

TCSS 558: APPLIED DISTRIBUTED COMPUTING

**Ch. 3 - Processes:
Servers**

Ch. 4 - Communication

Wes J. Lloyd
 School of Engineering
 & Technology (SET)
 University of Washington - Tacoma

OBJECTIVES – 2/16

- **Questions from 2/9**
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.2

ONLINE DAILY FEEDBACK SURVEY

- Daily Feedback Quiz in Canvas – Available After Each Class
- Extra credit available for completing surveys **ON TIME**
- Tuesday surveys: due by ~ Wed @ 10p
- Thursday surveys: due ~ Mon @ 10p

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.3

TCSS 558 - Online Daily Feedback Survey - 1/5

Due Jan 6 at 10pm Points 1 Questions 4
 Available Jan 5 at 1:30pm - Jan 6 at 11:59pm 1 day Time Limit None

Question 1 0.5 pts

On a scale of 1 to 10, please classify your perspective on material covered in today's class:

1	2	3	4	5	6	7	8	9	10
Mostly Review To Me			Equal New and Review				Mostly New To Me		

Question 2 0.5 pts

Please rate the pace of today's class:

1	2	3	4	5	6	7	8	9	10
Slow			Just Right				Fast		

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.4

MATERIAL / PACE

- Please classify your perspective on material covered in today's class (22 respondents):
- 1-mostly review, 5-equal new/review, 10-mostly new
- **Average – 6.32** (↑ - previous 6.10)
- Please rate the pace of today's class:
- 1-slow, 5-just right, 10-fast
- **Average – 5.41** (↓ - previous 5.71)

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.5

FEEDBACK FROM 2/9


- **How can you efficiently transfer a very large amount of data (TB or PB) from server to server? Do distributed systems play a role?**
- Number of days to transfer 240 TB via the internet at typical utilizations:

Utilization	Internet Connection Speed			
	1Gbps	500Mbps	300Mbps	150Mbps
25%	95	190	316	632
50%	47	95	158	316
75%	32	63	105	211

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.6

FEEDBACK - 2

- Best case to transfer 240TB is ~32 days @ 1 Gbps & 75% utilization
 - Can't dedicate entire company Internet to data transfer
- To address large scale data transfer issue, AWS offers rental hardware & shipping solution
- AWS snowball
 - Each device transfers 50 TB
 - Rental cost is \$300 per device for 10 days
 - Round trip shipping time is ~1 week
 - 5 devices will transfer 250TB (\$1500)
- See:
 - <http://blog.zorangatic.com/2016/09/aws-snowball.html>
 - <https://aws.amazon.com/snowball>



February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.7
-------------------	--	-------

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review**
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.8
-------------------	--	-------

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance**
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.10
-------------------	--	--------

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.11
-------------------	--	--------

ASSIGNMENT 1

- Extension to Sunday February 21st**
- Discussion Board created on Canvas
 - Answers to common questions posted online
- Team signup posted on Canvas under 'People'
- 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 or 8)
- Client node program interacts with server node to put, get, delete, or list items in a key/value store

February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.12
-------------------	--	--------

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021	TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.13
-------------------	--	--------



CH. 3.4: SERVERS

L11.14

SERVERS

- Cloud & Distributed Systems – rely on **Linux**
- <http://www.zdnet.com/article/it-runs-on-the-cloud-and-the-cloud-runs-on-linux-any-questions/>
- IT is moving to the cloud. And, what powers the cloud?
 - **Linux**
- Uptime Institute survey - 1,000 IT executives (2016)
 - 50% of IT executives – plan to migrate majority of IT workloads to off-premise to cloud or colocation sites
 - 23% expect the shift in 2017, 70% by 2020...
- Docker on Windows / Mac OS X
 - Based on **Linux**
 - Mac: Hyperkit Linux VM
 - Windows: Hyper-V Linux VM

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.15

SERVERS - 2

- Servers implement a specific service for a collection of clients
- Servers wait for incoming requests, and respond accordingly
- **Server types**
- **Iterative:** immediately handle client requests
- **Concurrent:** Pass client request to separate thread
- Multithreaded servers are concurrent servers
 - E.g. Apache Tomcat
- **Alternative:** fork a new process for each incoming request
- **Hybrid:** mix the use of multiple processes with thread pools

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.16

END POINTS

- Clients connect to servers via:
IP Address and **Port Number**
- How do ports get assigned?
 - Many protocols support “default” port numbers
 - Client must find IP address(es) of servers
 - A single server often hosts multiple end points (servers/services)
 - When designing new TCP client/servers must be careful not to repurpose ports already commonly used by others

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.17

COMMON PORTS packetlife.net

TCP/UDP Port Numbers			
7 Echo	554 RTSP	2745 Ssh	6891-6901 Windows Live
19 Chargen	546-547 DHCPv6	2967 Symantec AV	6970 Quicktime
20-21 FTP	560 rmonitor	3050 Interbase DB	7212 GhostSurf
22 SSH/SCP	563 NNTP over SSL	3074 XBOX Live	7648-7649 CU-SeeMe
23 Telnet	587 SMTP	3124 HTTP Proxy	8000 Internet Radio
25 SMTP	591 FileMaker	3127 MyDoom	8080 HTTP Proxy
42 WINS Replication	593 Microsoft DCOM	3128 HTTP Proxy	8086-8087 Kaspersky AV
43 WHOIS	631 Internet Printing	3222 GUPP	8118 Proxy
49 TACACS	636 LDAP over SSL	3260 iSCSI Target	8200 VMware Server
53 DNS	639 MSDP (PM)	3306 MySQL	8500 Adobe ColdFusion
67-68 DHCP/BOOTP	646 LDP (MPLS)	3389 Terminal Server	8767 TeamSpeak
69 FTP	691 MS Exchange	3689 iTunes	8866 Ssh
70 Gopher	860 iSCSI	3690 Subversion	9100 HP JetDirect
79 Finger	873 rsync	3724 World of Warcraft	9101-9103 Bacula
80 HTTP	902 VMware Server	3784-3785 Ventrilo	9119 XMGT
88 Kerberos	989-990 FTP over SSL	4333 mSQL	9800 WebDAV
102 MS Exchange	993 IMAP4 over SSL	4444 Raster	9898 Jabber
110 POP3	995 POP3 over SSL	4664 Google Desktop	9988 RealSociety
113 Ident	1025 Microsoft RPC	4672 Ssh	9999 Urchin
119 NNTP (Usenet)	1026-1029 Windows Messenger	4899 Radmin	10000 Webmin
123 NTP	1080 SOCKS Proxy	5000 Ushp	10000 BackupExec
135 Microsoft RPC	1080 Rdp	5001 Ssh	10113-10116 NetIQ
137-139 NetBIOS	1194 OpenVPN	5001 Iperf	11371 OpenPGP
143 IMAP4	1214 Kazaa	5004-5005 RTP	12035-12036 Second Life
161-162 SNMP	1241 Nessus	5050 Yahoo! Messenger	12345 NetBus
177 XDMCP	1311 Dell OpenManage	5060 SIP	13720-13721 NetBackup
178 RFB	1337 WASTE	5190 XMPP	14567 NetSarang

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.18

TYPES OF SERVERS

- **Daemon server**
 - Example: NTP server
- **Superserver**
- **Stateless server**
 - Example: Apache server
- **Stateful server**
- **Object servers**
- **EJB servers**

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.19

NTP EXAMPLE

- **Daemon servers**
 - Run locally on Linux
 - Track current server end points (outside servers)
 - Example: network time protocol (ntp) daemon
 - Listen locally on specific port (ntp is 123)
 - Daemons routes local client traffic to the configured endpoint servers
 - University of Washington: time.u.washington.edu
 - Example "ntpq -p"
 - Queries local ntp daemon, routes traffic to configured server(s)

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.20

SUPERSERVER

- **Linux inetd / xinetd**
 - Single superserver
 - Extended internet service daemon
 - Not installed by default on Ubuntu
 - Intended for use on server machines
 - Used to configure box as a server for multiple internet services
 - E.g. ftp, pop, telnet
 - inetd daemon responds to multiple endpoints for multiple services
 - Requests fork a process to run required executable program
- **Check what ports you're listening on:**
 - `sudo netstat -tap | grep LISTEN`

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.21

INTERRUPTING A SERVER

- **Server design issue:**
 - Active client/server communication is taking place over a port
 - How can the server / data transfer protocol support interruption?
- Consider transferring a 1 GB image, how do you pass a unrelated message in this stream?
 1. **Out-of-band** data: special messages sent in-stream to support interrupting the server (*TCP urgent data*)
 2. Use a separate connection (different port) for admin control info
- Example: sftp secure file transfer protocol
 - Once a file transfer is started, can't be stopped easily
 - Must kill the client and/or server

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.22

STATELESS SERVERS

- Data about state of clients is not stored
- Example: web application servers are typically stateless
 - Also function-as-a-service (FaaS) platforms
- Many servers maintain information on clients (e.g. log files)
- Loss of stateless data doesn't disrupt server availability
 - Losing log files typically has minimal consequences
- **Soft state:** server maintains state on the client for a limited time (*to support sessions*)
- Soft state information expires and is deleted

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.23

STATEFUL SERVERS

- Maintain persistent information about clients
- Information must be explicitly deleted by the server
- Example:
 - File server - allows clients to keep local file copies for RW
- Server tracks client file permissions and most recent versions
 - Table of (client, file) entries
- If server crashes data must be recovered
- Entire state before a crash must be restored
- Fault tolerance - Ch. 8

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.24

STATEFUL SERVERS - 2

- **Session state**
 - Tracks series of operations by a single user
 - Maintained temporarily, not indefinitely
 - Often retained for multi-tier client server applications
 - Minimal consequence if session state is lost
 - Clients must start over, reinitialize sessions
- **Permanent state**
 - Customer information, software keys
- **Client-side cookies**
 - When servers don't maintain client state, clients can store state locally in "cookies"
 - Cookies are not executable, simply client-side data

February 16, 2021

TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma

L11.25

OBJECT SERVERS

- **OBJECTIVE:** Host objects and enable remote client access
- Do not provide a specific service
 - Do nothing if there are no objects to host
- Support adding/removing hosted objects
- Provide a home where objects live
- Objects, *themselves*, provide "services"
- Object parts
 - State data
 - Code (methods, etc.)
- **Transient object(s)**
 - Objects with limited lifetime (< server)
 - Created at first invocation, destroyed when no longer used (i.e. no clients remain "bound").
 - Disadvantage: initialization may be expensive
 - Alternative: preinitialize and retain objects on server start-up

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.26
-------------------	---	--------

OBJECT SERVERS - 2

- **Should object servers isolate memory for object instances?**
 - Share neither code nor data
 - May be necessary if objects couple data and implementation
- Object server threading designs:
 - Single thread of control for object server
 - One thread for each object
 - Servers use separate thread for client requests
- Threads created on demand **vs.** Server maintains pool of threads
- **What are the tradeoffs for creating server threads on demand vs. using a thread pool?**

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.27
-------------------	---	--------

EJB – ENTERPRISE JAVA BEANS

- EJB- specialized Java object hosted by a EJB web container
- 4 types: stateless, stateful, entity, and message-driven beans
- Provides "middleware" standard (framework) for implementing back-ends of enterprise applications
- EJB web application containers integrate support for:
 - Transaction processing
 - Persistence
 - Concurrency
 - Event-driven programming
 - Asynchronous method invocation
 - Job scheduling
 - Naming and discovery services (JNDI)
 - Interprocess communication
 - Security
 - Software component deployment to an application server

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.28
-------------------	---	--------

APACHE WEB SERVER

- Highly configurable, extensible, platform independent
- Supports TCP HTTP protocol communication
- Uses hooks – placeholders for group of functions
- Requests processed in phases by hooks
- Many hooks:
 - Translate a URL
 - Write info to log
 - Check client ID
 - Check access rights
- Hooks processed in order enforcing flow-of-control
- Functions in replaceable modules

The diagram illustrates the Apache Web Server architecture. It shows a central 'Apache core' box. Above it are several 'Module' boxes, each containing a 'Function' box. 'Hooks' are shown as boxes that 'point to functions in modules'. A 'Request' enters from the bottom left, passes through the 'Apache core', and a 'Response' exits at the bottom right. A 'Link between function and hook' is also indicated.

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.29
-------------------	---	--------

SERVER CLUSTERS

- Hosted across an LAN or WAN
- Collection of interconnected machines
- Can be organized in tiers:
 - Web server → app server → DB server
 - App and DB server sometimes integrated

The diagram shows a three-tier architecture. On the left, 'Client requests' enter a 'Logical switch (possibly multiple)' labeled 'First tier'. This switch dispatches requests to a grid of 'Application/compute servers' labeled 'Second tier'. These servers are connected to a 'Distributed file/database system' labeled 'Third tier'.

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.30
-------------------	---	--------

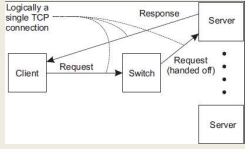
LAN REQUEST DISPATCHING

- Front end of three tier architecture (logical switch) provides distribution transparency – hides multiple servers
- Transport-layer switches: switch accepts TCP connection requests, hands off to a server
 - Example: hardware load balancer (F5 networks – Seattle)
 - HW Load balancer - OSI layers 4-7
- Network-address-translation (NAT) approach:
 - All requests pass through switch
 - Switch sits in the middle of the client/server TCP connection
 - Maps (rewrites) source and destination addresses
- Connection hand-off approach:
 - **TCP Handoff:** switch hands off connection to a selected server

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.31
-------------------	---	--------

LAN REQUEST DISPATCHING - 2

- Who is the best server to handle the request?
- Switch plays important role in distributing requests
- Implements load balancing
- **Round-robin** – routes client requests to servers in a looping fashion
- **Transport-level** – route client requests based on TCP port number
- **Content-aware request distribution** – route requests based on inspecting data payload and determining which server node should process the request



February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.32

WIDE AREA CLUSTERS

- Deployed across the internet
- Leverage resource/infrastructure from Internet Service Providers (ISPs)
- Cloud computing simplifies building WAN clusters
- Resource from a single cloud provider can be combined to form a cluster
- **For deploying a cloud-based cluster (WAN), what are the implications of deploying nodes to:**
 - (1) a single availability zone (e.g. us-east-1e)?
 - (2) across multiple availability zones?

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.33

WAN REQUEST DISPATCHING

- Goal: minimize network latency using WANs (e.g. Internet)
- Send requests to nearby servers
- Request dispatcher: routes requests to nearby server
- **Example:** Domain Name System
 - Hierarchical decentralized naming system
- Linux: find your DNS servers:

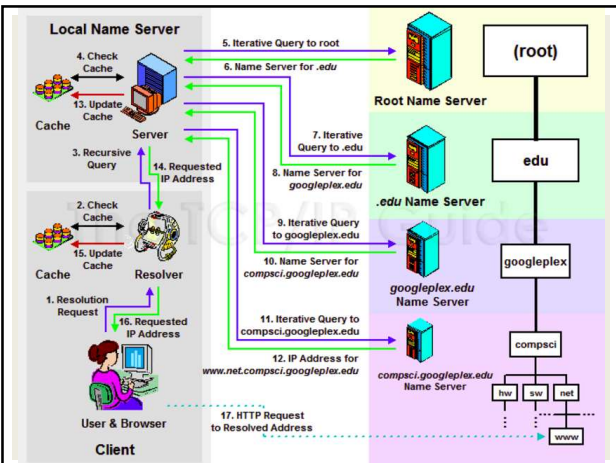

```
# Find you device name of interest
nmcli dev
# Show device configuration
nmcli device show <device name>
```

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.34

DNS LOOKUP

- First query local server(s) for address
- Typically there are (2) local DNS servers
 - One is backup
- Hostname may be cached at local DNS server
 - E.g. www.google.com
- If not found, local DNS server routes to other servers
- Routing based on components of the hostname
- DNS servers down the chain mask the client IP, and use the originating DNS server IP to identify a local host
- **Weakness:** client may be far from DNS server used. Resolved hostname is close to DNS server, but not necessarily close to the client

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.35



DNS: LINUX COMMANDS

- `nslookup <ip addr / hostname>`
- Name server lookup – translates hostname or IP to the inverse
- `traceroute <ip addr / hostname>`
- Traces network path to destination
- By default, output is limited to 30 hops, can be increased

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.37

DNS EXAMPLE – WAN DISPATCHING

- Ping www.google.com in WA from wireless network:
 - nslookup: 6 alternate addresses returned, choose (74.125.28.147)
 - Ping 74.125.28.147: Average RTT = **22.458 ms (11 attempts, 22 hops)**
- Ping www.google.com in VA (us-east-1) from EC2 instance:
 - nslookup: 1 address returned, choose 172.217.9.196
 - Ping 172.217.9.196: Average RTT = 1.278 ms (11 attempts, 13 hops)
- From VA EC2 instance, ping WA [www.google](http://www.google.com) server
 - Ping 74.125.28.147: Average RTT 62.349ms (11 attempts, 27 hops)
 - Pinging the WA-local server is ~60x slower from VA
- From local wireless network, ping VA us-east-1 google :
 - Ping 172.217.9.196: Average RTT=81.637ms (11 attempts, 15 hops)

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.38

DNS EXAMPLE – WAN DISPATCHING

- Ping www.google.com in WA from wireless network:
 - nslookup: 6 alternate addresses returned, choose (74.125.28.147)

Latency to ping VA server in WA: ~3.63x
 VA client: local-google 22.458ms to VA-google 81.637ms

Latency to ping WA server in VA: ~48.7x
 VA client: local-google 1.278ms to WA-google 62.349!

- From local wireless network, ping VA us-east-1 google :
- Ping 172.217.9.196: Average RTT=81.637ms (11 attempts, 15 hops)

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.39

CH 3.2 - EXAMPLE: PLANETLAB

- Unstructured heterogeneous cluster of servers
- Similar to grid but organized as cluster (no grid middleware)
- Testbed established in 2002 for computer networking and distributed systems research
- Organizations share nodes in the cluster

Leverages Linux Vservers
 Early "containers" similar to Docker

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.40

PLANETLAB - 2

- **Slices:** set of Vservers running across PlanetLab
- Acts as a virtual server cluster (similar to Amazon VPC)
- **Node manager:** manages Vservers running on a host
- **Slice creation service (SCS):** To create virtual server clusters
- Clients must be **slice authorities** to create cluster
- **Rspec:** resource specification
 - Specifies resource requirements for a slice
- **Rcap:** resource capability
 - Specifies resource capabilities of nodes

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.41

VSERVERS

- Early container based approach
- Vservers share a single operating system kernel
- Primary task is to support a group of processes
- Provides separation of name spaces
- Linux kernel maps process IDs: host OS → Vservers
- Each Vserver has its own set of libraries and file system
- Similar name separation as the "chroot" command
- Additional isolation provided to prevent unauthorized access among Vservers directory trees

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.42

VSERVERS - 2

- **Advantages of Vservers (containers) vs. VMs:**
- Simpler resource allocation
- Possible to overbook resources by leveraging dynamic resource allocation - **Example: CPU or RAM** ([assignment 0, config 2](#))
- VMs reserve a block of memory
- Containers can oversubscribe memory
 - Memory not formally reserved
 - Linux kernel shares memory among processes
 - Swap filesystem can use disk as extended RAM
- Memory sharing important for PlanetLab
 - Early nodes had limited memory (e.g. 4 GB)
- Vserver hogging most memory reset when out of swap space

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.43

WE WILL RETURN AT
2:44PM




OBJECTIVES – 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.45

CH. 3.5: RESOURCE (CODE) MIGRATION



L11.46

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.47

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.48

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.49

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.50

MOTIVATIONS FOR MIGRATION



- Can migrate **processes** or entire **virtual machines**
- **Goals:**
 - Off-loading machines: reduce load on oversubscribed servers
 - Loading machine: ensure machine has enough work to do
 - 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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.51

LINUX CRUI

- 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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.52

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 than CPU capacity

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.53

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.54

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.55

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

The diagram shows a Client and a Server connected to a Code repository. Step 1: Client fetches code from the repository. Step 2: Client and server communicate.

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.56

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.57

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.58

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**

The diagram shows four scenarios: CS (Client-Server), REV (Remote Evaluation), CoD (Code-on-demand), and MA (Mobile agents). Each scenario shows the state of code, execution (exec), and resources on both Client and Server before and after execution. Asterisks indicate what is modified. CS: everything runs remotely. REV: client provides code for remote exec. CoD: client obtains & runs code. MA: client moves code and exec to server.

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.59

MIGRATION OF HETEROGENEOUS SYSTEMS

- Assumption: code will always work at new node
- Invalid if node architecture is different (*heterogeneous*)

- 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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.60

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.61

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

L11.62

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - Chapter 4.3: Message Oriented Communication

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.63



CH. 4 COMMUNICATION

L11.64

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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.65



CH. 4.1: FOUNDATIONS

L11.66

LAYERED PROTOCOLS

- Distributed systems lack shared memory
- All distributed system communication is based on sending and receiving low-level messages
 - P → Q
- **Open Systems Interconnection 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 16, 2021 TCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma L11.67

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 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.68

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

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.69

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)

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.70

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*)

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.71

OSI MODEL - 4

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

Offset		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

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.72

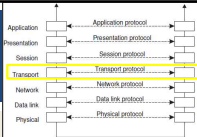
OSI MODEL - 5

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

- Breaks application layer protocol messages into pieces to transmit
- Assigns messages sequence numbers
- Sends all messages

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.73

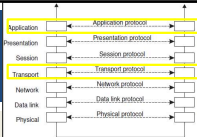
OSI MODEL - 6



- 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)

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.74

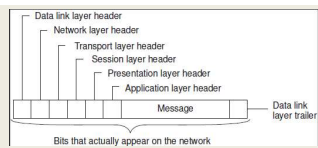
OSI MODEL - 7



- 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.

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.75

OSI MODEL - 8

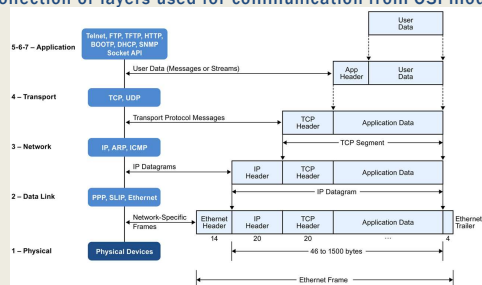


- Each OSI layer contributes overhead bits to the message
- Layers append data to front (and maybe end) of the message
- Receiver strips off headers as the message goes up the OSI model stack:
physical → data-link → network → transport → application

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.76

PROTOCOL STACK

- Collection of layers used for communication from OSI model



February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.77

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”

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.78

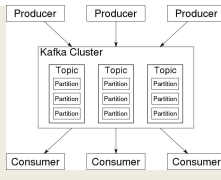
MIDDLEWARE PROTOCOLS - 2

- **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

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.79

MIDDLEWARE PROTOCOLS - 3

- **Message queuing 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



February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.80


MIDDLEWARE PROTOCOLS - 3

- **Message queuing services**

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

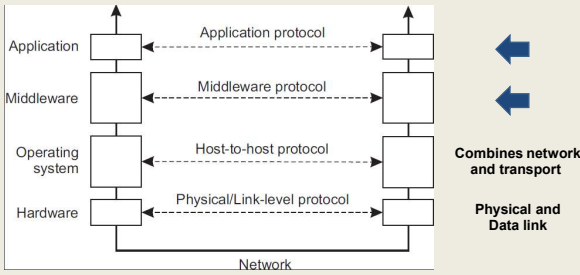
Middleware functions are general, application-independent in nature

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



February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.81

ADAPTED REFERENCE MODEL



February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.82

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.83

TYPES OF COMMUNICATION - 2

- **Asynchronous communication**
 - Client does not block, continues doing other work
- **Synchronous communication**
 - Client blocks and waits
- **Three types of blocking**
 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-queuing systems
 - Block until message delivered to queue
- **Consider each type of blocking (1, 2, 3). Are these modes connectionless (UDP)? connection-oriented (TCP)?**

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.84

OBJECTIVES - 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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**
 - Chapter 4.3: Message Oriented Communication

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.85

CH. 4.2: RPC

L11.86

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 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.87

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 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.88

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 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.89

RPC STEPS

- Client procedure calls client stub
- Client stub builds message and calls OS
- Client's OS send message to remote OS
- Server OS gives message to server stub
- Server stub unpacks parameters, calls server
- Server performs work, returns results to server-side stub
- Server stub packs results in messages, calls server OS
- Server OS sends message to client's OS
- Client's OS delivers message to client stub
- Client stub unpacks result, returns to client

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.90

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)

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.91

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

BIG-ENDIAN									
...	00	01	02	03	04	05	06	07	...
	a	a+1	a+2	a+3	a+4	a+5	a+6	a+7	
LITTLE-ENDIAN									
...	07	06	05	04	03	02	01	00	...
	a	a+1	a+2	a+3	a+4	a+5	a+6	a+7	

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.92

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


February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.93

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

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.94

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?

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.95

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

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.96

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

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.97

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

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.98

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...

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.99

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

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.100

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

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.101

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

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.102

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

February 16, 2021 TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma L11.103

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

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.104

EXTRA: DCE – BINDING CLIENT TO SERVER

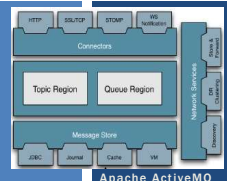
- For a client to call a server, server must be registered
 - *Java: uses RMI registry*
- Client process to search for RMI server:
 1. Locate the server's host machine
 2. Locate the server (i.e. process) on the host
- Client must discover the server's RPC port
- **DCE daemon:** maintains table of (server,port) 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

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.105

OBJECTIVES – 2/16

- Questions from 2/9
- Midterm Review
- Assignment 0: Load Balancing & Performance
- Assignment 1: 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
 - **Chapter 4.3: Message Oriented Communication**

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.106



Apache ActiveMQ

CH. 4.3: MESSAGE-ORIENTED COMMUNICATION

L11.10
7

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

February 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.108

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 16, 2021 TCCS558: Applied Distributed Computing [Winter 2021]
 School of Engineering and Technology, University of Washington - Tacoma L11.109

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.110

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.111

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.112

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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.113

ZEROMQ – SOCKET LIBRARY

- (OMQ) High performance intelligent **socket library**
- **zero broker, zero latency, zero admin, zero cost, zero waste**
- Provides a message queue
- **Buils upon** functionality of traditional sockets
- Implementation in C++
 - 30+ language bindings provided
- Enables support for various messaging patterns
- Can support brokered (centralized) and broker-less topologies

February 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.114

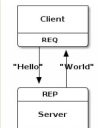
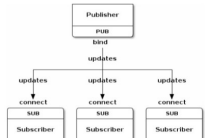
ZEROMQ – 2

- 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 16, 2021
TCCS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.115

ZEROMQ - PATTERNS

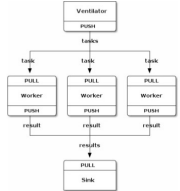
- **Request-reply pattern**
 - Traditional client-server communication (e.g. RPC)
 - Client: request socket (**REQ**)
 - Server: reply socket (**REP**)
- **Publish-subscribe pattern**
 - Clients **subscribe** to messages **published** by servers
 - As in event-based coordination (Ch. 1)
 - Supports multicasting messages from server to multiple
 - Client: subscribe socket (**SUB**)
 - Server: publish socket (**PUB**)

February 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.116

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 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.117

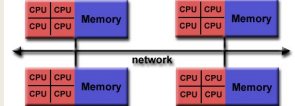
QUEUEING ALTERNATIVES

- **Cloud services**
 - Amazon Simple Queueing Service (SQS)
 - Azure service bus
- **Open source frameworks**
 - Nanomsg
 - ZeroMQ

February 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.118

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 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.119

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 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.120

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 16, 2021 | TCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma | L11.121

MPI FUNCTIONS / DATATYPES

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

C datatype	MPI datatype
signed char	MPI_CHAR
signed short int	MPI_SHORT
signed int	MPI_INT
signed long int	MPI_LONG
unsigned char	MPI_UNSIGNED_CHAR
unsigned short int	MPI_UNSIGNED_SHORT
unsigned int	MPI_UNSIGNED
unsigned long int	MPI_UNSIGNED_LONG
float	MPI_FLOAT
double	MPI_DOUBLE
long double	MPI_LONG_DOUBLE
	MPI_BYTE
	MPI_PACKED

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.122

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, wait until transmission starts
MPI_sendrecv	Send message, wait for reply
MPI_issend	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 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.123

MESSAGE-ORIENTED-MIDDLEWARE

- Message-queuing 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 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.124

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 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.125

MESSAGE QUEUEING SYSTEMS

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

February 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.126

MESSAGE QUEUEING SYSTEMS - 2

- Key:** Truly persistent messaging
- Message queuing 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 Interface**
 - 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 16, 2021
TCSS558: Applied Distributed Computing [Winter 2021]
School of Engineering and Technology, University of Washington - Tacoma
L11.127

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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.128
-------------------	---	---------

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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.129
-------------------	---	---------

MESSAGE BROKER ORGANIZATION

The diagram illustrates the Message Broker Organization. It shows three main components: Source, Message broker, and Destination, each running on its own Local OS. The Source contains an Application and an Interface. The Message broker contains Broker plugins, Rules, and a Queuing layer. The Destination contains an Application and an Interface. Arrows indicate the flow of messages from the Source Application through the Message broker's Queuing layer to the Destination Application. Labels at the bottom indicate 'Plugins to convert messages between APPs' and 'Application-level Queues'.

February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.130
-------------------	---	---------

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 queuing 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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.131
-------------------	---	---------

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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.132
-------------------	---	---------

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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.133
-------------------	---	---------

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 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.134
-------------------	---	---------

QUESTIONS



February 16, 2021	TCCS558: Applied Distributed Computing [Winter 2021] School of Engineering and Technology, University of Washington - Tacoma	L11.135
-------------------	---	---------