံစံ
  - မျှတော်မူနဲ့ → return OK
  - မျှတော်မူနဲ့ စီရင် → queue the request
  - မျှတော်မူနဲ့ စီရင် → queue the request
  - မျှတော်မူနဲ့ စီရင် → return OK
- Assume no message losts

- မျှတော်မူနဲ့ စီရင် → return OK to every process in queue and clear queue

- Disadvantage: All processes involved in all permission requests, N-point of failure

## Token-ring Algorithm

- Logical ring
- ကြော်လွှာ → စွမ်းဆောင်ရွက် Token → ဖျော်လွှာ နဲ့ Token ပေါ်
- ပို့မယ်လွှာ → ထပ်မံပို့
- Disadvantage: 1 ရှိနိုင်ပေးသူ မျှတော်မူနဲ့ လျှို့ဝှက် delay များ

Election Algorithm: choosing a unique process to play a particular role  
assume unique ID, select highest ID to be the role

### Ring-based Election Algorithm

- notices that current coordinator is down
- send **election** message (this. processID)
- receives and compare receive.ID, this.ID
- send max processID until receive.ID = this.ID (maxID)
- send **elected** message (elected. processID)
- receives and set new coordinator

### Bully Algorithm

- knows ID of all processes.
- finds ID winner
- if running & return OK to all others
- if not OK + sends **coordinator** message to other process

## ⑯ Fault Tolerance

- partial failure, یا یوں نہ نہیں
- presence of faults →  $\text{mean time to failure}$  = fault tolerance
- Availability, Reliability, Safety, Maintainability

### Failure Metrics

- Mean Time To Failure = avg(uptime)
- Mean Time To Repair = avg(downtime)
- Mean Time Between Failure = MTTF + MTTR
- Availability =  $\frac{\text{MTTF}}{\text{MTBF}} = \frac{\text{uptime}}{\text{total}}$

### Failure Model

Crash failure	کھنڈا
Timing failure	میلے میلے
Omission failure	بھاٹاک
(Receive/send)	لے لے
Response failure	میلے میلے
Arbitrary failure	کوئی نہ پیدا

### Failure Handling

- Detect - Checksum
- Masking - Hide (Retransmit, Replication, Dropping)
- Tolerating - inform users
- Recovery - rollback data
- Redundancy - 2 replica of servers and routes

k-fault tolerant w/  $M \geq k$  components  $\rightarrow$  no failures

### (8) Agreement Problems

- Uncoordinated processes w/o coordinator

Consensus Problem (centralized)

- finding a consensus

- centralized coordinator

- Majority → finding an agreed value (min, max, majority)  
→ return decision

-  $2k+1$  process with  $k$ -fault can have majority

Byzantine Generals Problem (noisy - wrong)

- accept arbitrary failure

- cross check and select majority (commander may send wrong message)  
(any node may have response failure)

-  $3k+1$  process with  $k$ -fault

- Correct processes agree on the same value

- Unreliable = Arbitrary + Noisy + Unresponsive

### CAP Theorem

Consistency : some node must reject requests (no most recent writes)

Availability : responses will be inconsistent