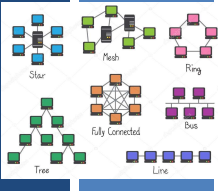



TCCS 558: APPLIED DISTRIBUTED COMPUTING

Communication

Wes J. Lloyd
 Institute of Technology
 University of Washington - Tacoma

OBJECTIVES

- Assignment #1 Questions
- Assignment #2
- Ch. 4 - Communications
 - Message-oriented communication:
 - Zeromq, MPI,
 - Message Queueing Systems
 - Multicast communication

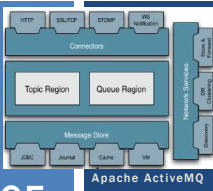
November 14, 2017
TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma
L13.2

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

November 14, 2017
TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma
L13.3

CH. 4.3: MESSAGE-ORIENTED COMMUNICATION



Apache ActiveMQ

L13.4

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

November 14, 2017
TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma
L13.5

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

November 14, 2017
TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma
L13.6

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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.7

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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.8

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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.9

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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.10

ZEROMQ

- (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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.11

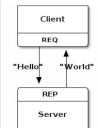
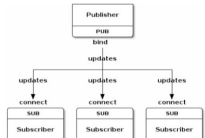
ZEROMQ - 2

- 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

November 14, 2017 TCS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma L13.12

ZEROMQ - PATTERNS

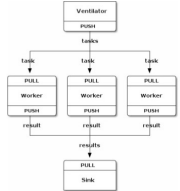
- **Request-reply pattern**
 - Traditional client-server communication (e.g. RPC)
 - Client: request socket (**REQ**)
 - Server: reply socket (**REP**)
- **Publish-subscribe pattern**
 - Clients **subscribe** to messages **published** by servers
 - As in event-based coordination (Ch. 1)
 - Supports multicasting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

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.13

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



November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.14

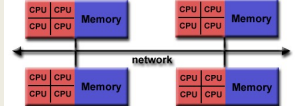
QUEUEING ALTERNATIVES

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

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.15

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



November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.16

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



November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.17

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

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.18

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

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.25

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

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.26

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 and Destination each contain an Application and an Interface. The Message broker component contains Broker plugins, Rules, and a Queuing layer. Arrows indicate the flow of messages from the Source Application through the Interface to the Message broker, and then from the Message broker through the Interface to the Destination Application. A note below the diagram states: 'Plugins to convert messages between APPs' and 'Application-level Queues'.

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.27

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+

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.28

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

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.29

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)

November 14, 2017
TCSS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma
L13.30

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

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.31



Multicast

one to many
x = subscriber

CH. 4.4: MULTICAST COMMUNICATION

Apache ActiveMQ

L13.32

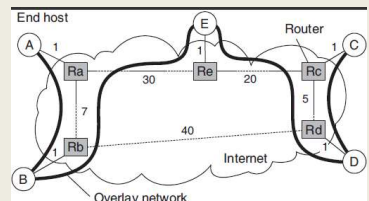
MULTICAST COMMUNICATION

- Sending data to multiple receivers
- Many **failed** proposals for network-level / transport-level protocols to support multicast communication
- **Problem:** How to set up communication paths for information dissemination?
- **Solutions:** require huge management effort, human invention
- Focus shifted more recently to **peer-to-peer** networks
 - Structured overlay networks can be setup easily and provide efficient communication paths
 - Application-level multicasting techniques more successful
 - Gossip-based dissemination: unstructured p2p networks

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.33

NETWORK STRUCTURE

- **Overlay network**
 - Virtual network implemented on top of an actual physical network
- **Underlying network**
 - The actual physical network that implements the overlay



November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.34

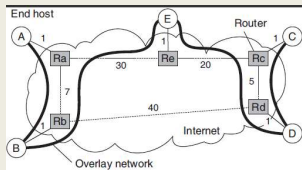
APPLICATION LEVEL TREE-BASED MULTICASTING

- Application level multi-casting
 - Nodes organize into an overlay network
 - Network routers not involved in group membership
 - Group membership is managed at the application level (A2)
- **Downside:**
 - Application-level routing likely less efficient than network-level
 - Necessary tradeoff until having better multicasting protocols at lower layers
- **Overlay topologies**
 - **TREE:** top-down, unique paths between nodes
 - **MESH:** nodes have multiple neighbors; multiple paths between nodes

November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.35

MULTICAST TREE METRICS

- Measure quality of application-level multicast tree
- **Link stress:** is defined per link, counts how often a packet crosses same link (**ideally not more than 1**)
- **Stretch:** ratio in delay between two nodes in the **overlay** vs. the **underlying** networks



November 14, 2017 TCCS558: Applied Distributed Computing [Fall 2017]
 Institute of Technology, University of Washington - Tacoma L13.36

MULTICAST TREE METRICS - 2

- **Stretch (Relative Delay Penalty RDP)** for B to C routes:
- **Overlay:** $B \rightarrow Rb \rightarrow Ra \rightarrow Re \rightarrow E \rightarrow Re \rightarrow Rc \rightarrow Rd \rightarrow D \rightarrow Rd \rightarrow Rc \rightarrow C$
= 73
- **Underlying:** $B \rightarrow Rb \rightarrow Rd \rightarrow Rc \rightarrow C$ = 47
- $73 / 47 = 1.55$
- **Tree cost:** Overall cost of the overlay network
- Ideally would like to minimize network costs
- Find a minimal spanning tree which minimizes total time for disseminating information

November 14, 2017

TCCS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma

L13.37

QUESTIONS



November 14, 2017

TCCS558: Applied Distributed Computing [Fall 2017]
Institute of Technology, University of Washington - Tacoma

L13.38

EXTRA SLIDES



39