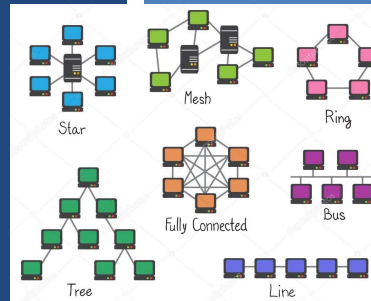# TCSS 558:
# APPLIED DISTRIBUTED COMPUTING

## 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**

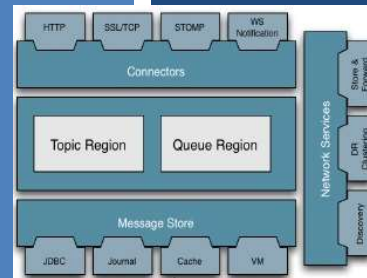| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.2 |
|---|---|---|

# 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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.3 |



**Apache ActiveMQ**

# CH. 4.3: MESSAGE-ORIENTED COMMUNICATION

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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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 | TCSS558: 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
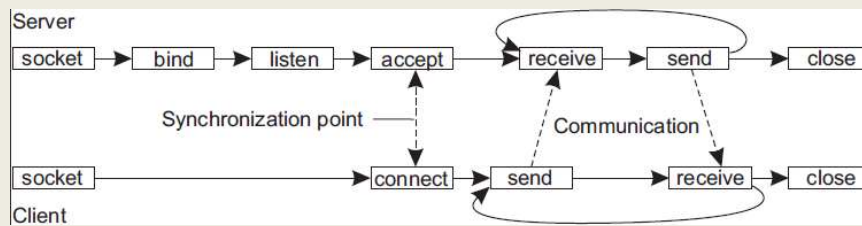
| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma | L13.8 |
|---|---|---|

# 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

# 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
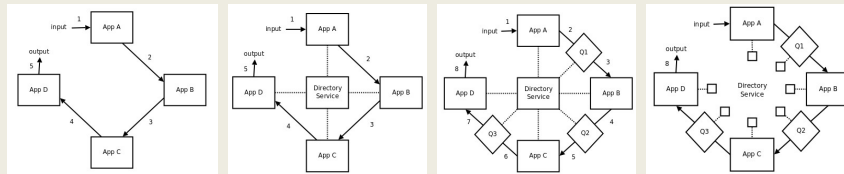  - E.g. frameworks with support common desirable functionality

## ZEROMQ

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

| November 14, 2017 | TCSS558: 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

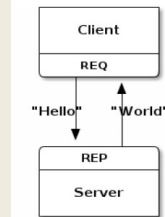- Supports pairing sockets to support communication patterns

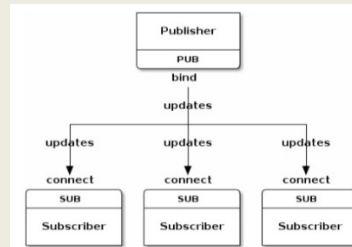| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma | L13.12 |

## ZEROMQ - PATTERNS

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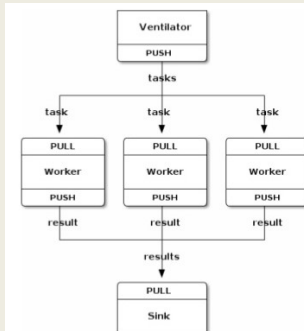
- **Publish-subscribe pattern**
  - Clients **subscribe** to messages **published** by servers
  - As in event-based coordination (Ch. 1)
  - Supports multicasting messages from server to multiple
  - Client: subscribe socket (**SUB**)
  - Server: publish socket (**PUB**)

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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**
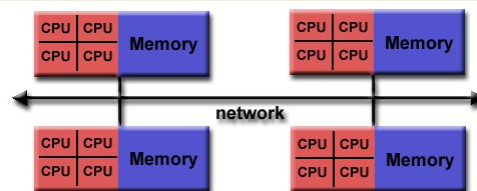  - **Consumers pull results from pipe and push results downstream**

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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**
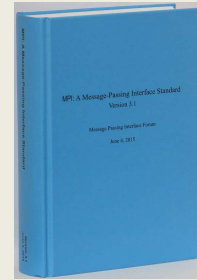


| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.16 |
|---|---|---|

# 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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.18 |
|---|---|---|

# MPI FUNCTIONS / DATATYPES

- **Very large library, v1.0 (1994) 128 functions**

- **Version 3 (2015) 440+**

- **MPI data types:**
- **Provide common mappings**

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

| | | | |
|---|---|---|---|
| MPI_ABORT | MPI_ADDRESS | MPI_ALLGATHER | MPI_ALLGATHERV |
| MPI_ALLREDUCE | MPI_ALLTOALL | MPI_ALLTOALLV | MPI_ATTR_DELETE |
| MPI_ATTR_GET | MPI_ATTR_PUT | MPI_BARRIER | MPI_BCAST |
| MPI_BSEND | MPI_BSEND_INIT | MPI_BUFFER_ATTACH | MPI_BUFFER_DETACH |
| MPI_CANCEL | MPI_CARTDIM_GET | MPI_CART_COORDS | MPI_CART_CREATE |
| MPI_CART_GET | MPI_CART_MAP | MPI_CART_RANK | MPI_CART_SHIFT |
| MPI_CART_SUB | MPI_COMM_COMPARE | MPI_COMM_CREATE | MPI_COMM_DUP |
| MPI_COMM_FREE | MPI_COMM_GROUP | MPI_COMM_RANK | MPI_COMM_REMOTE_GROUP |
| MPI_COMM_REMOTE_SIZE | MPI_COMM_SIZE | MPI_COMM_SPLIT | MPI_COMM_TEST_INTER |
| MPI_DIMS_CREATE | MPI_ERRHANDLER_CREATE | MPI_ERRHANDLER_FREE | MPI_ERRHANDLER_GET |
| MPI_ERRHANDLER_SET | MPI_ERROR_CLASS | MPI_ERROR_STRING | MPI_FINALIZE |
| MPI_GATHER | MPI_GATHERV | MPI_GET_COUNT | MPI_GET_ELEMENTS |
| MPI_GET_PROCESSOR_NAME | MPI_GRAPHDIMS_GET | MPI_GRAPH_CREATE | MPI_GRAPH_GET |
| MPI_GRAPH_MAP | MPI_GRAPH_NEIGHBORS | MPI_GRAPH_NEIGHBORS_COUNT | MPI_GROUP_COMPARE |
| MPI_GROUP_DIFFERENCE | MPI_GROUP_EXCL | MPI_GROUP_FREE | MPI_GROUP_INCL |
| MPI_GROUP_INTERSECTION | MPI_GROUP_RANGE_EXCL | MPI_GROUP_RANGE_INCL | MPI_GROUP_RANK |
| MPI_GROUP_SIZE | MPI_GROUP_TRANSLATE_RANKS | MPI_GROUP_UNION | MPI_IBSEND |
| MPI_INIT | MPI_INITIALIZED | MPI_INTERCOMM_CREATE | MPI_INTERCOMM_MERGE |
| MPI_IPROBE | MPI_IRECV | MPI_IRSEND | MPI_ISEND |
| MPI_ISSEND | MPI_KEYVAL_CREATE | MPI_KEYVAL_FREE | MPI_OP_CREATE |
| MPI_OP_FREE | MPI_PACK | MPI_PACK_SIZE | MPI_PCONTROL |
| MPI_PROBE | MPI_RECV | MPI_RECV_INIT | MPI_REDUCE |
| MPI_REDUCE_SCATTER | MPI_REQUEST_FREE | MPI_RSEND | MPI_RSEND_INIT |
| MPI_SCAN | MPI_SCATTER | MPI_SCATTERV | MPI_SEND |
| MPI_SENDRECV | MPI_SENDRECV_REPLACE | MPI_SEND_INIT | MPI_SSEND |
| MPI_SSEND_INIT | MPI_START | MPI_STARTALL | MPI_TEST |
| MPI_TESTALL | MPI_TESTANY | MPI_TESTSOME | MPI_TEST_CANCELLED |
| MPI_TOPO_TEST | MPI_TYPE_COMMIT | MPI_TYPE_CONTIGUOUS | MPI_TYPE_EXTENT |
| MPI_TYPE_FREE | MPI_TYPE_HINDEXED | MPI_TYPE_HVECTOR | MPI_TYPE_INDEXED |
| MPI_TYPE_LB | MPI_TYPE_SIZE | MPI_TYPE_STRUCT | MPI_TYPE_UB |
| MPI_TYPE_VECTOR | MPI_UNPACK | MPI_WAIT | MPI_WAITALL |
| MPI_WAITANY | MPI_WAITSOME | MPI_WTICK | MPI_WTIME |

| | | |
|---|---|---|
| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.19 |

# COMMON MPI FUNCTIONS

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

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

| | | |
|---|---|---|
| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.20 |

# MESSAGE-ORIENTED-MIDDLEWARE

- **Message-queueing systems**
  - Provide extensive support for *persistent* asynchronous communication
  - In contrast to transient systems
  - Temporally decoupled: messages are eventually delivered to recipient queues

- Message transfers may take minutes vs. sec or ms

- Each application has its own private queue to which other applications can send messages

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.21 |
|---|---|---|

# 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

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.22 |
|---|---|---|

# MESSAGE QUEUEING SYSTEMS

- **Scenarios:**
  (a) Sender/receiver both running

  (b) Sender running, receiver offline

  (c) Sender offline, receiver running

  (d) Sender/receiver both offline

- **Queue persists msgs, and attempts to send them but no one may be available to receive them...**



| Sender running | Sender running | Sender passive | Sender passive |
|---|---|---|---|
| *SENDS* | | | |
| *READS* | | | |
| Receiver running | Receiver passive | Receiver running | Receiver passive |
| (a) | (b) | (c) | (d) |

# MESSAGE QUEUEING SYSTEMS - 2

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

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

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

## 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]<br>Institute of Technology, University of Washington - Tacoma | L13.25 |
|---|---|---|

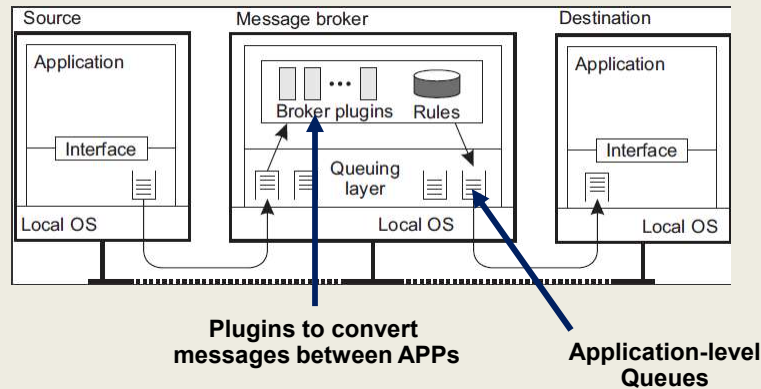## MESSAGE QUEUEING SYSTEMS ARCHITECTURE - 2

- **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]<br>Institute of Technology, University of Washington - Tacoma | L13.26 |
|---|---|---|

# MESSAGE BROKER ORGANIZATION



**Plugins to convert
messages between APPs**

**Application-level
Queues**

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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 queueing protocol (AMQP), 2006**
- **Address openness/interoperability of proprietary solutions**
- **Open wire protocol for messaging with powerful routing capabilities**
- **Help *abstract* messaging and application interoperability by means of a generic open protocol**
- **Suffer from incompatibility among protocol versions**
- **pre-1.0, 1.0+**

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

# AMQP - 2

- Consists of: Applications, Queue managers, Queues

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

- <u>Channels:</u> support short-lived one-way communication

- <u>Sessions:</u> bi-directional communication across two channels

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

| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017]<br>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:
- <u>Messages</u> 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
- <u>Source/target</u> nodes can be marked *durable*
- Track what is durable (node state, node+msgs)

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

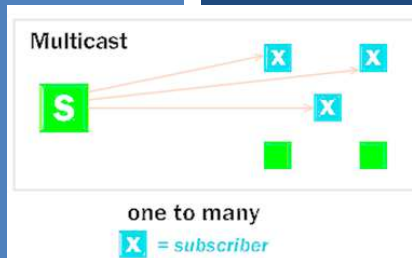## MESSAGE-ORIENTED-MIDDLEWARE EXAMPLES:

- **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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.31 |

---

Multicast

one to many

x = subscriber

Apache ActiveMQ

# CH. 4.4: MULTICAST COMMUNICATION

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
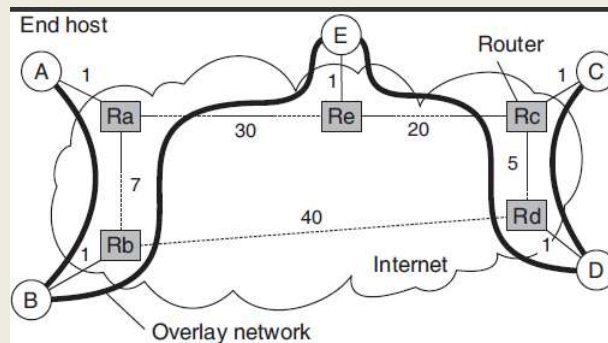  - Gossip-based dissemination: unstructured p2p networks

| November 14, 2017 | TCSS558: 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 | TCSS558: 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**

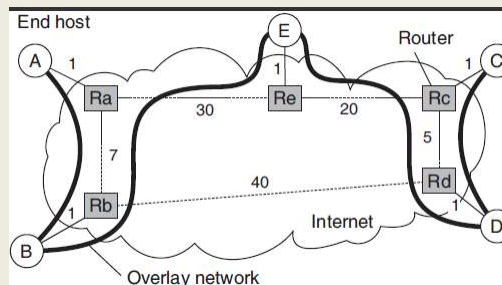| November 14, 2017 | TCSS558: Applied Distributed Computing [Fall 2017] Institute of Technology, University of Washington - Tacoma | L13.35 |
|---|---|---|

---

# 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 <u>overlay</u> vs. the <u>underlying</u> networks**



| November 14, 2017 | TCSS558: 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→Rb→Ra→Re→E→Re→Rc→Rd→D→Rd→Rc→ C = 73
- *Underlying:* B→Rb→Rd→Rc→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 | TCSS558: Applied Distributed Computing [Fall 2017]<br>Institute of Technology, University of Washington - Tacoma | L13.37 |
|---|---|---|

## QUESTIONS

# EXTRA SLIDES