

LECTURE 6: MESSAGE-ORIENTED COMMUNICATION II: MESSAGING IN DISTRIBUTED SYSTEMS

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

1

Lecture Contents

- Middleware in Distributed Systems
- Types of Distributed Communications
 - Remote Procedure Call (RPC):
 - Parameter passing, Example: DCE
 - Registration & Discovery in DCE
 - Message Queuing Systems:
 - Basic Architecture, Role of Message Brokers
 - Example: IBM Websphere
 - Advanced Message Queuing Protocol (AMQP)
 - Example: Rabbit MQ
 - Multicast Communications:
 - Application Layer Messaging
 - Epidemic Protocols

Lecture 6: Messaging on Distributed Systems

CA4006 Lecture Notes (Martin Crane 2017)

2

SECTION 6.1: MIDDLEWARE IN DISTRIBUTED SYSTEMS

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

3

Role of Middleware

- *Observation*
 - Role to provide common services/protocols in Distributed Systems
 - Can be used by many different distributed applications
- *Middleware Functionality*
 - (Un)marshalling of data: necessary for integrated systems
 - Naming protocols: to allow easy sharing, discovery of resources
 - Security protocols: for secure communication
 - Scaling mechanisms, such as for replication & caching (e.g. decisions on where to cache etc.)
 - A rich set of comms protocols: to allow applications to transparently interact with other processes regardless of location.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

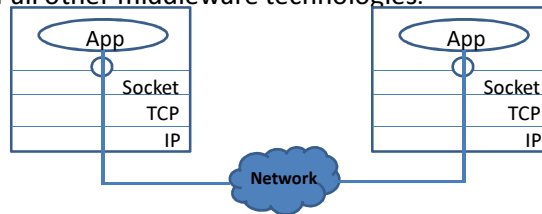
4

Classification of Middleware

- Classify middleware technologies into the following groups:

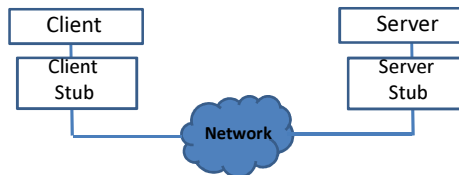
1. *Bog-standard Sockets*

- The basis of all other middleware technologies.



2. *RPC – Remote Procedure Call (more later)*

- RPCs provide a simple way to distribute application logic on separate hosts



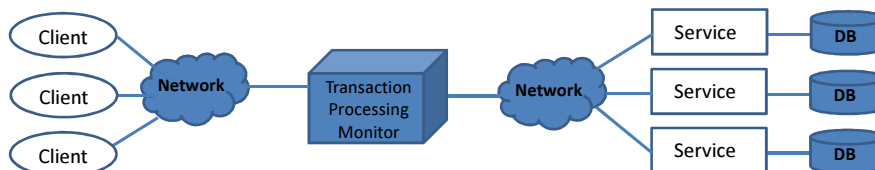
Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

5

Classification of Middleware (/2)

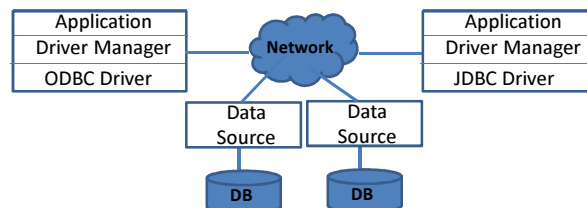
3. *TPM - Transaction Processing Monitors:*

- TPMs are a special form of MW targeted at distributed transactions.



4. *DAM - Database Access Middleware:*

- DBs can be used to share & communicate data between distributed applications.



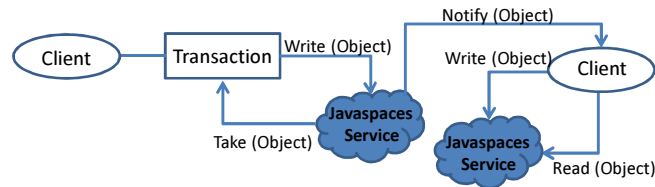
Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

6

Classification of Middleware (/3)

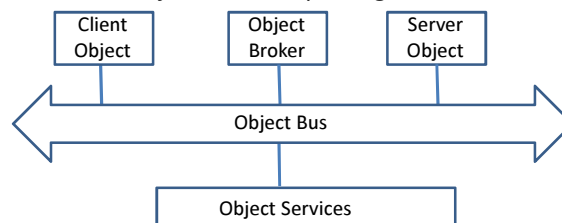
5. *Distributed Tuple:*

- Distributed tuple spaces implement a distributed shared memory space.



6. *DOT (Dist Object Technology) / OOM (Object-Oriented M/w):*

- DOT extends the object-oriented paradigm to distributed applications.



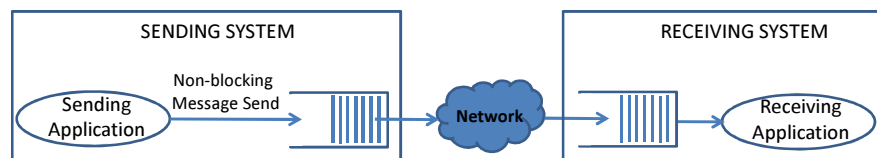
Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

7

Classification of Middleware (/4)

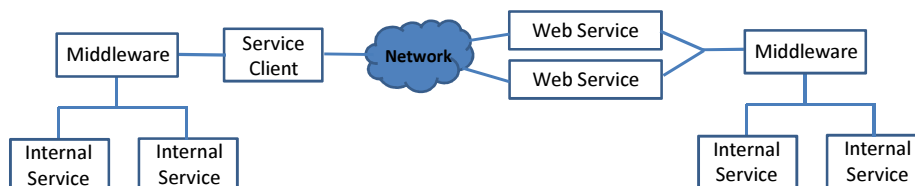
7. *MOM (Message Oriented Middleware):*

- In MOM, messages are exchanged asynchronously between distributed applications (senders and receivers).



8. *Web services:*

- Web services expose services (functionality) on a defined interface, typically accessible through the web protocol HTTP.



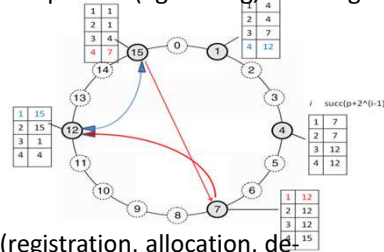
Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

8

Classification of Middleware (/5)

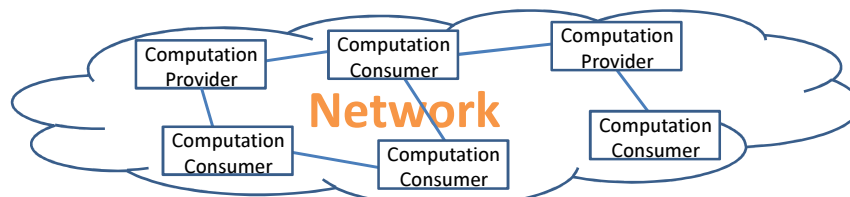
9. Peer-to-peer middleware:

- Have seen above how MW often follows particular *architectural style*.
- In P2P, each peer has equal role in comms pattern (eg routing, node mgmt)
- More on this later...



10. Grid middleware:

- Provides computation power services (registration, allocation, de-allocation) to consumers.

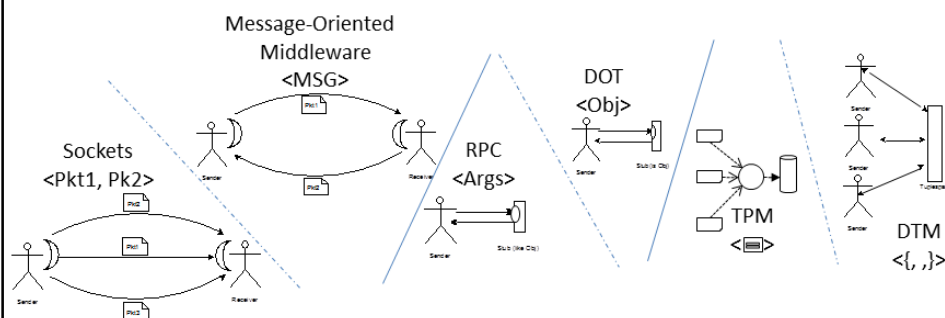


Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

9

Summary of Communications Middleware

- Essentially a range of types of communications middleware
- All can be used to implement others, all are suited to different cases
 - All carry some payload from one side to another <with details>
 - Some of these payloads are 'active' and some are 'passive'
 - Also differ in granularities and whether synchronous or not.



Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

10

SECTION 6.2: COMMUNICATION IN DISTRIBUTED SYSTEMS

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

11

Terminology for Distributed Communications

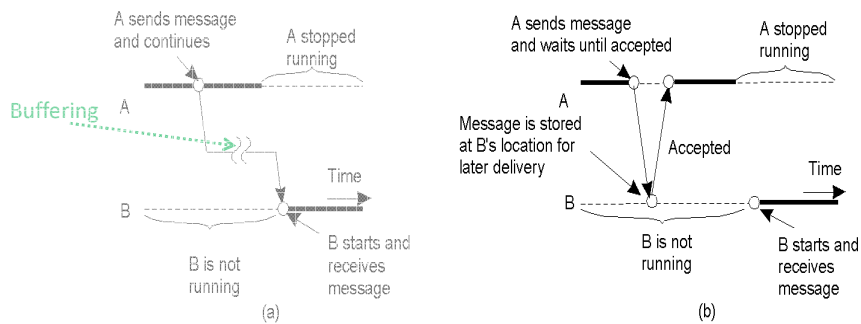
- *Terminology for Distributed Communications*
 - *Persistent Communications:*
 - Once sent, the “sender” stops executing.
 - “Receiver” need not be in operation – communications system buffers message as required until delivery can occur.
 - *Transient Communications:*
 - Message only stored as long as “sender” & “receiver” are executing.
 - If problems occur either deal with them (sender is waiting) or message is simply discarded ...

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

12

Persistence & Synchronicity in Communications

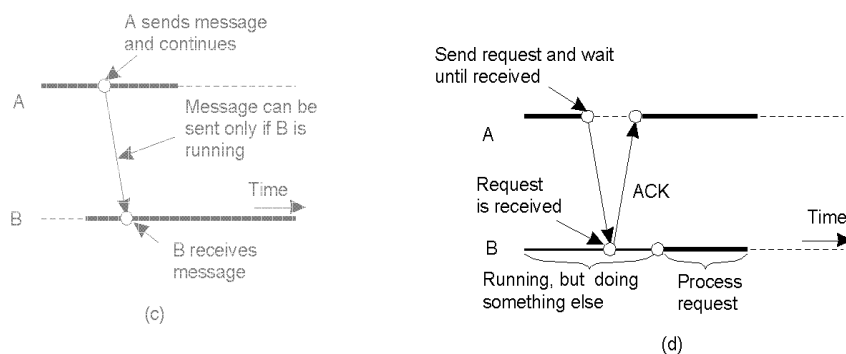
- a) *Persistent asynchronous communication*
- b) *Persistent synchronous communication*



Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

13

Persistence & Synchronicity in Communications (/2)

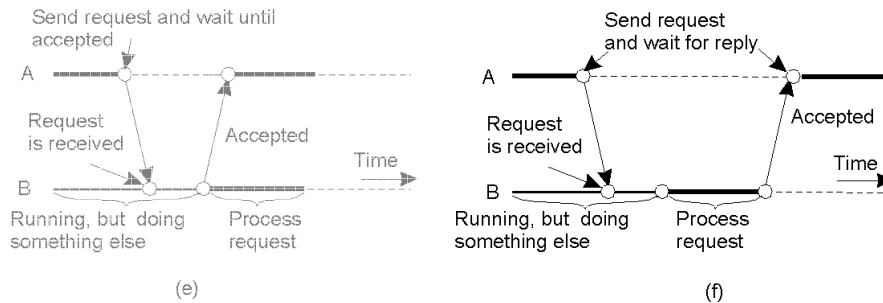


- c) *Transient asynchronous communication*
- d) *Receipt-based transient synchronous communication*

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

14

Persistence & Synchronicity in Communications (/3)



e) *Delivery-based transient synchronous communication at message delivery*

f) *Response-based transient synchronous communication*

SECTION 6.3: REMOTE PROCEDURE CALL (RPC)

Remote Procedure Call (RPC)

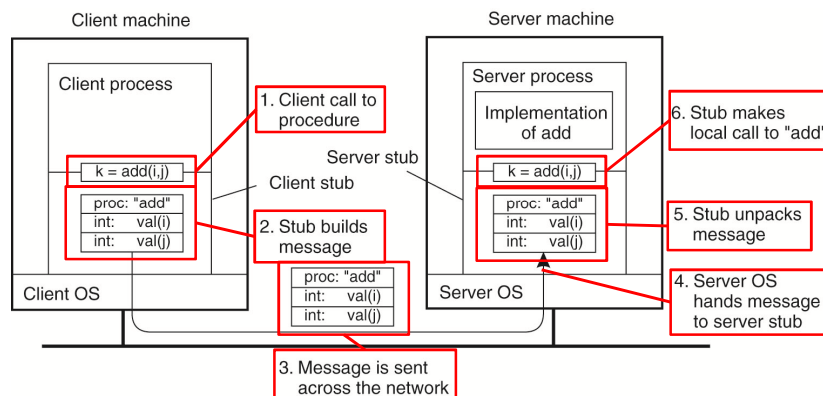
- **Rationale:** *Why RPC?*
- **Distribution Transparency:**
 - Send/Receive don't conceal comms at all – need to achieve *access* transparency.
- **Answer:** *Totally New 'Communication' System:*
 - RPC allows programs to communicate by calling procedures on other machines.
- **Mechanism**
 - When a process on machine A calls a procedure on machine B, calling process on A is suspended,
 - Execution of the called procedure takes place on B.
 - Info 'sent' from caller to callee in parameters & comes back in result.
 - No message passing at all is visible to the programmer.
 - Application developers familiar with simple procedure model.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

17

Basic RPC Operation

1. Client procedure calls client stub
2. Stub builds message, calls local OS.
3. OS sends message to remote OS.
4. Remote OS gives message to stub.
5. Stub unpacks parameters, calls server.
6. Server works, returns result to stub.
7. Stub builds message, calls local OS.
8. OS sends message to client's OS.
9. Client OS gives message to client stub.
10. Stub unpacks result, returns to client.



Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

18

18

RPC: Parameter Passing

- *Parameter marshalling*

More than just wrapping parameters into a message:

- Client/server machines may have different data representations (e.g. byte ordering)
- Wrapping parameter means converting value into byte sequence
- Client and server have to agree on the same encoding:
 - How are basic data values represented (integers, floats, characters)?
 - How are complex data values represented (arrays, unions)?
- Client and server need to properly interpret messages, transforming them into machine-dependent representations.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

19

RPC: Parameter Passing (/2)

- *Assumptions Regarding RPC Parameter Passing:*

- Copy in/copy out semantics: while procedure is executed, nothing can be assumed about parameter values.
- All data that is to be operated on is passed by parameters. Excludes passing references to (global) data.

- *Conclusion*

- Full access transparency cannot be realized

- *Observation:*

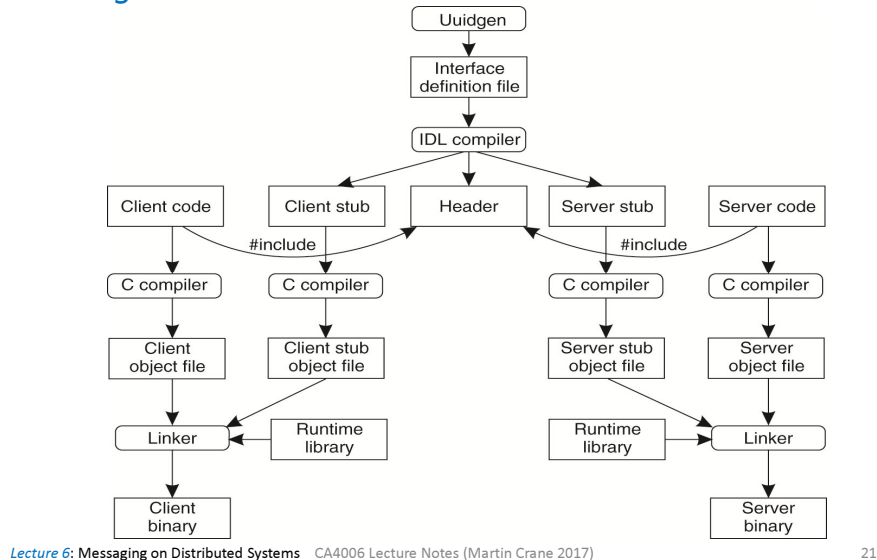
- A remote reference mechanism enhances access transparency: Remote reference offers unified access to remote data
- Remote references can be passed as parameter in RPCs

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

20

RPC Example: Distributed Computing Environment (DCE)

- Writing A Client and Server in DCE:*

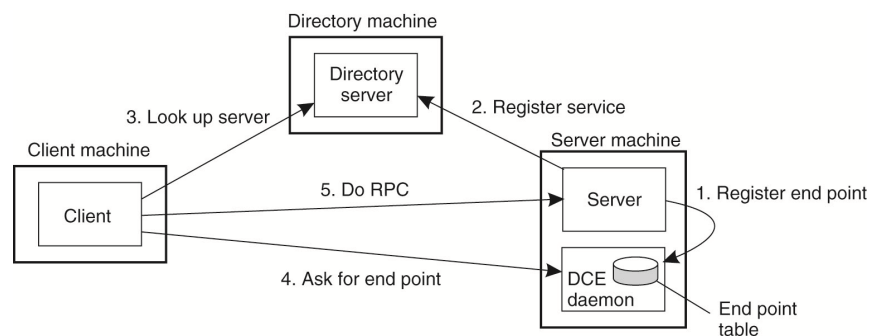


21

DCE Client to Server Binding

- Registration & Discovery:*

- Server registration enables client to locate server and bind to it.
- Server location is done in two steps:
 1. Locate the server's machine.
 2. Locate the server on that machine.



Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

22

SECTION 6.4: MESSAGE QUEUING SYSTEMS

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

23

Message-Oriented Persistent Comms

- Rationale: *Why Another Messaging System?:*
- Scalability:
 - “Transient” messaging systems, do not scale well geographically.
- Granularity:
 - MPI supports messaging $O(\text{ms})$. Distributed message transfer can take minutes
- What about RPC?:
 - In DS can’t assume receiver is “awake” => default “synchronous, blocking” nature of RPC often too restrictive.
- How about Sockets, then?:
 - *Wrong level of abstraction (only “send” and “receive”).*
 - *Too closely coupled to TCP/IP networks – not diverse enough*
- Answer: Message Queueing Systems:
 - MQS give extensive support for *Persistent Asynchronous Communication*.
 - Offer medium-term storage for messages – don’t require sender/receiver to be active during message transmission.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

24

Message-Oriented Persistent Comms. (/2)

- **Message Queuing Systems:**

- *Basic idea:* applications communicate by putting messages into and taking messages out of “message queues”.
- Only guarantee: your message will eventually make it into the receiver’s message queue => “loosely-coupled” communications.
- Asynchronous persistent communication thro middleware-level queues.
- Queues correspond to buffers at communication servers.

- **Four Commands:**

Primitive	Meaning
Put	Append a message to a specified queue.
Get	Block until the specified queue is nonempty, and remove the first message.
Poll	Check a specified queue for messages, and remove the first. Never block.
Notify	Install a handler to be called when a message is put into the specified queue.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

25

Message-Queuing System Architecture

- **Operation:**

- Messages are “put into” a *source queue*.
- They are then “taken from” a *destination queue*.
- Obviously, a mechanism has to exist to move a message from a source queue to a destination queue.
- This is the role of the *Queue Manager*.
- These are message-queuing “relays” that interact with the distributed applications and with each other.
- Not unlike routers, these devices support the notion of a DS “overlay network”.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

26

Role of Message Brokers

- Rationale:

Often need to integrate new/existing apps into a “single, coherent *Distributed Information System* (DIS)”.

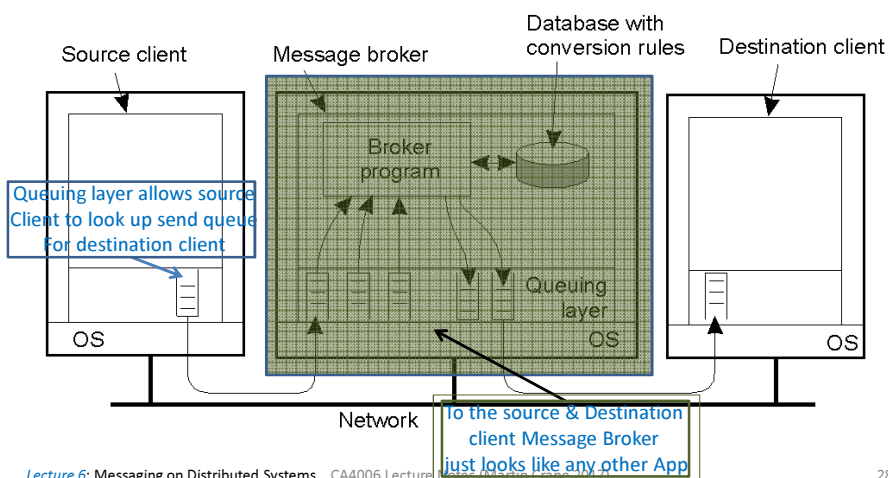
- Problem: different message formats exist in legacy systems
- Can’t “force” legacy systems into single, global message format.
- “Message Broker” allows us to live with different formats
- Centralized component that takes care of application heterogeneity in an MQ system:
 - Transforms incoming messages to target format
 - Very often acts as an application gateway
 - May provide subject-based routing capabilities ⇒ *Enterprise Application Integration*

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

27

Message Broker Organization

- General organization of message broker in a MQS – also known variously as an “interface engine”.



Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

28

IBM's WebSphere MQ

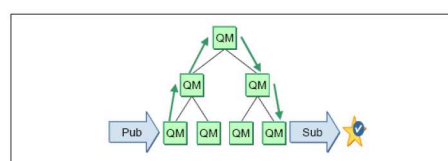
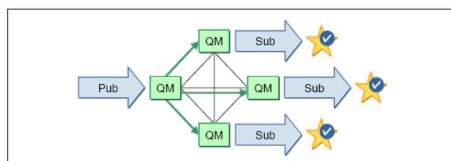
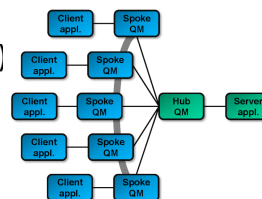
- *Basic concepts:*
 - Application-specific messages are put into, removed from queues
 - Queues reside under the regime of a queue manager
 - Processes can put messages only in local queues, or thro an RPC
- *Message transfer*
 - Messages are transferred between queues
 - Message transfer btw process queues requires a channel
 - At each endpoint of channel is a *message channel agent*
 - Message channel agents are responsible for:
 - Setting up channels using lower-level n/w comm facilities (e.g. TCP/IP)
 - (Un)wrapping messages from/in transport-level packets
 - Sending/receiving packets

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

29

IBM's WebSphere MQ (/2)

- Supported Topologies are:
 1. *Hub/spoke* topology (point-to-point queues)
 - Apps subscribe to "their" QM.
 - Routes to hub QM def'd in spoke QMs.
 2. *Distributed Publish/Subscribe*:
 - Apps subscribe to topics & publish messages to multiple receivers.
 - 2 Topologies: *Clusters* and *Trees*:
 - Cluster*: Cluster of QMs connected by channels. Published messages sent to all connected QMs of the published topic.
 - Tree*: Trees allow reducing number of channels between QMs.



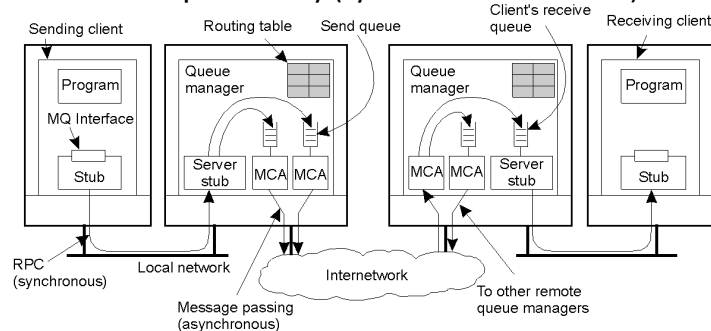
Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

30

IBM's WebSphere MQ (/2)

- *Principles of Operation:*

- Channels are inherently unidirectional
- Automatically start MCAs when messages arrive
- Any network of queue managers can be created
- Routes are set up manually (system administration)



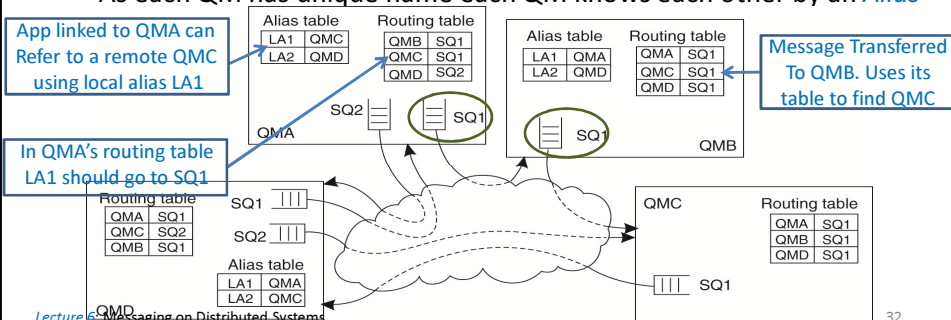
General organization of IBM's WebSphere Message-Queuing System

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

31

IBM's WebSphere MQ (/3)

- *Routing:* Using logical names, in combination with name resolution to local queues, possible to route message to remote queue
 - Sending message from one QM to another (possibly remote) QM, each message needs destination address, so a transmission header is used
 - MQ Address has two parts:
 1. Part 1 is the *Destination QM Name* (say QMX)
 2. Part 2 is the *Name of the Destination Queue* (i.e. QMX's destination Queue)
 - As each QM has unique name each QM knows each other by an *Alias*



Lecture 6: Messaging on Distributed Systems

32

Advanced Message Queuing Protocol (AMQP)

• Why AMQP?

1. Lack of standardization:

- Little standardization in MOM products (mostly proprietary solutions).
 - E.g. 1: JMS Java- dependent, doesn't specify wire protocol only an API.
=> different JMS providers not directly interoperable on wire level.
 - E.g. 2: IBM Websphere clunky and expensive

2. Need for bridges¹ for interoperability:

- To achieve interoperability between different queueing systems, 3rd party vendors offer *bridges*.
- These complicate the architecture / topology, increase costs while reduce performance (additional delay).

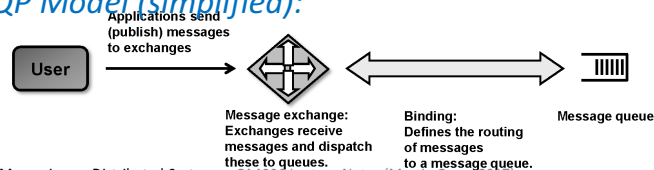
¹Entities that help in different stages of message mediation

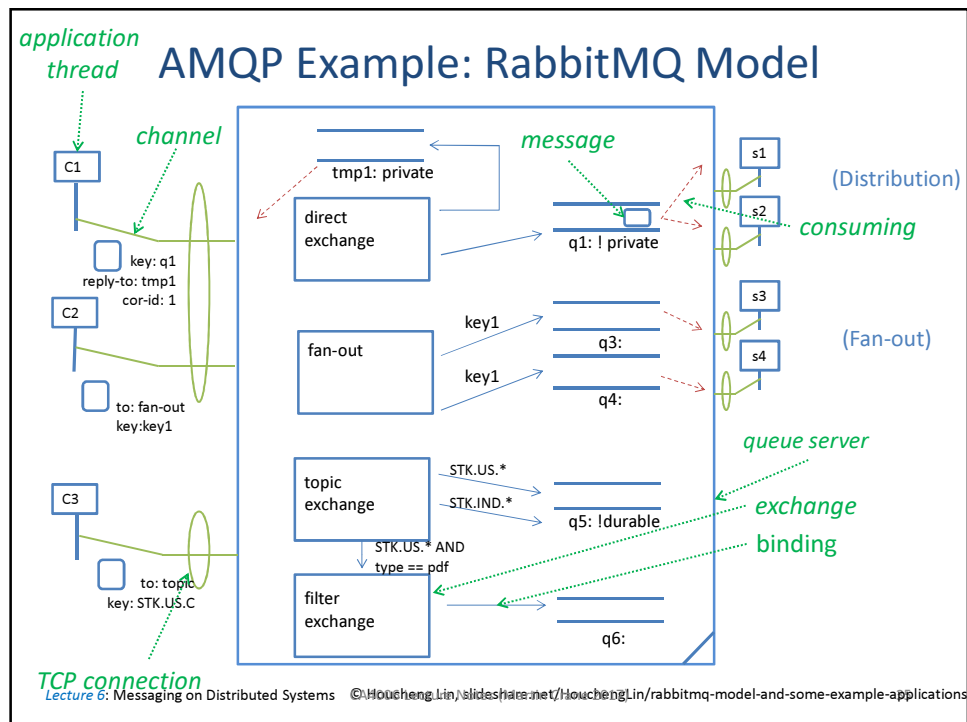
AMQP (/2)

• Characteristics of AMPQ:

- What is it? Open protocol for enterprise messaging, supported by industry (JP Morgan, Cisco, Microsoft, Red Hat, Microsoft etc.).
- Open/ Multi-platform / language messaging system.
- AMQP defines:
 1. Messaging capabilities (called *AMQP model*)
 2. *Wire-level protocol* for interoperability
- AMQP messaging patterns:
 1. Request-response: messages delivered to a specific queue
 2. Publish/Subscribe: messages delivered to a set of receiver queues
 3. Round-robin: message distribution to set of receivers based on availability

• AMQP Model (simplified):





RabbitMQ Model

- **Virtual Host**
- **Exchange**
 - direct
 - fan-out
 - topic
- **Binding**
 - topic
 - cascading
 - message select
- **queue**
 - flags: private, durable
- **Connection**
 - channel: every thread work with one channel
- **Message**
 - content header
 - Properties: Reply-To, Cor-Id, Message-Id, Key
 - queue server may add properties , wont remove/modify
 - content body (won't modify)
 - binary/ file/ stream
- **Application**
 - client/ server

© Houcheng Lin, slideshare.net/HouchengLin/rabbitmq-model-and-some-example-applications

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

36

Hello World in RabbitMQ

```
#!/usr/bin/env ruby
# encoding: utf-8
require "bunny"

conn = Bunny.new(:automatically_recover => false)
conn.start

ch = conn.create_channel
q = ch.queue("hello") # create a message queue called "hello"

ch.default_exchange.publish("Hello World!", :routing_key => q.name)
# default_exchange is a direct exchange with no name
# main advantage is every queue is automatically bound to it with routing key same as queue name
puts "[x] Sent 'Hello World!'"

conn.close # close off the connection
```

```
#!/usr/bin/env ruby
# encoding: utf-8
require "bunny"

conn = Bunny.new(:automatically_recover => false)
conn.start # if conn fails, reconnect tried every 5 secs, this disables automatic connection recovery

ch = conn.create_channel
q = ch.queue("hello") # create a message queue with same name as above

begin
  puts "[*] Waiting for messages. To exit press CTRL+C"
  q.subscribe(:block => true) do |delivery_info, properties, body|
    puts "[x] Received #{body}"
  end
rescue Interrupt => _ # exception handling if Interrupt happens (i.e. if CTRL+C hit)
  conn.close # close off the connection
end
```

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017) 37

RabbitMQ

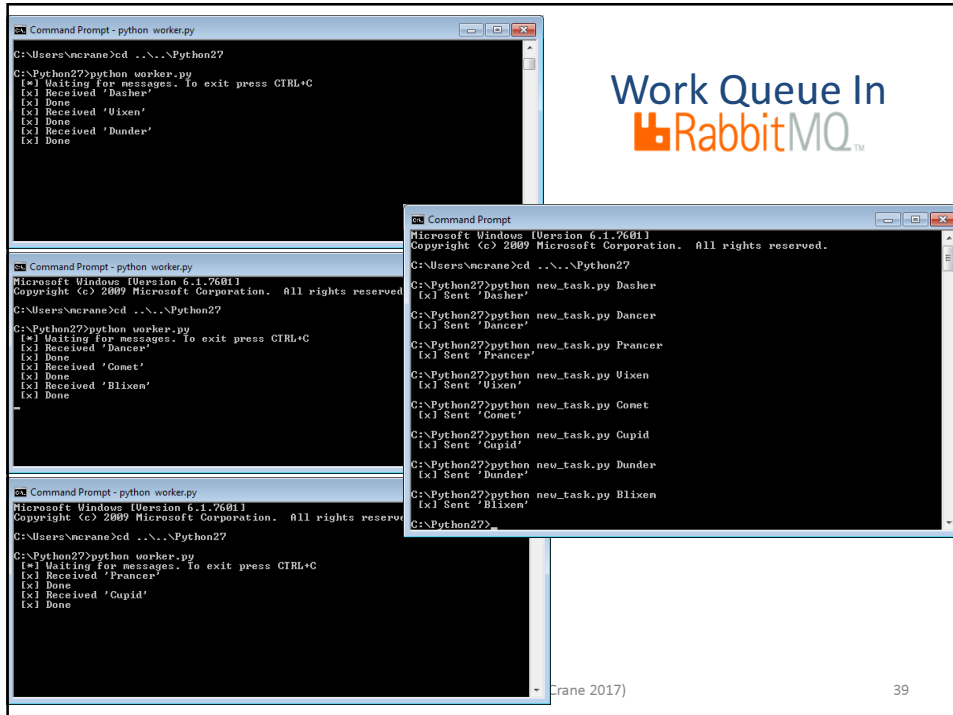
- Afterwards should see something like this:

The screenshot shows the RabbitMQ Management interface. The 'Queues' tab is active, displaying a list of queues. The 'hello' queue is selected, showing its details and message rates. The interface includes a navigation bar with 'Overview', 'Connections', 'Channels', 'Exchanges', 'Queues', and 'Admin'. Below the queue list, there is a table with columns for Name, Features, State, Ready, Unacked, Total, Incoming, deliver / get, and ack. The 'hello' queue is shown with a state of 'idle' and message rates of 0.00/s for incoming, deliver, and ack.

The screenshots show the execution of the Ruby scripts in a Windows Command Prompt. The top prompt shows the execution of 'python send.py' which outputs '[x] Sent 'Hello World!'' and 'C:\Python27>'. The bottom prompt shows the execution of 'python receive.py' which outputs '[*] Waiting for messages. To exit press CTRL+C' and '[x] Received 'Hello World!''. Both prompts are in a Windows environment.

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017) 38

Work Queue In

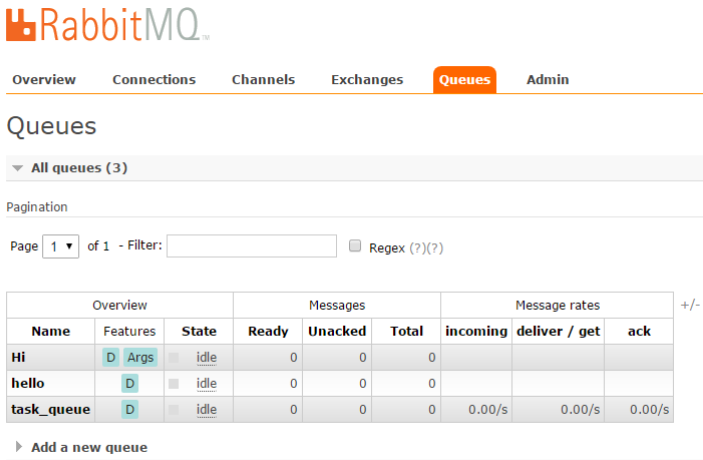


Crane 2017)

39

Work Queue In (/2)

- Afterwards should see something like this:



HTTP API | Command Line

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

40

```
C:\Command Prompt
C:\Python27>python new_task.py Francer
[+] Sent 'Francer'
C:\Python27>python new_task.py Uixen
[+] Sent 'Uixen'
C:\Python27>python new_task.py Comet
[+] Sent 'Comet'
C:\Python27>python new_task.py Cupid
[+] Sent 'Cupid'
C:\Python27>python new_task.py Dunder
[+] Sent 'Dunder'
C:\Python27>python new_task.py Blitzen
[+] Sent 'Blitzen'
C:\Python27>python cmd list bindings
pythoncmd is not recognized as an internal or external command,
operable program or batch file.
C:\Python27>python exit_log.py
[+] Sent 'infer Hello World!'
C:\Python27>python exit_log.py 0! CH4006!
[+] Sent '0! CH4006!'
```

41

```

graph LR
    P((P)) --> X((X))
    X -- "type=topic *orange.*" --> Q1[Q1]
    X -- "*.*rabbit laez,y.#" --> Q2[Q2]
    Q1 --> C1((C1))
    Q2 --> C2((C2))
  
```

[illegible]

42

a

b Multicast on Chord Network¹

¹from Talia & Trunfrio, J. Parallel & Dist Computing Vol(70(12)) pp1254 - 1265, 2010

2. Epidemic Algorithms

- *Essence:*
- Epidemic algorithms used to rapidly spread info in large P2P systems without setting up a multicast tree
- Assumptions:
 - All updates for specific data item are done at a single node (i.e., no write-write conflict)
 - Can distinguish old from new data as data is time stamped or versioned
- Operation:
 - Node receives an update, forwards it to randomly chosen peers (akin to spreading a contagious disease)
 - Eventually, each update should reach every node
 - Update propagation is lazy

Lecture 6: Messaging on Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

47

2. Epidemic Algorithms (/2)

- *Glossary of Terms:*
 - Node is *infected* if it has an update & wants to send to others
 - Node is *susceptible* if it has not yet been updated/infected
 - Node is *removed* if it is not willing or able to spread its update or can no longer send to others for some reason.
- We study two propagation models here:
 - *Anti-entropy*
Each replica regularly chooses another randomly & exchanges state differences, giving identical states at both afterwards.
 - *Gossiping:*
A replica which has just been updated (i.e., has been infected), tells other replicas about its update (infecting them as well).

Lecture 6: Messaging in Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

48

2. Epidemic Algorithms (/3)

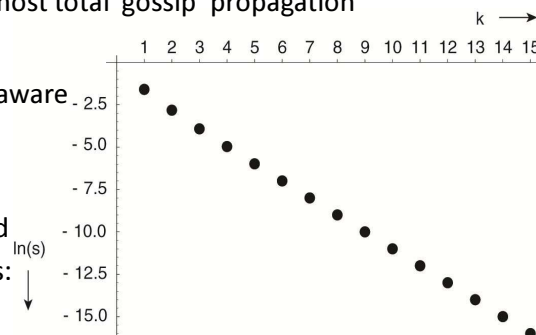
- *Principal Operations of Anti-Entropy:*
 - A node P selects another node Q from the system at random.
 - *Push:* P only sends its updates to Q
 - *Pull:* P only retrieves updates from Q
 - *Push-Pull:* P and Q exchange mutual updates (after which they hold the same information).
- *Observations*
 - For push-pull it takes $O(\log(N))$ rounds to disseminate updates to all N nodes (*round*= when every node has initiated an exchange).
 - Anti-Entropy is reliable but costly (each replica must regularly choose another randomly)

Lecture 6: Messaging in Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

49

2. Epidemic Algorithms (/4)

- *Basic model of Gossiping:*
 - A server S having an update to report, contacts other servers.
 - If a server is contacted to which update has already propagated, S stops contacting other servers with probability $1/k$.
 - i.e. increasing k ensures almost total 'gossip' propagation
- *Observations*
 - If s is fraction of servers unaware of update, can show that with many servers, the equation $s = e^{-(k+1)(1-s)}$ is satisfied
 - Example: for 10,000 servers: when $k = 4, s < 0.007$
 - If need 100% propagation, gossiping alone is not enough, maybe need to run one round of anti-entropy.



Lecture 6: Messaging in Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

50

2. Epidemic Algorithms (/5)

- *The Deletion Problem in Epidemic Algorithms:*
 - Cannot remove old value from a server, expecting removal to propagate.
 - Instead, mere removal will be undone in time using epidemic algorithms
- *Solution:* Must register removal as special update by inserting a death cert
- *Next problem:*
 - When to remove a death certificate (it is not allowed to stay for ever)?
 - Run a global algorithm to detect if removal is known everywhere, and then collect the death certificates (looks like *garbage collection*) or
 - Assume death certificates propagate in finite time, and associate max lifetime for a certificate (can be done at risk of not reaching all servers)
 - Note: It is necessary that a removal actually reaches all servers.
- *Applications of Epidemic Algorithms:*
 - (Obviously) data dissemination
 - Data aggregation: each node with value x_i . Two nodes gossiping should reset their variable to $(x_i + x_j)/2$. What final value will nodes possess?

Lecture 6: Messaging in Distributed Systems CA4006 Lecture Notes (Martin Crane 2017)

51

Lecture Summary

- Middleware enables much functionality in DS
- Especially the many types of interaction/communications necessary
- With rational reasons for every one!
 - *Remote Procedure Call* (RPC) enables transparency
 - But *Message Queuing Systems* necessary for persistent communications
 - IBM Websphere is ok but a bit old, clunky & tired at this stage?
 - AMQP open source, more flexible, better Industrial support?
 - *Multicast Communications* are often necessary in DS:
 - Application Layer Messaging (ALM)
 - Epidemic Protocols

Lecture 6: Messaging in Distributed Systems

CA4006 Lecture Notes (Martin Crane 2017)

52