1



2



3



4



5



6

**Slide 7**

## OBJECTIVES – 2/13

- Questions from 2/6
- **Midterm Grading In Progress - Targeting Review Thursday**
- Assignment 2: Key/Value Store
  - Java Maven project template files posted
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - Chapter 3.5: Resource (Code) Migration (*light-review*)
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - Chapter 4.2: RPC (light-review)
  - Chapter 4.3: Message Oriented Communication

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.7

**Slide 8**

## OBJECTIVES – 2/13

- Questions from 2/6
- Midterm Grading In Progress - Targeting Review Thursday
- **Assignment 2: Key/Value Store**
  - **Java Maven project template files posted**
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - Chapter 3.5: Resource (Code) Migration (*light-review*)
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - Chapter 4.2: RPC (light-review)
  - Chapter 4.3: Message Oriented Communication

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.8

**Slide 9**

## ASSIGNMENT 2

- **Find Teammates**: signup posted on Canvas under 'People'
- GenericNode.tar.gz includes Dockerfile examples
- GenericNode.tar.gz assumes Java 11
- TCP/UDP/RMI Key Value Store
- Implement a 'GenericNode' project which assumes the role of a client or server for a Key/Value Store
- Recommended in Java 11 LTS
- Client node program interacts with server node to put, get, delete, or list items in a key/value store

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.9

**Slide 10**

## USING JAVA 11 IN NETBEANS

- In Netbeans IDE, under Tools menu, 'Java Platforms', be sure to install and select JDK 11



Java Platform Manager

Use the Javadoc tab to register the API documentation for your JDK in the IDE.
Click Add Platform to register other Java platform versions.

Platforms:
- Java SE
  - JDK 11
  - JDK 11 (Default)

Platform Name: JDK 11 (Default)
Platform Folder: /usr/lib/jvm/java-11-openjdk-amd64

- On left-hand Project menu, right-click on 'GenericNode' project
- Select Properties
- Under Build | Compile, be sure Java Platform is JDK 11
- Under Sources, be sure Source/Binary Format is 11

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.10

**Slide 11**



**Slide 12**

## OBJECTIVES – 2/13

- Questions from 2/6
- Midterm Grading In Progress - Targeting Review Thursday
- Assignment 2: Key/Value Store
  - Java Maven project template files posted
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - **Chapter 3.5: Resource (Code) Migration (*light-review*)**
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - Chapter 4.2: RPC (light-review)
  - Chapter 4.3: Message Oriented Communication

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.12

## CH. 3.5: RESOURCE (CODE) MIGRATION

L11.13

13

---

## RESOURCE MIGRATION

- To support on-the-fly reorganization of distributed systems, at times there is interest in resource migration

- Can consider various types of resource migration
  - **Code migration**: source code, libraries
  - **Process migration**: a running job/task
  - **VM migration**: an entire virtual server!

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.14

14

---

## TYPES OF CODE MIGRATION

- Distributed systems can support more than **passing data**

- Some situations call for **passing programs** (e.g. *code*)

- **Live migration** – moving code while it is executing

- **Portability** – transferring code (running or not) across heterogeneous systems:
  Mac OS X → Windows 10 → Linux

- Code migration enables **flexibility** of distributed systems
  - Topologies can be dynamically reconfigured on-the-fly

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.15

15

---

## PROCESS MIGRATION

- Move an entire process from one node to another

- Motivation is always to address performance

- Process migration is slow, costly, and intricate
  - Need to pause, save intermediate state, move, resume
  - Consider application *specific* vs. *agnostic* approaches

- What would be:
  an **application agnostic** approach to migration?
  an **application specific** approach?

- What are advantages and disadvantages of each?

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.16

16

---

## PROCESS MIGRATION - 2

- **Move processes**:
  from heavily loaded → lightly loaded nodes

- When do we consider a node as heavily loaded?
  - Load average
  - CPU utilization
  - CPU queue length

- Which process(es) should be moved?
  - Must consider *resource requirements* for the task

- Where should process(es) be moved to?

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.17

17

---

## MOTIVATIONS FOR MIGRATION

- Can migrate **processes** or entire **virtual machines**

- **Goals:**
  o Off-loading machines: reduce load on oversubscribed servers
  o Loading machine: ensure machine has enough work to do
  o Minimize total hosts/servers in use to save energy/cost

- **VM migration:**
- Migrate complete VMs with apps to lightly loaded hosts
- Generally, VM migration is easier than process migration

- **Is VM migration application specific or agnostic?**

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.18

18

---

## LINUX CRIU

- Linux (CRIU) Checkpoint restore in userspace
- Linux tool: https://www.criu.org/
- Supports freezing a running application (or part of it) to create a checkpoint to persistent storage (e.g. disk) as a collection of files.
  - This means saving the state of RAM to disk
- Can use checkpoint files to restore and run the application from the point it was frozen at.
- Distinctive feature of CRIU is that it can be run in the user space (CPU user mode), rather than in kernel mode.
- CRIU can save a Docker container's state for migration elsewhere

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.19

19

## LOAD DISTRIBUTION ALGORITHMS

- Make decisions concerning allocation and redistribution of tasks across machines
- Provide resource management for compute intensive systems
- Often CPU centric
  - Algorithms should also account for other resources
  - Network capacity may be larger bottleneck that CPU capacity

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.20

20

## WHEN TO MIGRATE?

- Decisions to migrate code often based on qualitative reasoning or adhoc decisions vs. formal mathematical models
  - Difficult to formalize solutions due to heterogeneous composition and state of systems and networks
- Is it better to migrate code or data?
- What factors should be considered?
  - Size of code
  - Size of data
  - Available network transfer speed
  - Cost of data transfer
  - Processing power of nodes
  - Cost of processing
  - Are there security requirements for the data?

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.21

21

## APPROACHES TO CODE MIGRATION

- Traditional clients
  - Client interacts with server using specific protocol
  - Tight coupling of client->server limits system flexibility
  - Difficult to change protocol when there are *many* clients
- Dynamic web clients
  - Web browser downloads client code immediately before use
  - New versions can readily be distributed

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.22

22

## DYNAMIC WEB CLIENTS

- Advantages
  - Client code loaded in as necessary
  - Discarded when no longer needed
  - Can easily change the client/server protocol
- Disadvantages
  - Security: we have to trust the code
  - Downloading client requires network bandwidth & time

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.23

23

## CODE MIGRATION

- Sender-initiated: (upload the code)... e.g. Github
- Receiver-initiated: (download the code)... e.g. web browser
- Remote cloning
  - Produce a copy of the process on another machine while parent runs

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.24

24

## CODE MIGRATION - 2

- What is migrated?
  - *Code* segment
  - *Resource* segment (device info)
  - *Execution* segment (process info: data, state, stack, PC)
- **Weak mobility**
  - Only *code* segment, no state
  - Code always restarts
- **Strong mobility**
  - *Code* + *execution* segment
  - Process stopped, state saved, moved, resumed
  - Represents true *process migration*

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.25

25

## CODE MOBILITY TYPES

- \* indicates what is modified
- **CS: Client-Server**
- **REV: Remote Evaluation**
- **CoD: Code-on-demand**
- **MA: Mobile agents**

- Where does state get modified?

- State is stored in **exec**



February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.26

26

## MIGRATION OF HETEROGENEOUS SYSTEMS

- Assumption: code will always work at new node
- Invalid if node architecture is different (*heterogeneous*)
  - X86, ARM, MAC, etc.

- What approaches are available to migrate code across heterogeneous systems?

- Intermediate code
  - 1970s Pascal: generate machine-independent intermediate code
  - Programs could then run anywhere
  - **Today:** web languages: Javascript, Java

- VM Migration

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.27

27

## VIRTUAL MACHINE MIGRATION

- Four approaches:

1. **PRECOPY**: Push all memory pages to new machine *(slow)*, resend modified pages later, transfer control
2. **STOP-AND-COPY**: Stop the VM, migrate memory pages, start new VM
3. **ON DEMAND**: Start new VM, copy memory as needed
4. **HYBRID**: PRECOPY followed by brief STOP-AND-COPY

- **What are some advantages and disadvantages of 1-4?**

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.28

28

1. **PRECOPY**: Push all memory pages to new machine *(slow)*, resend modified pages later, transfer control
2. **STOP-AND-COPY**: Stop the VM, migrate memory pages, start new VM
3. **ON DEMAND**: Start new VM, copy memory pages as needed
4. **HYBRID**: PRECOPY and followed by brief STOP-AND-COPY

- **What are some advantages and disadvantages of 1-4?**
  - (+) 1/3: no loss of service
  - (+) 4: fast transfer, minimal loss of service
  - (+) 2: fastest data transfer
  - (+) 3: new VM immediately available

  - (-) 1: must track modified pages during full page copy
  - (-) 2: longest downtime - unacceptable for live services
  - (-) 3: prolonged, slow, migration
  - (-) 3: original VM must stay online for quite a while
  - (-) 1/3: network load while original VM still in service

L11.29

29



30

## OBJECTIVES – 2/13

- Questions from 2/6
- Midterm Grading In Progress - Targeting Review Thursday
- Assignment 2: Key/Value Store
  - Java Maven project template files posted
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - Chapter 3.5: Resource (Code) Migration (*light-review*)
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - Chapter 4.2: RPC (light-review)
  - Chapter 4.3: Message Oriented Communication

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.31 |

31

---

## CH. 4 COMMUNICATION

L11.32

32

---

## CHAPTER 4

- 4.1 Foundations
  - Protocols
  - Types of communication
- 4.2 Remote procedure call
- 4.3 Message-oriented communication
  - Socket communication
  - Messaging libraries
  - Message-Passing Interface (MPI)
  - Message-queueing systems
  - Examples
- 4.4 Multicast communication
  - Flooding-based multicasting
  - Gossip-based data dissemination

*Reviews and builds on content from Ch. 2/3*

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.33 |

33

---

## CH. 4.1: FOUNDATIONS

L11.34

34

---

## LAYERED PROTOCOLS

- Distributed systems lack shared memory
- All distributed system communication is based on sending and receiving low-level messages
  - P → Q
- **O**pen **S**ystems **I**nterconnection Reference Model (OSI Model)
  - Open systems communicate with any other open system
  - Standards govern format, contents, meaning of messages
  - Formalization of rules forms a **communication protocol**

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.35 |

35

---

## LAYERED PROTOCOLS - 2

- Protocols provide a **communication service**
- **Two service types:**
  - **Connection-oriented**: sender/receiver establish connection, negotiate parameters of the protocol, close connection when done
  - Physical example: telephone
  - **Connectionless**: No setup.  Sender sends. Receiver receives.
  - Physical example: Mailing a letter

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.36 |

36

## OSI MODEL REVISITED

- Physical layer: just sends bits → … 0 0 0 1 0 1 1 0 1 1 …
- Data link layer: Groups bits into frames
  - Provides error correction via *checksum*
  - Special bit pattern at start/end of frame

37

## OSI MODEL - 2

- Data link layer:
  - **Checksum**: computed by adding all bytes in frame in particular way
  - Added to message
  - Receiver removes checksum, recomputes checksum, and compares
  - If receiver and sender agree, frame is considered correct
  - Receiver can request failed frames to be resent
  - Frames assigned sequence numbers *in the header*
- Network layer:
  - Sometimes referred to as the *Internet layer*
  - On WANs sending msgs between client/server requires routing
  - Provides addressing using IPV4 (32-bit), IPV6 (64-bit)

38

## OSI MODEL - 3

- Network layer:
  - Helps with routing network traffic
  - Shortest route (# of hops) may not be the best route
  - Minimizing delay (latency) is paramount
  - Routing algorithms: use long-term average network conditions, or try to adapt to changing conditions
  - ICMP Protocol: Internet Control Message Protocol
  - Not typically for sending data, used for diagnostic/control purposes
  - ICMP Examples: (*ping*, *traceroute*)

39

## OSI MODEL - 4

- Internet Control Message Protocol (ICMP)
  - 8 bytes header: 4 fixed, 4 variable

| ICMP Header Format | | | | |
|---|---|---|---|---|
| Offsets | Octet | 0 | 1 | 2 | 3 |
| Octet | Bit | 0 1 2 3 4 5 6 7 | 8 9 10 11 12 13 14 15 | 16 17 18 19 20 21 22 23 | 24 25 26 27 28 29 30 31 |
| 0 | 0 | Type | Code | Checksum | |
| 4 | 32 | Rest of Header | | | |

  - Example message types:
  - 0- echo reply (**PING**), 3- destination unreachable, 4- source quench (congestion control), 5- redirect message, 8- echo request (**PING**), 9- router advertisement
  - Others: 10 (router solicitation), 11 (time exceeded), 12 (parameter problem), 13 (timestamp), 15 (info request), 16 (info reply), 17 (address mask request), 18 (address mask reply), 30-39 (**traceroute**), 40 (security failures), 42 (ext echo request)…255

40

## OSI MODEL - 5

- Transport layer:
  - Provides reliable connections
  - Reorganizes packets arriving out of sequence
  - Requests delivery of missing packets

  1. Breaks application layer protocol messages into pieces to transmit
  2. Assigns messages sequence numbers
  3. Sends all messages

41

## OSI MODEL - 6

- Transport layer provides an infallible "message pipe"
  - Put messages in
  - Always come out undamaged, in correct order

- Transport layer protocols:
  - TCP: Transmission Control Protocol (connection-oriented)
  - UDP: Universal Datagram Protocol (connectionless)

42

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.37

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.38

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.39

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.40

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.41

February 13, 2024    TCSS558: Applied Distributed Computing [Winter 2024]
School of Engineering and Technology, University of Washington - Tacoma    L11.42

## OSI MODEL - 7

- Other transport protocols
  - Real-time transport protocol (RTP): real-time data, no data delivery guarantee
  - Streaming Control Transmission Protocol (SCTP): alternative to TCP
- Higher-level protocols:
- **Session layer**: mechanisms for opening, closing, managing session between communicating processes
- **Presentation layer**: deals with syntactical meaning of messages
  - Presentation services convert data among formats, for example:
  - from extended binary coded decimal interchange code (EBCDIC) to ASCII
- **Application layer**: protocols that don't fit into other layers
  - Many protocols: FTP, SFTP, HTTP, etc. etc.

43

## OSI MODEL - 8

- Each OSI layer contributes overhead bits to the message
- Layers append data to front (and maybe end) of the message
- Receiver strips off headers as the message goes up the OSI model stack:

*physical → data-link → network → transport → application*

44

## PROTOCOL STACK

- Collection of layers used for communication from OSI model

45

## MIDDLEWARE PROTOCOLS

- Middleware is reused by many applications
- Provide needed functions applications are built and depend upon
  - For example: communication frameworks/libraries
- Middleware offer many general-purpose protocols
- Middleware protocol examples:
  - **Authentication protocols**: supports granting users and processes access to authorized resources
  - Doesn't fit as an "application specific" protocol
  - Considered a "Middleware protocol"

46

## MIDDLEWARE PROTOCOLS - 2

- **Distributed commit protocols**
  - Coordinate a group of processes (nodes)
  - Facilitate all nodes carrying out a particular operation
  - Or abort transaction
  - Provides distributed atomicity (all-or-nothing) operations
- **Distributed locking protocols**
  - Protect a resource from simultaneous access from multiple nodes
- **Remote procedure call**
  - One of the oldest middleware protocols

47

## MIDDLEWARE PROTOCOLS - 3

- **Message queueing services**
  - Support synchronization of data streams
  - Transfer real-time data
  - Distributed and scalable implementation
- **Multicast services**
  - Scale communication to thousands of receivers spread across the Internet

48

## MIDDLEWARE PROTOCOLS - 3

- **Message queueing services**
- ~~Support synchronization of data~~

**KEY:** middleware protocols offer functionality to satisfy the software requirements of _many_ applications

Middleware functions are general, application-independent in nature

Functions are so commonly needed they are offered in reusable frameworks / libraries

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.49

49

## ADAPTED REFERENCE MODEL



February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.50
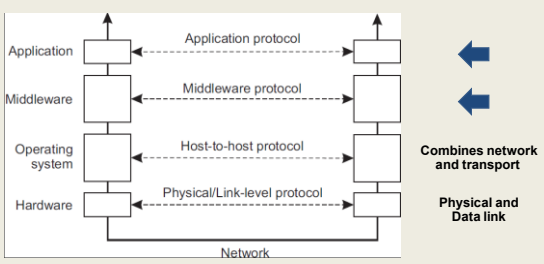
50

## TYPES OF COMMUNICATION

- **Persistent communication**
  - Message submitted for transmission is stored by communication middleware as long as it takes to deliver it
  - Example: email system (SMTP)
  - Receiver can be offline when message sent
  - Temporal decoupling (delayed message delivery)

- **Transient communication**
  - Message stored by middleware only as long as sender/receiver applications are running
  - If recipient is not active, message is dropped
  - Transport level protocols typically are transient (_no msg storage_)

- **What OSI protocol level is the SMTP Protocol?**

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.51

51

## TYPES OF COMMUNICATION - 2

- **Asynchronous communication**
  - Client does not block, continues doing other work
- **Synchronous communication**
  - Client blocks and waits
- **Three types of blocking (_synchronous_)**
  1. Until middleware notifies it will take over delivering _request_
  2. Sender may block until _request_ has been delivered
  3. Sender waits until _request_ is processed and result is returned
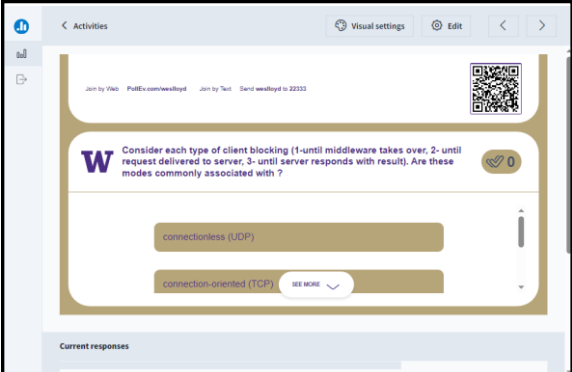- **Persistence + synchronization (blocking)**
  - Common scheme for message-queueing systems
  - _**Publish message to queue**_: block until message delivered to queue

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.52

52



53

## OBJECTIVES – 2/13

- Questions from 2/6
- Midterm Grading In Progress - Targeting Review Thursday
- Assignment 2: Key/Value Store
  - Java Maven project template files posted
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - Chapter 3.5: Resource (Code) Migration (_light-review_)
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - **Chapter 4.2: RPC (light-review)**
  - Chapter 4.3: Message Oriented Communication

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington  - Tacoma | L11.54

54

## WE WILL RETURN AT 4:55 PM

55

## CH. 4.2: RPC (LIGHT-REVIEW)

56

---

## RPC – REMOTE PROCEDURE CALL

- In a nutshell,
- Allow programs to call procedures on other machines
- Process on **machine A** calls procedure on **machine B**
- Calling process on **machine A** is suspended
- Execution of the called procedure takes place on **machine B**
- Data transported from caller **(A)** to provider **(B)** and back **(A)**.
- No message passing is visible to the programmer
- **Distribution transparency**: make remote procedure call look like a local one
- `newlist = append(data, dbList)`

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.57

57

## RPC - 2

- Transparency enabled with client and server "stubs"
- Client has "stub" implementation of the server-side function
- Interface exactly same as server side
- But client **DOES NOT HAVE THE IMPLEMENTATION**
- **Client stub**: packs parameters into message, sends *request* to server. Call blocks and waits for reply
- **Server stub**: transforms incoming *request* into local procedure call
- Blocks to wait for *reply*
- Server stub unpacks *request*, calls server procedure
- *It's as if the routine were called locally*

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.58

58

---

## RPC - 3

- Server packs procedure *results* and sends back to client.
- Client "*request*" call unblocks and data is unpacked
- Client can't tell method was called remotely over the network... *except for network latency...*
- Call abstraction enables clients to invoke functions in alternate languages, on different machines
- Differences are handled by the RPC "framework"

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.59

59

## RPC STEPS

1. Client procedure calls client stub
2. Client stub builds message and calls OS
3. Client's OS send message to remote OS
4. Server OS gives message to server stub
5. Server stub unpacks parameters, calls server
6. Server performs work, returns results to server-side stub
7. Server stub packs results in messages, calls server OS
8. Server OS sends message to client's OS
9. Client's OS delivers message to client stub
10. Client stub unpacks result, returns to client

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.60

60

## PARAMETER PASSING

- **STUBS**: take parameters, pack into a message, send across network

- Parameter marshaling:
- `newlist = append(data, dbList)`
- Two parameters must be sent over network and correctly interpreted

- Message is transferred as a series of bytes
- Data is serialized into a "stream" of bytes
- Must understand how to unmarshal (unserialize) data

- Processor architectures vary with how bytes are numbered: Intel (right→left), older ARM (left→right)

61

## RPC: BYTE ORDERING

- Big-Endian: write bytes left to right (ARM)

- Little-endian: write bytes right to left (Intel)

- Networks: typically transfer data in Big-Endian form

- Solution: transform data to machine/network independent format

- Marshaling/unmarshaling: transform data to neutral format

62

## RPC: PASS-BY-REFERENCE

- Passing by value is straightforward
- Passing by reference is challenging
- Pointers only make sense on local machine owning the data
- Memory space of client and server are different

- Solutions to **RPC pass-by-reference**:
1. Forbid pointers altogether
2. Replace pass-by-reference with pass-by-value
   - Requires transferring entire object/array data over network
   - **Read-only optimization**: don't return data if unchanged on server
3. Passing global references
   - Example: file handle to file accessible by client and server via shared file system

63

## RPC: DEVELOPMENT SUPPORT

- Let developer specify which routines will be called remotely
  - Automate client/server side stub generation for these routines

- Embed remote procedure call mechanism into the programming language
  - E.g. Java RMI

64

## STUB GENERATION



- `void func(char x; float y; int z[5])`
- 1-byte character transmits with 3-padded bytes
- Float sent as whole word (4-bytes)
  - Array as group of words, proceed by word describing length
  - Client stub must package data in specific format
  - Server stub must receive and unpackage in specific format

- Client and server must agree on representation of simple data structures: int, char, floats w/ little endian
- RPC clients/servers: must agree on protocol
  - TCP? UDP?

65

## STUB GENERATION - 2

- Interfaces are specified using an Interface Definition Language (IDL)

- Interface specifications in IDL are used to generate language specific stubs

- IDL is compiled into client and server-side stubs

- Much of the plumbing for RPC involves maintaining boilerplate-code

66

## LANGUAGE BASED SUPPORT

- Leads to simpler application development

- Helps with providing access transparency
  - Differences in data representation, and how object is accessed
  - Inter-language parameter passing issues resolved:
    → *Just 1 language*

- Well known example: *Java Remote Method Invocation*
  RPC equivalent embedded in Java

67

## RPC VARIATIONS

- RPC: client typically blocks until reply is returned
- Strict blocking *unnecessary* when there is no result

- *Asynchronous RPCs*
  - When no result, server can immediately send reply



Client/server synchronous RPC    Client/server asynchronous RPC

68

## RPC VARIATIONS – 2

- What are tradeoffs for synchronous vs. asynchronous procedure calls?
  - For a local program
  - For a distributed program (system)

- Use cases for asynchronous procedure calls
  - Long running jobs allow client to perform alternate work in background (in parallel)
  - Client may need to make multiple service calls to multiple server backends at the same time…

69

## TYPES OF ASYNCHRONOUS RPC

- **Deferred synchronous RPC**
  - Server performs *CALLBACK* to client
  - Client, upon making call, spawns separate thread which blocks and waits for call



- **One-way RPCs**
  - Client **does not wait** for **any** server acknowledgement – it just goes…
- **Client polling**
  - Client (*using separate thread*) continually polls server for result

70

## MULTICAST RPC

- Send RPC request *simultaneously* to group of servers
- Hide that multiple servers are involved
- Consideration:
  *Does the client need all results or just one?*
- Use cases:
  - Fault tolerance – wait for just one
  - Replicate execution – verify results, *use first result*
  - Divide and conquer - multiple RPC calls work in parallel on different parts of dataset, client aggregates results

71

## RPC EXAMPLE: DISTRIBUTED COMPUTING ENVIRONMENT (DCE)

- **DCE**: basis for Microsoft's distributed computing object model (DCOM)
- Used in Samba, *cross-platform* file and print sharing via RPC
- Middleware system – provides layer of abstraction between OS and distributed applications
- Designed for Unix, ported to **all** major operating systems
- Install DCE middleware on set of heterogeneous machines – distributed applications can then access shared resources to:
  - Mount a windows file system on Linux
  - Share a printer connected to a Windows server
- Uses client/server model
- All communication via RPC
- DCE daemon tracks participating machines, ports

72

## DCE CLIENT-TO-SERVER BINDING



- Server name comes from directory server
- Server port comes from DCE daemon
  - DCE daemon has a well known port # client already knows

73

## EXTRA: DCE – CLIENT/SERVER DEVELOPMENT

1. Create Interface definition language (IDL) files
   - IDL files contain Globally unique identifier (GUID)
   - GUIDs must match: client and server compare GUIDs to verify proper versions of the distributed object
   - 128-bit binary number
2. Next, add names of remote procs and params to IDL
3. Then compile the IDL files
   _Compiler generates:_
   - Header file (interface.h in C)
   - Client stub
   - Server stub

74

## EXTRA: DCE – BINDING CLIENT TO SERVER

- For a client to call a server, server must be registered
  - _Java: uses RMI registry_
- Client process to search for RMI server:
  1. Locate the server's host machine
  2. Locate the server (i.e. process) on the host
- Client must discover the server's RPC port

- **DCE daemon:** maintains table of (server,port) pairs

- When servers boot:
  1. Server asks OS for a port, registers port with DCE daemon
  2. Also, server registers with directory server, separate server that tracks DCE servers

75

## OBJECTIVES – 2/13

- Questions from 2/6
- Midterm Grading In Progress - Targeting Review Thursday
- Assignment 2: Key/Value Store
  - Java Maven project template files posted
- Chapter 3: Processes
  - Chapter 3.4: Servers
  - Chapter 3.5: Resource (Code) Migration (_light-review_)
- Chapter 4: Communication
  - Chapter 4.1: Foundations
  - Chapter 4.2: RPC (light-review)
  - **Chapter 4.3: Message Oriented Communication**

76



Apache ActiveMQ

# CH. 4.3: MESSAGE-ORIENTED COMMUNICATION

77

## MESSAGE ORIENTED COMMUNICATION

- RPC assumes that the _client_ and _server_ are running **at the same time…**  (temporally coupled)
- RPC communication is typically **synchronous**

- When client and server are not running at the same time
- Or when communications should not be **blocked**…

- **This is a use case for message-oriented communication**
  - Synchronous vs. asynchronous
  - Messaging systems
  - Message-queueing systems

78

## SOCKETS

- **Communication end point**
- **Applications can read / write data to**
- **Analogous to file streams for I/O, but _network streams_**

| Operation | Description |
|-----------|-------------|
| socket | Create a new communication end point |
| bind | Attach local address to socket (IP / port) |
| listen | Tell OS what max # of pending connection requests should be |
| accept | Block caller until a connection request arrives |
| connect | Actively attempt to establish a connection |
| send | Send some data over the connection |
| receive | Receive some data over the connection |
| close | Release the connection |

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.79

79

## SOCKETS - 2

- **Servers execute 1st - 4 operations (socket, bind, listen, accept)**
- **Methods refer to C API functions**
- **Mappings across different libraries will vary (_e.g. Java_)**

| Operation | Description |
|-----------|-------------|
| socket | Create a new communication end point |
| bind | Attach local address to socket (IP / port) |
| listen | Tell OS what max # of pending connection requests should be |
| accept | Block caller until a connection request arrives |
| connect | Actively attempt to establish a connection |
| send | Send some data over the connection |
| receive | Receive some data over the connection |
| close | Release the connection |

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.80

80

## SERVER SOCKET OPERATIONS

- **Socket**: creates new communication end point
- **Bind**: associated IP and port with end point
- **Listen**: for connection-oriented communication, non-blocking call reserves buffers for specified number of pending connection requests server is willing to accept
- **Accept**: blocks until connection request arrives
  - Upon arrival, new socket is created matching original
  - Server spawns thread, or forks process to service incoming request
  - Server continues to wait for new connections on original socket

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.81

81

## CLIENT SOCKET OPERATIONS

- **Socket**: Creates socket client uses for communication
- **Connect**: Server transport-level address provided, client blocks until connection established
- **Send**: Supports sending data (to: server/client)
- **Receive**: Supports receiving data (from: server/client)
- **Close**: Closes communication channel
  - Analogous to closing a file stream



February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.82

82

## SOCKET COMMUNICATION

- Sockets provide primitives for implementing your own TCP/UDP communication protocols
- Directly using sockets for transient (non-persisted) messaging is very basic, can be brittle
  - Easy to make mistakes…
- Any extra communication facilities must be implemented by the application developer
- More advanced approaches are desirable
  - E.g. frameworks with support common desirable functionality

February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.83

83

## ZEROMQ – SOCKET LIBRARY

- (0MQ) High performance intelligent **socket library**
- _zero broker, zero latency, zero admin, zero cost, zero waste_
- Provides a message queue
- _Builds upon_ functionality of traditional sockets  **ØMQ**
- Implementation in C++
  - 30+ language bindings provided
- Enables support for various messaging patterns
- Can support brokered (centralized) and broker-less topologies



February 13, 2024 — TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma — L11.84

84

## ZEROMQ – 2

- ZeroMQ is **TCP-connection-oriented communication**

- Provides socket-like primitives with more functionality
  - Basic socket operations *abstracted* away
  - Supports many-to-one, one-to-one, and one-to-many connections
  - *Multicast* connections (one-to-many – single server socket simultaneously "connects" to multiple clients)

- Asynchronous messaging

- Supports pairing sockets to support communication patterns

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.85

85

## ZEROMQ - PATTERNS

- **Request-reply pattern**
  - Traditional client-server communication (e.g. RPC)
  - Client: request socket (**REQ**)
  - Server: reply socket (**REP**)

- **Publish-subscribe pattern**
  - Clients **subscribe** to messages **published** by servers
  - As in event-based coordination (Ch. 1)
  - Supports multicasting messages from server to multiple
  - Client: subscribe socket (**SUB**)
  - Server: publish socket (**PUB**)

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.86

86

## ZEROMQ – PATTERNS - 2

- **Pipeline pattern (FIFO-queue)**
  - Analogous to a producer/consumer bounded buffer
  - Producing processes generate results, push to pipe
  - Consuming processes consume results, pull from pipe
  - Producers: push socket (**PUSH** socket)
  - Consumers: pull socket (**PULL** socket)
  - Push- distributes messages to all pull clients evenly
  - Consumers pull results from pipe and push results downstream

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.87

87

## QUEUEING ALTERNATIVES

- Cloud services
  - Amazon Simple Queueing Service (SQS)
  - Azure service bus

- Open source frameworks
  - Nanomsg
  - ZeroMQ

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.88

88

## MESSAGE PASSING INTERFACE (MPI)

- MPI introduced – version 1.0 March 1994
- Message passing API for parallel programming: *supercomputers*

- Communication protocol for parallel programming for: Supercomputers, High Performance Computing (HPC) clusters

- Point-to-point and collective communication

- Goals: high performance, scalability, portability

- Most implementations in C, C++, Fortran

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.89

89

## MOTIVATIONS FOR MPI

- Motivation: sockets insufficient for interprocess communication on large scale HPC compute clusters and super computers

  - Sockets at the wrong level of abstraction
  - Sockets designed to communicate over the network using general purpose TCP/IP stacks
  - Not designed for proprietary protocols
  - Not designed for high-speed interconnection networks used by supercomputers, HPC-clusters, etc.
  - Better buffering and synchronization needed

February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.90

90

## MOTIVATIONS FOR MPI - 2

- Supercomputers had proprietary communication libraries
  - Offer a wealth of efficient communication operations
- All libraries mutually incompatible
- Led to significant portability problems developing parallel code that could migrate across supercomputers
- Led to development of MPI
  - To support transient (non-persistent) communication for parallel programming

February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.91

91

## MPI FUNCTIONS / DATATYPES

- Very large library, v1.0 (1994) 128 functions
- Version 3 (2015) 440+
- MPI data types:
- Provide common mappings

| MPI datatype | C datatype |
|---|---|
| MPI.CHAR | signed char |
| MPI.SHORT | signed short int |
| MPI.INT | signed int |
| MPI.LONG | signed long int |
| MPI.UNSIGNED_CHAR | unsigned char |
| MPI.UNSIGNED_SHORT | unsigned short int |
| MPI.UNSIGNED | unsigned int |
| MPI.UNSIGNED_LONG | unsigned long int |
| MPI.FLOAT | float |
| MPI.DOUBLE | double |
| MPI.LONG_DOUBLE | long double |
| MPI.BYTE | |
| MPI.PACKED | |

February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.92

92

## COMMON MPI FUNCTIONS

- MPI - no recovery for process crashes, network partitions
- Communication among grouped processes: (groupID, processID)
- IDs used to route messages in place of IP addresses

| Operation | Description |
|---|---|
| MPI_bsend | Append outgoing message to a local send buffer |
| MPI_send | Send message, wait until copied to local/remote buffer |
| MPI_ssend | Send message, wat until transmission starts |
| MPI_sendrecv | Send message, wait for reply |
| MPI_isend | Pass reference to outgoing message and continue |
| MPI_issend | Pass reference to outgoing messages, wait until receipt start |
| MPI_recv | Receive a message, block if there is none |
| MPI_irecv | Check for incoming message, **do not block!** |

February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.93

93

## MESSAGE-ORIENTED-MIDDLEWARE

- **Message-queueing systems**
  - Provide extensive support for *persistent* asynchronous communication
  - In contrast to transient systems
  - Temporally decoupled: messages are eventually delivered to recipient queues
- Message transfers may take minutes vs. sec or ms
- Each application has its own private queue to which other applications can send messages

February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.94

94

## MESSAGE QUEUEING SYSTEMS: USE CASES

- Enables communication between applications, or sets of processes
  - User applications
  - App-to-database
  - To support distributed real-time computations
- Use cases
  - Batch processing, Email, workflow, groupware, routing subqueries

February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.95

95

## MESSAGE QUEUEING SYSTEMS

- Scenarios:
  - (a) Sender/receiver both running
  - (b) Sender running, receiver offline
  - (c) Sender offline, receiver running
  - (d) Sender/receiver both offline
- Queue persists msgs, and attempts to send them but no one may be available to receive them...



February 13, 2024  TCSS558: Applied Distributed Computing [Winter 2024]  School of Engineering and Technology, University of Washington - Tacoma  L11.96

96

## MESSAGE QUEUEING SYSTEMS - 2

- **Key:** Truly persistent messaging
- Message queueing systems can persist messages for awhile and senders and receivers can be offline

- **Messages**
- Contain *any* data, may have size limit
- Are properly addressed, to a destination queue

- **Basic Inteface**
- PUT: called by sender to append msg to specified queue
- GET: blocking call to remove oldest msg from specified queue
  - Blocked if queue is empty
- POLL: Non-blocking, gets msg from specified queue

February 13, 2024 TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.97

97

## MESSAGE QUEUEING SYSTEMS ARCHITECTURE

- **Basic Interface cont'd**
- NOTIFY: install a callback function, for when msg is placed into a queue. Notifies receivers

- **Queue managers**: manage individual message queues as a separate process/library
- Applications get/put messages only from *local* queues
- Queue manager and apps share local network
- **ISSUES:**
- How should we reference the destination queue?
- How should names be resolved (looked-up)?
  - Contact address (host, port) pairs
  - Local look-up tables can be stored at each queue manager

February 13, 2024 TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.98

98

## MESSAGE QUEUEING SYSTEMS ARCHITECTURE - 2

- **ISSUES**:
- How do we route traffic between queue managers?
  - How are name-to-address mappings efficiently kept?
  - Each queue manager should be known to all others

- **Message brokers**
- Handle message conversion among different users/formats
- Addresses cases when senders and receivers don't speak the same protocol (language)
- Need arises for message protocol converters
  - "Reformatter" of messages
- Act as application-level gateway

February 13, 2024 TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.99

99

## MESSAGE BROKER ORGANIZATION



99

## AMQP PROTOCOL

- Message-queueing systems initially developed to enable legacy applications to interoperate
- Decouple inter-application communication to "open" messaging-middleware
- Many are proprietary solutions, *so not very open*
- e.g. Microsoft Message Queueing service, Windows NT 1997
- **Advanced message queueing protocol (AMQP)**, 2006
- Address openness/interoperability of proprietary solutions
- Open wire protocol for messaging with powerful routing capabilities
- Help *abstract* messaging and application interoperability by means of a generic open protocol
- Suffer from incompatibility among protocol versions
- pre-1.0, 1.0+

February 13, 2024 TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.101

101

## AMQP - 2

- Consists of: Applications, Queue managers, Queues

- **Connections:** set up to a queue manager, TCP, with potentially many channels, stable, reused by many channels, long-lived

- **Channels:** support short-lived one-way communication

- **Sessions:** bi-directional communication across two channels

- **Link:** provide fine-grained flow-control of message transfer/status between applications and queue manager

February 13, 2024 TCSS558: Applied Distributed Computing [Winter 2024] School of Engineering and Technology, University of Washington - Tacoma | L11.102

102

## AMQP MESSAGING

- AMQP nodes: producer, consumer, queue
- Producer/consumer: represent regular applications
- Queues: store/forward messages

- Persistent messaging:
- **Messages** can be marked *durable*
- These messages can only be delivered by nodes able to recover in case of failure
- Non-failure resistant nodes must reject durable messages
- **Source/target** nodes can be marked *durable*
- Track what is durable (node state, node+msgs)

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L11.103 |

103

## MESSAGE-ORIENTED-MIDDLEWARE EXAMPLES:

- **Some examples:**
- RabbitMQ, Apache QPid
  - Implement Advanced Message Queueing Protocol (AMQP)
- Apache Kafka
  - **Dumb broker** (message store), similar to a distributed log file
  - **Smart consumers** – intelligence pushed off to the clients
  - Stores stream of records in categories called topics
  - Supports voluminous data, many consumers, with minimal O/H
  - Kafka **does not track** which messages were read by each consumer
  - Messages are removed after timeout
  - Clients must track their own consumption (*Kafka doesn't help*)
  - Messages have key, value, timestamp
  - Supports high volume pub/sub messaging and streams

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L11.104 |

104

## QUESTIONS

| February 13, 2024 | TCSS558: Applied Distributed Computing [Winter 2024]<br>School of Engineering and Technology, University of Washington - Tacoma | L11.10 5 |

105