# Communication in Distributed Systems – Fundamental Concepts

## Hong-Linh Truong

**Note:**

The content is based on the latest version that I lectured in 2014 in TU Wien

Some old concepts will be updated later.

# What is this lecture about?

- Understanding basic terminologies in communication in distributed systems

- Understanding key concepts in communication in distributed systems
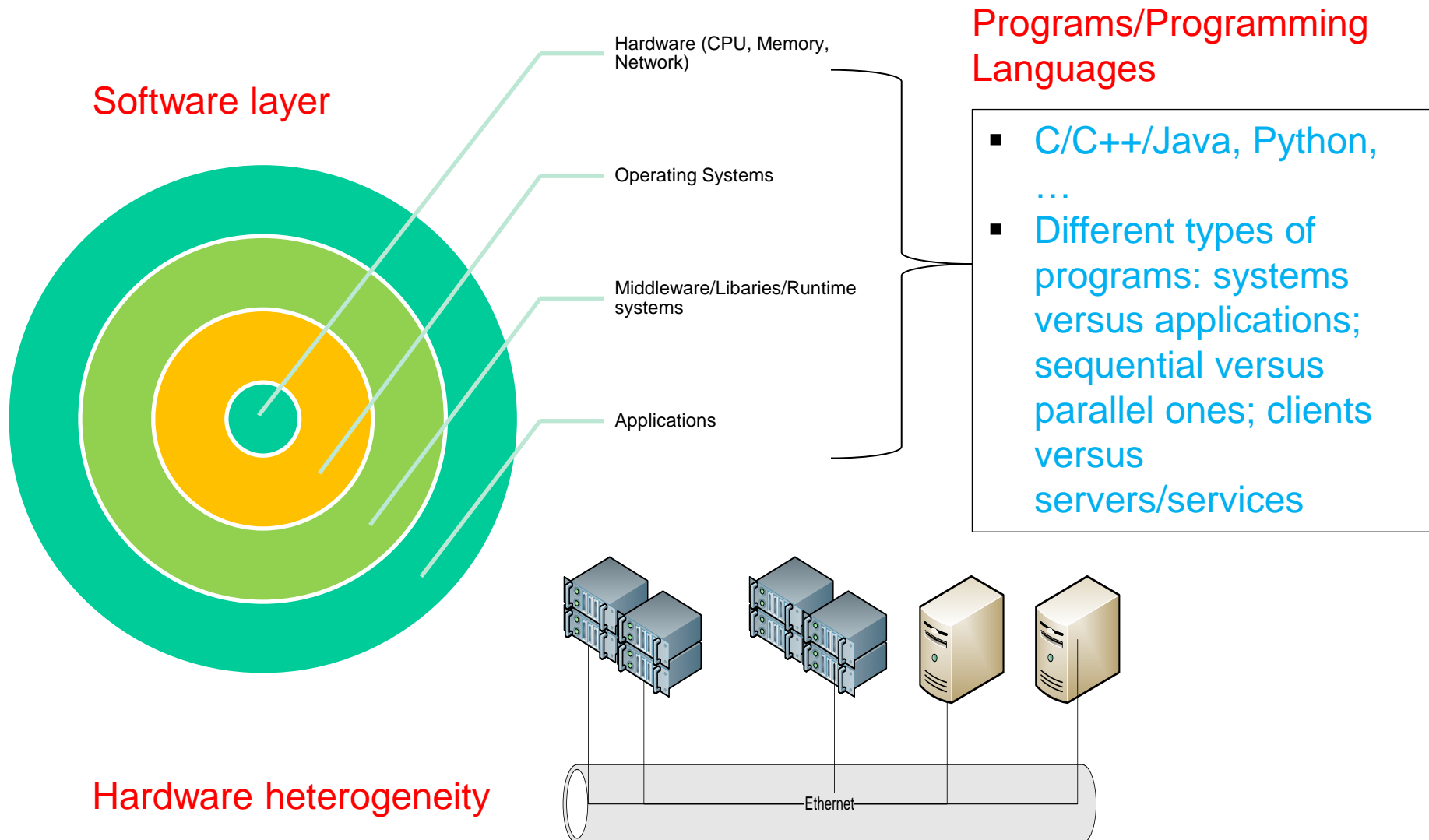
# Learning Materials

- Main readings/references for the lecture:
  - Tanenbaum & Van Steen, Distributed Systems: Principles and Paradigms, 2e, (c) 2007 Prentice-Hall
    - Chapters 3 & 4
  - George Coulouris, Jean Dollimore, Tim Kindberg, Gordon Blair„Distributed Systems – Concepts and Design", 5nd Edition
    - Chapters 2,3, 7.
  - Craig Hunt, TCP/IP Network Administration, 3edition, 2002, O'Reilly.
- Test the examples in the lecture

# Outline

- Communication entities, paradigm, roles/responsibilities
- Key issues in communication in distributed systems
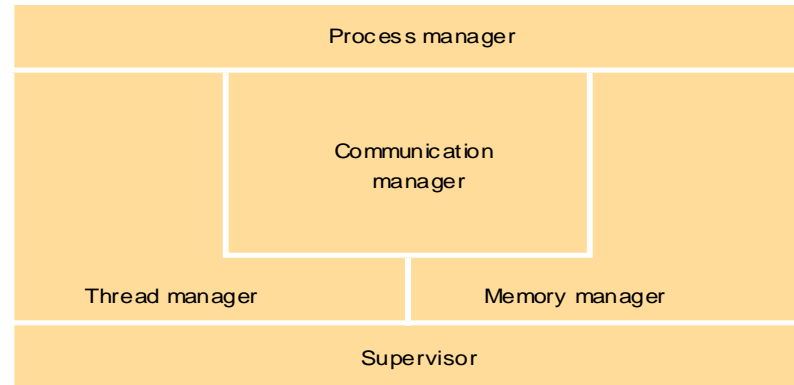- Protocols
- Processing requests
- Summary

# COMMUNICATION ENTITIES, PARADIGM, AND ROLES
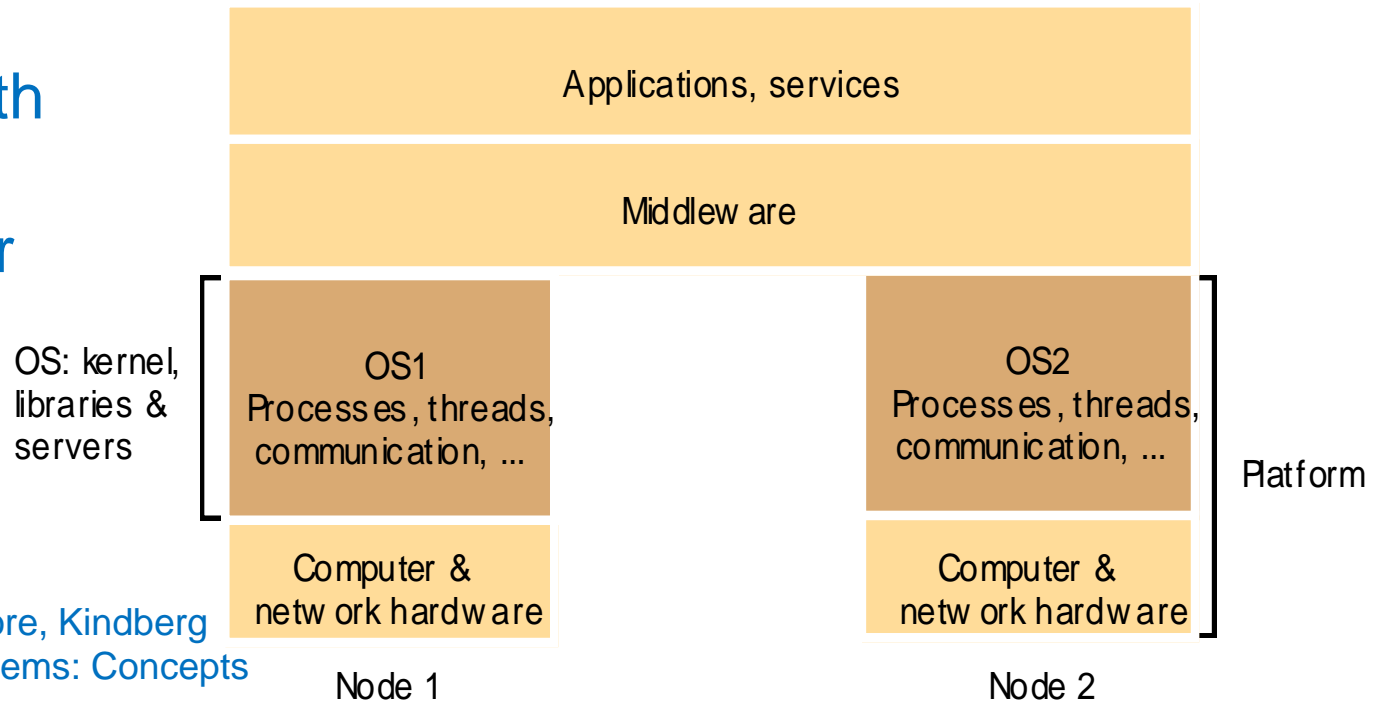
# Hardware, software layer, programs

Software layer

Hardware (CPU, Memory, Network)

Operating Systems

Middleware/Libaries/Runtime systems

Applications

Programs/Programming Languages

- C/C++/Java, Python, …
- Different types of programs: systems versus applications; sequential versus parallel ones; clients versus servers/services

Hardware heterogeneity

Ethernet

Distributed Systems, Bachelor Study                    7

# System Layers and Core OS functionality

Core OS functionality

Different OSs with a common middleware layer



OS: kernel, libraries & servers

Platform

Process manager

Communication manager

Thread manager | Memory manager

Supervisor

Applications, services

Middleware

OS1 Processes, threads, communication, ...

OS2 Processes, threads, communication, ...

Computer & network hardware

Computer & network hardware

Node 1

Node 2

Distributed Systems, Bachelor Study

# Process versus thread

Program → executed → P P P P P

Distributed processes of the same service (coded in the same program)

Thread of execution
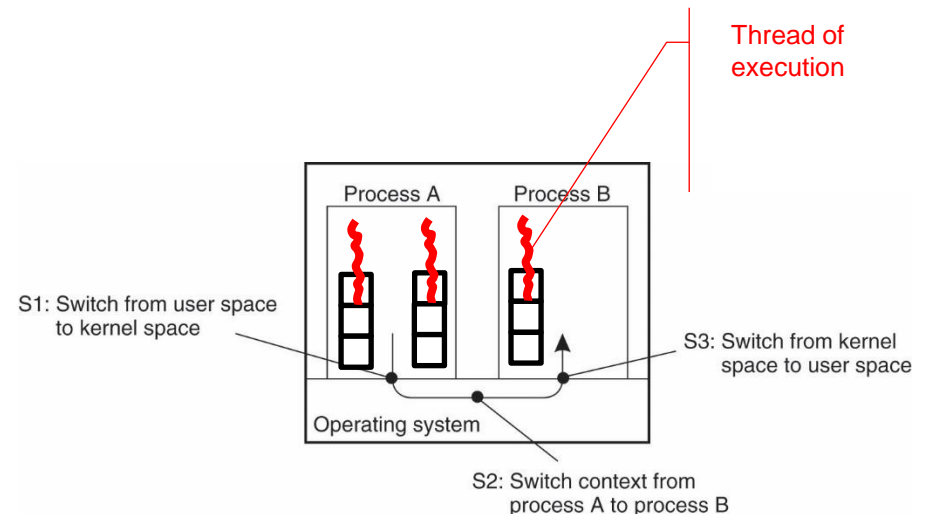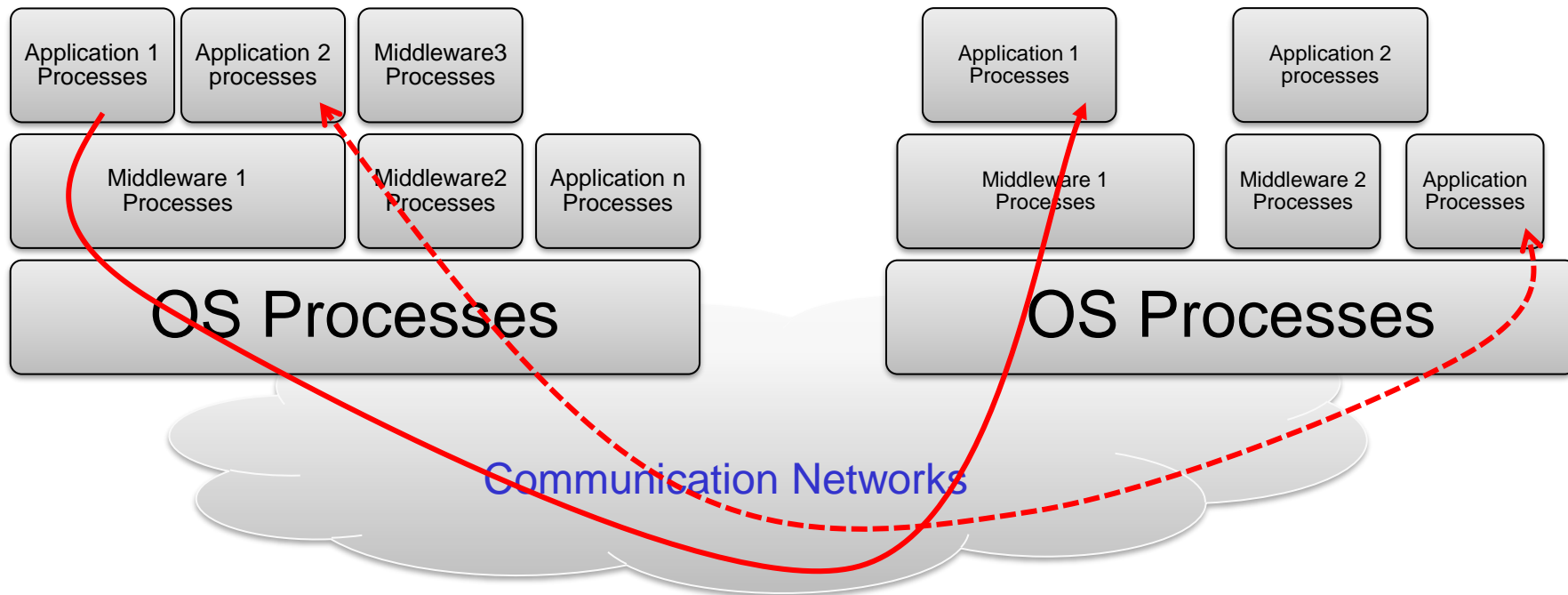
## Within a non distributed OS

- Process – the program being executed by the OS

- Threads within a process

- Switching thread context is much cheaper than that for the process context

- Blocking calls in a thread do not block the whole process

Process A | Process B

S1: Switch from user space to kernel space

S3: Switch from kernel space to user space

Operating system

S2: Switch context from process A to process B

Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

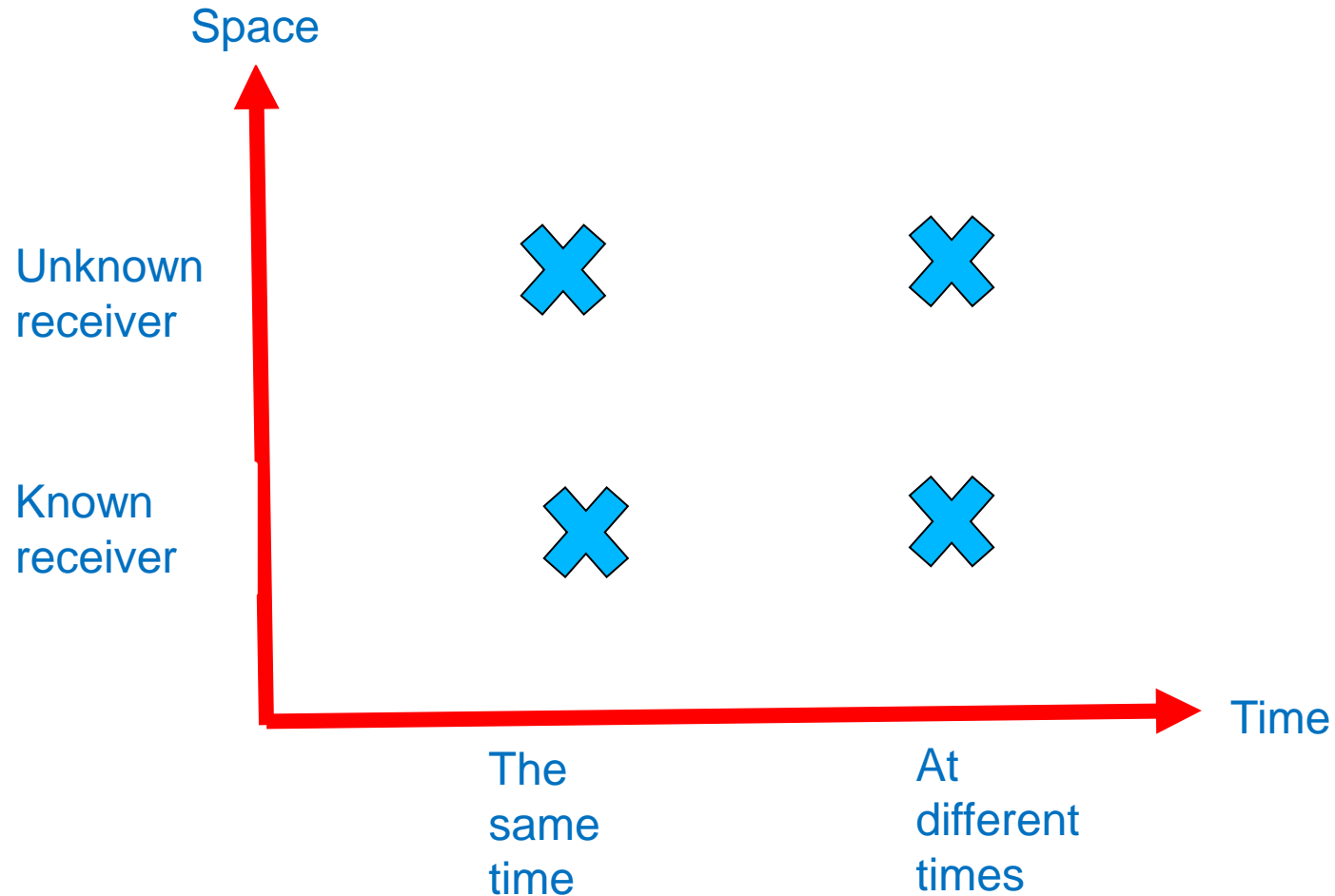Distributed Systems, Bachelor Study          9

# Communication entities



Communication in distributed systems
- between processes within a single application/middleware/service
- among processes belonging to different applications/middleware/services
- Among computing nodes which have no concept of processes (e.g. sensors)

Q: Identify some concrete types of communication entities in real-world distributed systems (e.g., in a parallel cluster system)

# Space and Time in communication

Space

Unknown receiver

Known receiver

Time

The same time

At different times

Q: why is understanding time and space uncoupling important for implementing communication in distributed systems?

# Communication Paradigm

- Interprocess communication
  - Low-level message-based communication, e.g., when communication entities are processes

- Remote invocation
  - (direct) calling of remote functions (of services/objects)

- Indirect communication
  - Communication carried out through third parties

# Communication roles and responsibilities

- Several terms indicating communication entities
    - Objects, components, processes  or services, clients, servers
    - forms versus roles/responsibilities

- Roles
    - Client/Server: client requests  - server serves!
    - Sender/Receiver: w.r.t send/receive operation
    - Service:   w.r.t. offering functionality
        - Network service, software-as-a-service,

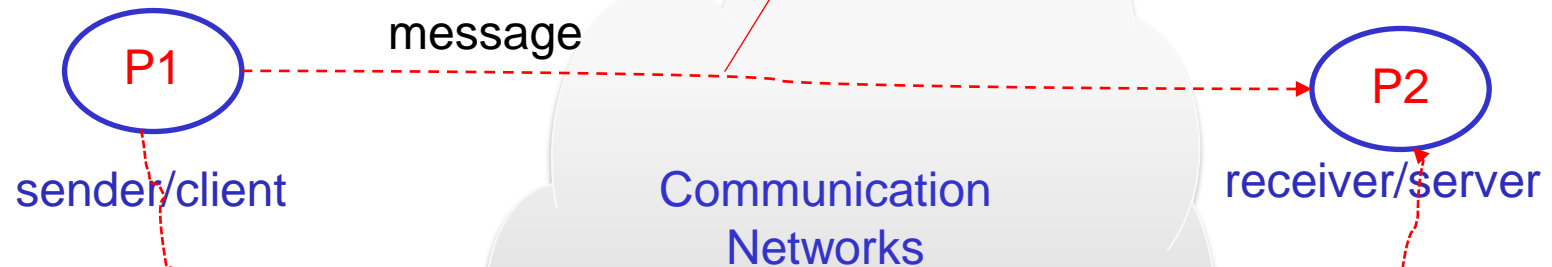Q:  Can a service have multiple servers placed in different machines?

# Communication networks in distributed systems

- Maybe designed for specific types of environments
    - High performance computing, M2M (Machine-to-Machine), building/home/city management, etc.
    - Events, voices, documents, image data, etc.
- Distributed, different network spans
    - Personal area networks (PANs), local area networks (LANs), campus area networks (CANs), metropolitan area networks (MANs), and wide area networks (WANs)
    - Communication entities are placed in different locations
- Different layered networks for distributed systems
    - Physical versus overlay network topologies (virtual network topologies atop physical networks)

# Layered communication

In the view of P1 and P2

End-to-end process-to-process communication
e.g., email abc@tuwien.ac.at to ab@gmail.com

P1

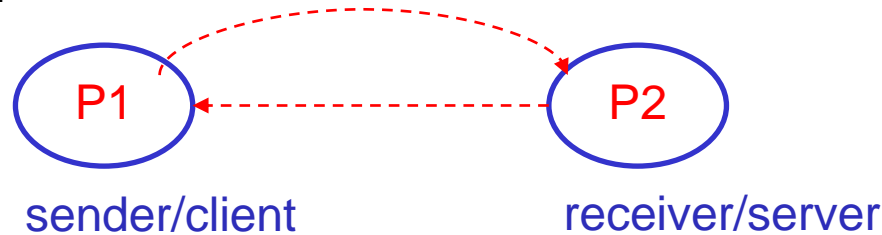message

P2

sender/client

receiver/server

Communication Networks

Holistic system view

Pk

Pm

End-to-end process-to-process communication

# Communication Patterns

One-to- one/client-server

P1 ← ⟶ P2

sender/client   receiver/server

Group communication

P1 → P2
P1 → P3
P1 → P4

P2 → P1
P3 → P1
P4 → P1

Q:  What are the benefits of group communication? Give some concrete examples (e.g., in P2P and social networks).
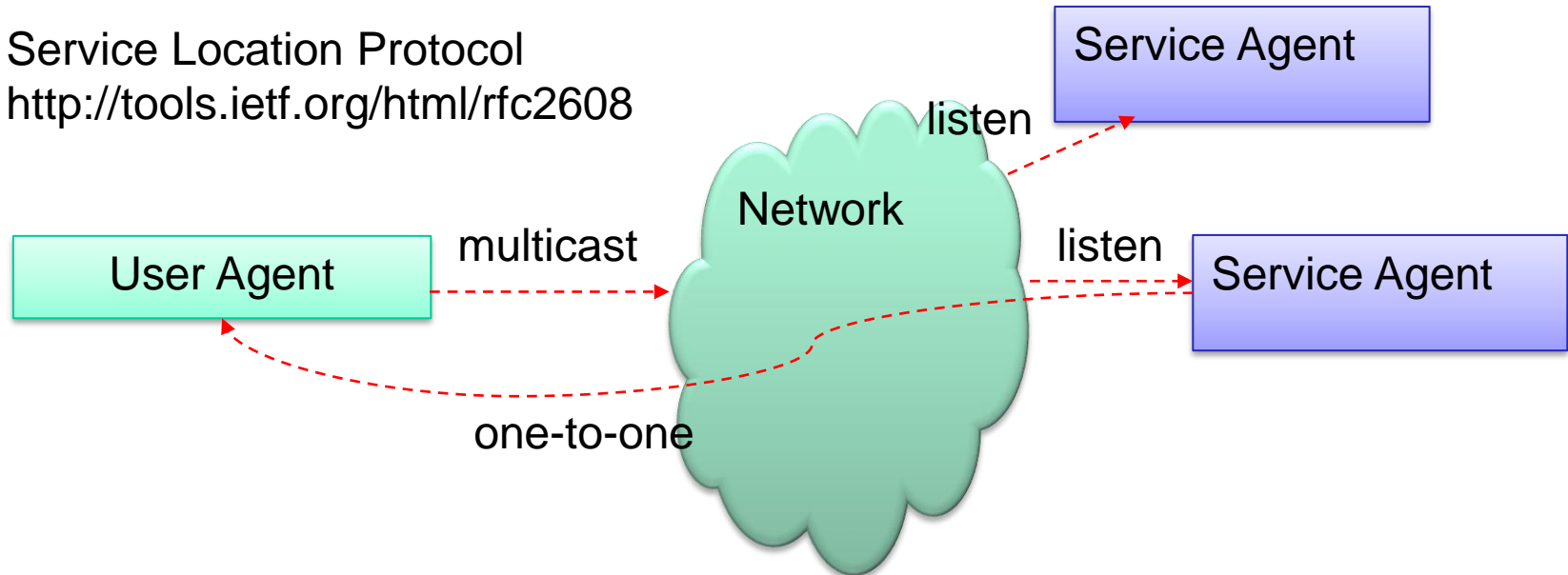
# Identifiers of entities participating in communication

- Communication cannot be done without knowing identifiers (names) of participating entities
  - Local versus global identifier
  - Individual versus group identifier
- Multiple layers/entities → different forms of identifiers
  - Process ID in an OS
  - Machine ID: name/IP address
  - Access point: (machine ID, port number)
  - A unique communication ID in a communication network
  - Emails for humans
  - Group ID

# Examples of communication patterns (1)

Service Location Protocol
http://tools.ietf.org/html/rfc2608

Service Agent

listen

Network

Service Agent

User Agent

multicast

listen

one-to-one

- ■ A User Agent  wants to find a Service Agent

- ■ Different roles and different communication patterns

- ■ Get http://jslp.sourceforge.net/ and play samples to see how it works

# Examples of communication patterns (2)



Communication world

source

- ## MPI (Message Passing Interface)

```
MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
  MPI_Comm_rank(MPI_COMM_WORLD,&myid);
  source=0;
  count=4;
  if(myid == source){
   for(i=0;i<count;i++)
     buffer[i]=i;
  }
MPI_Bcast(buffer,count,MPI_INT,source,MPI_COMM_WORLD);
```

$sudo apt-get install mpich
$mpicc c_ex04.c
$mpirun –np 4 ./a.out

http://geco.mines.edu/workshop/cla
ss2/examples/mpi/c_ex04.c

Distributed Systems, Bachelor Study          19

# Connection-oriented or connection less communication

The message: „there is a party tonight"

P1 - - - - - - - - - - - - - - - - - - - - - - -> P2

Write the message in a letter

Go to the post office

Send the letter to P2

Find the phone number of P2

Call P2

Tell P2 the message

Connection-oriented communication between P1 and P2 requires the setup of communication connection between them first – no setup in connectionless communication

Q: What are the pros/cons of connection-oriented/connectionless communications? Is it possible to have a connectionless communication between (P1,P2) through some connection-oriented connections?

# Blocking versus non-blocking communication calls

Individual process/machine boundary

Sending message buffer

Send operation → msg

Individual process/machine boundary

Receiving message buffer

Receive operation

Send: transmitting a message is finished, it does not necessarily mean that the message reaches its final destination.

- Blocking: the process/thread execution is suspended until the message transmission finishes

- Non-blocking: the process/thread execution continues without waiting until the finish of the message transmission

Q: Analyze the benefits of non-blocking communication. How does non-blocking receive() work?
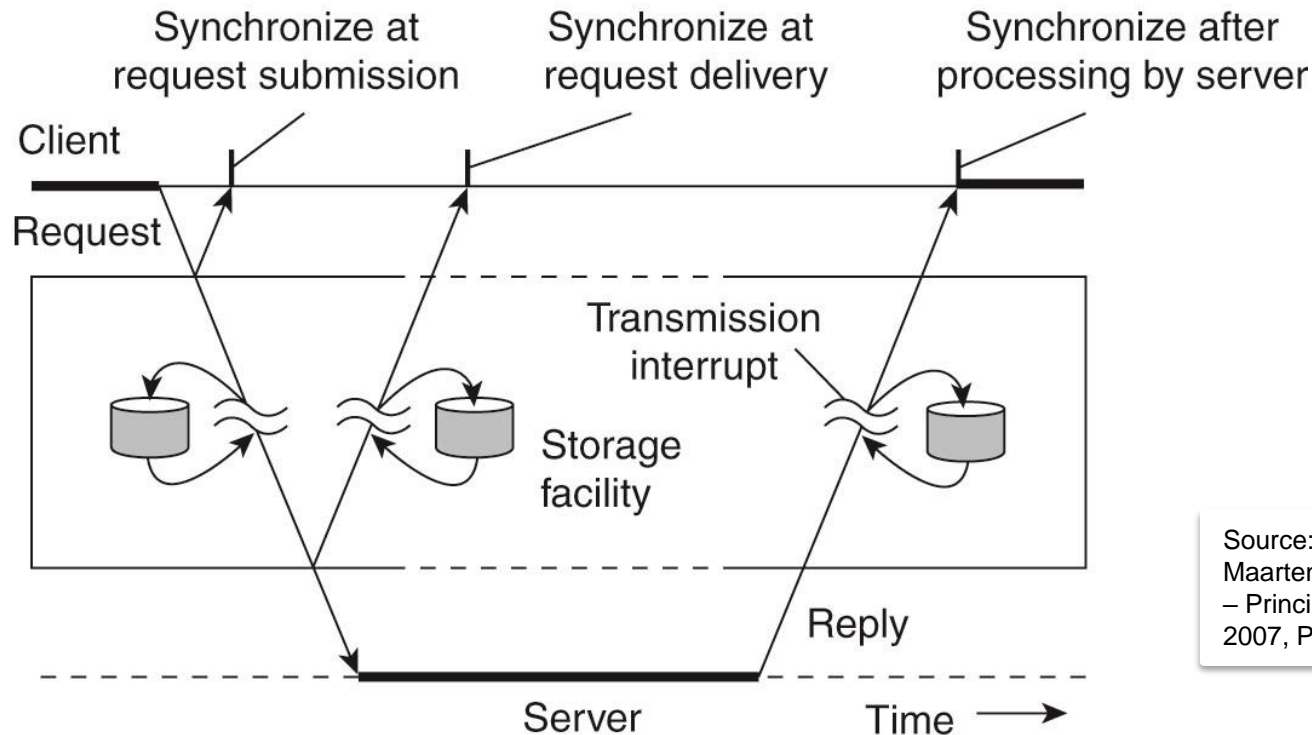
# Persistent and transient communication

- Persistent communication
    - <span style="color:red">Messages are kept</span> in the communication system <span style="color:red">until</span> they are delivered to the receiver
    - Often storage is needed
- Transient communication
    - <span style="color:red">Messages are kept</span> in the communication temporary <span style="color:red">only if</span> both the sender and receiver are live

# Asynchronous versus synchronous communication

- Asynchronous: the process continues after as soon as sending messages have been copied to the local buffer
    - Non blocking send; receive may/may not be blocking
    - Callback mechanisms
- Synchronous: the sender waits until it knows the messages have been delivered to the receiver
    - Blocking send/blocking receive
    - Typically utilize connection-oriented and keep-alive connection
    - Blocking request-reply styles
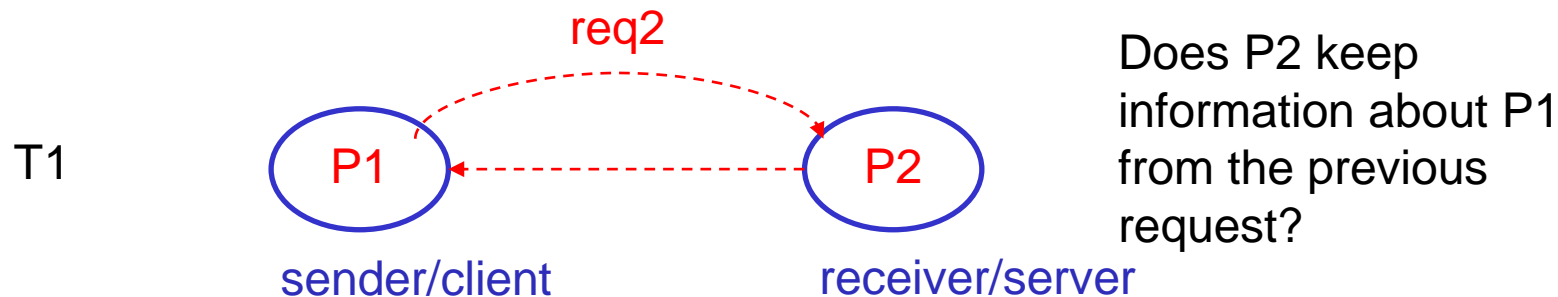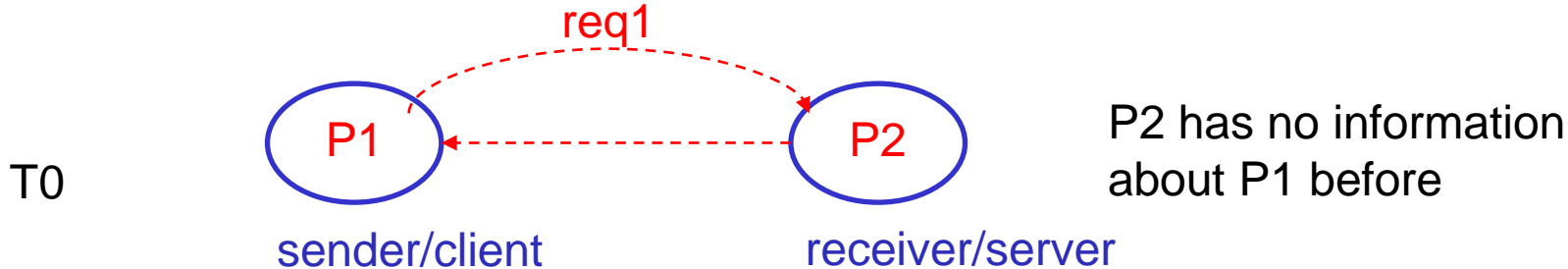
# Different forms of communication



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

Q: How can we achieve the „persistent communication"? What are possible problems if a server sends an accepted/ACK message before processing the request?

# Stateful versus Stateless Server

**req1**

P1 (sender/client) ⟷ P2 (receiver/server)

**T0**

P2 has no information about P1 before

---

**req2**

P1 (sender/client) ⟷ P2 (receiver/server)

**T1**

Does P2 keep information about P1 from the previous request?
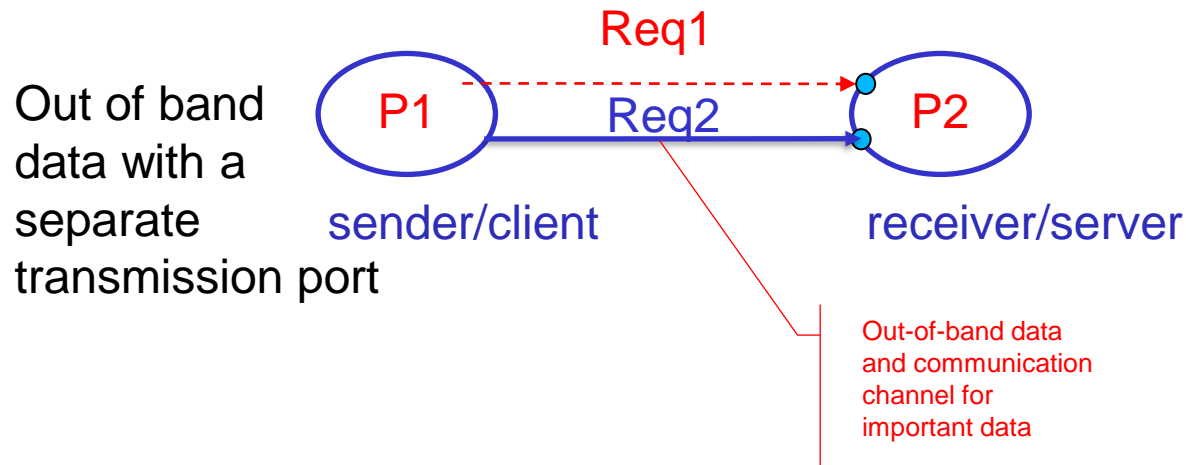
| Stateless server | Soft State | Stateful Server |
|---|---|---|
| Does not keep client's state information | Keep some limited client's state information in a limited time | Maintain client's state information permanently |

Q: Give an example of a stateless communication built atop stateless communication. Analyze "web cookie" w.r.t. stateless/stateful support.

# Handling out of band data

Req2, Req1

Normal case

P1 -----> P2

sender/client          receiver/server

All messages come to P2 in the same port, no clear information about priority

Out of band data with a separate transmission port

Req1

P1        Req2        P2

sender/client          receiver/server

Out-of-band data and communication channel for important data

Q: How can out-of-band data and normal data be handled by using the same transmission channel?

# COMMUNICATION PROTOCOLS

# Some key questions – Protocols

The message: „there is a party" tonight

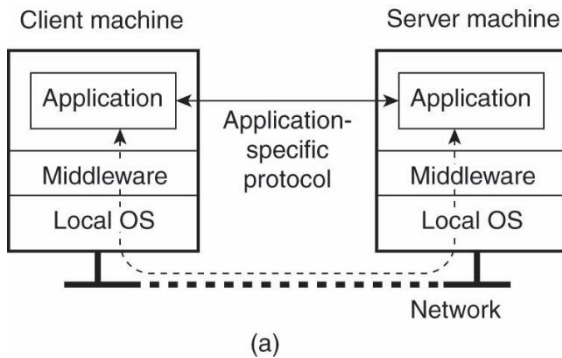P1  - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - →  P2

- **Communication patterns**
  - Can I use a single sending command to send the message to multiple people?
- **Identifier/Naming/Destination**
  - How do I identify the guys I need to send the message
- **Connection setup**
  - Can I send the message without setting up the connection
- **Message structure**
  - Can I use German or English to write the message
- **Layered communication**
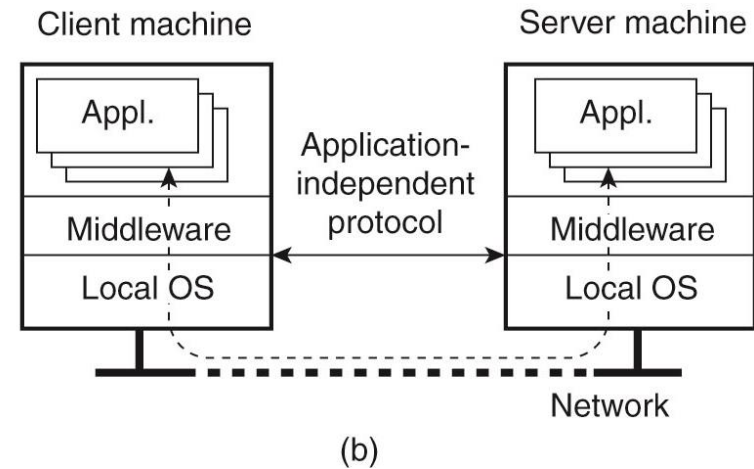  - Do I need other intermediators to relay the message?
- ...

A communication protocol will describe rules addressing these issues

# Applications and Protocols

Application-specific protocols

Application-independent protocols



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

# Layered Communication Protocols

- Complex and open communication requires <span style="color:red">multiple communication protocols</span>

- Communication protocols are typically organized into differ layers: layered protocols/protocol stacks

- Conceptually: each layer has a set of different <span style="color:red">protocols for certain communication functions</span>
  - Different protocols are designed for different environments/criteria

- <span style="color:red">A protocol suite</span>: usually a set of protocols used together in a layered model
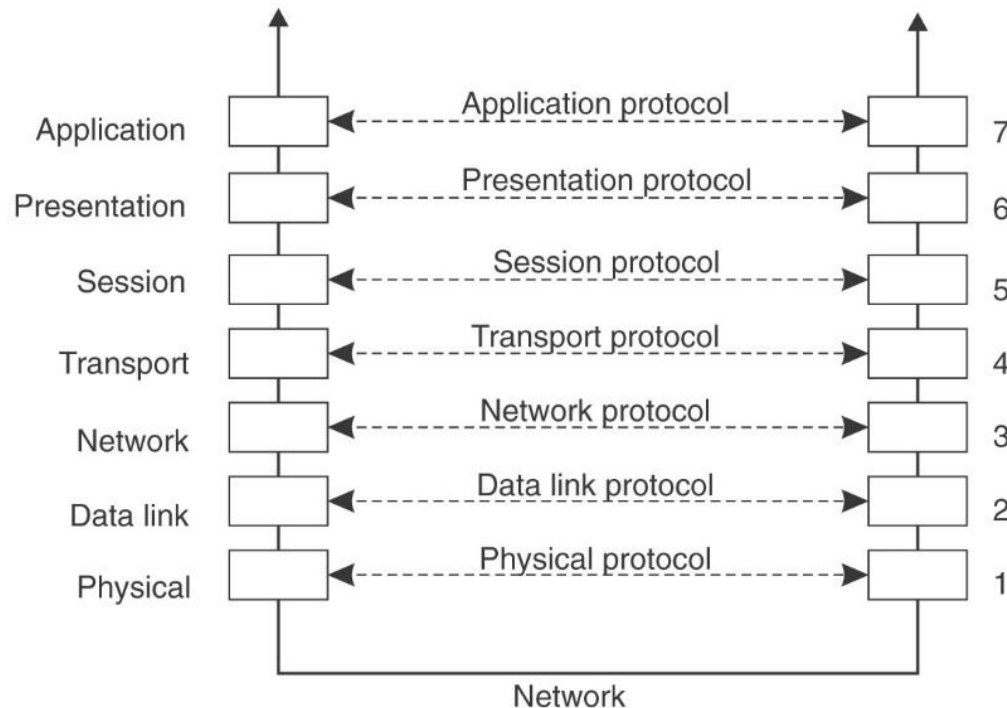
# OSI – Open Systems Interconnection **Reference** Model



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall
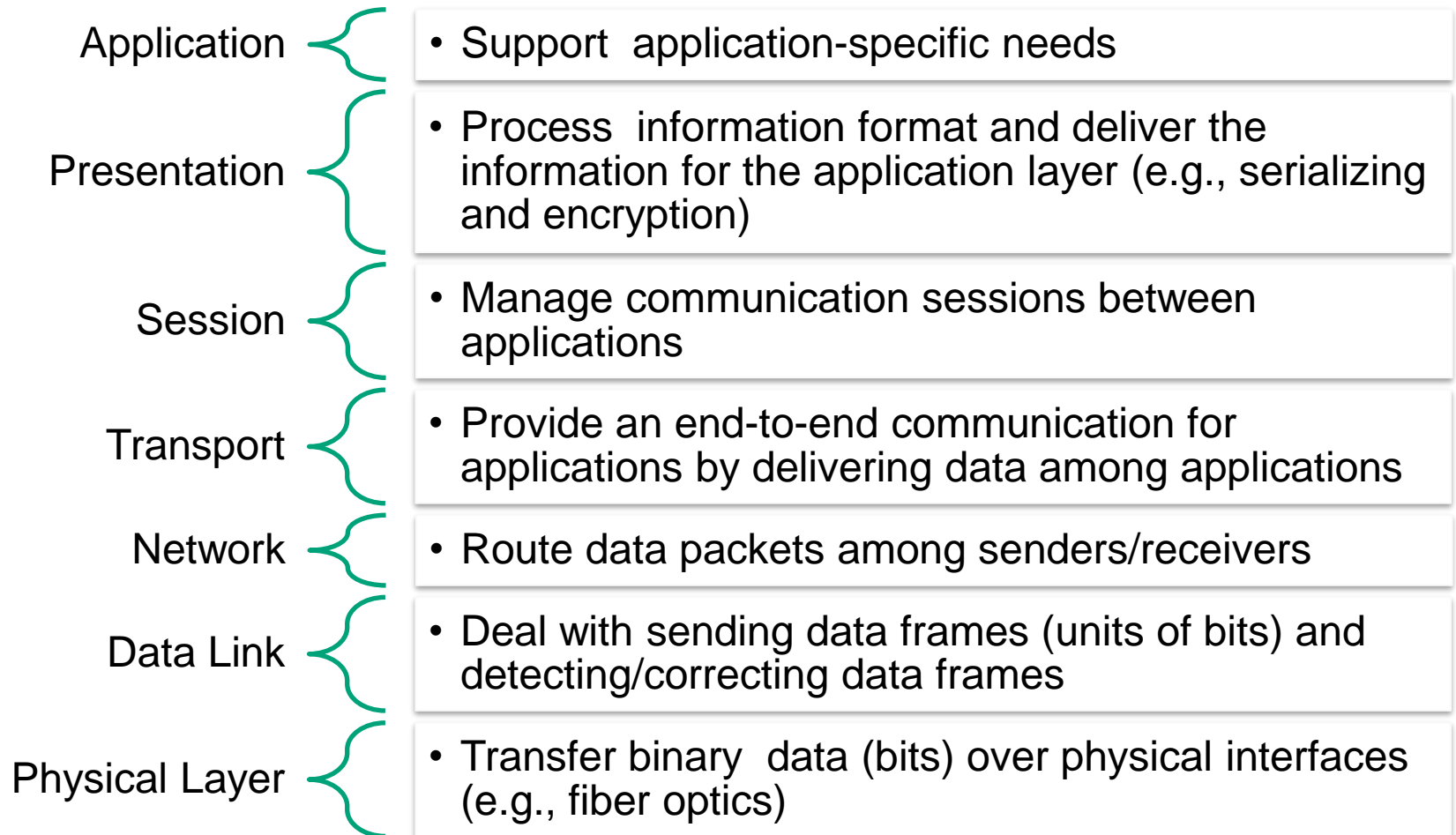
# OSI Layers

Application
- Support application-specific needs

Presentation
- Process information format and deliver the information for the application layer (e.g., serializing and encryption)

Session
- Manage communication sessions between applications

Transport
- Provide an end-to-end communication for applications by delivering data among applications

Network
- Route data packets among senders/receivers

Data Link
- Deal with sending data frames (units of bits) and detecting/correcting data frames

Physical Layer
- Transfer binary data (bits) over physical interfaces (e.g., fiber optics)
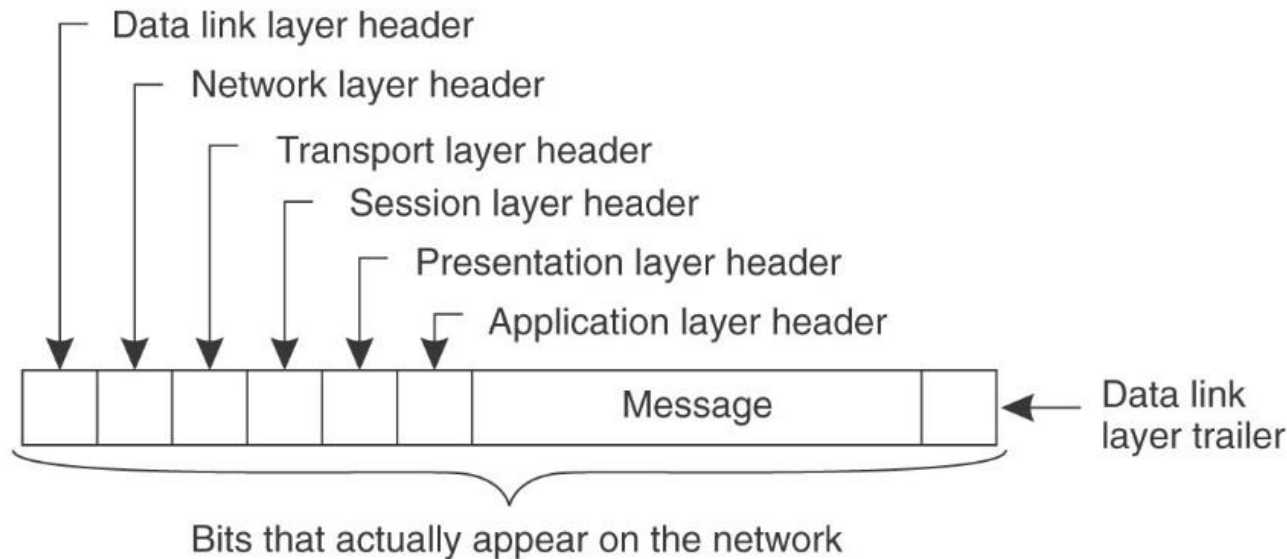
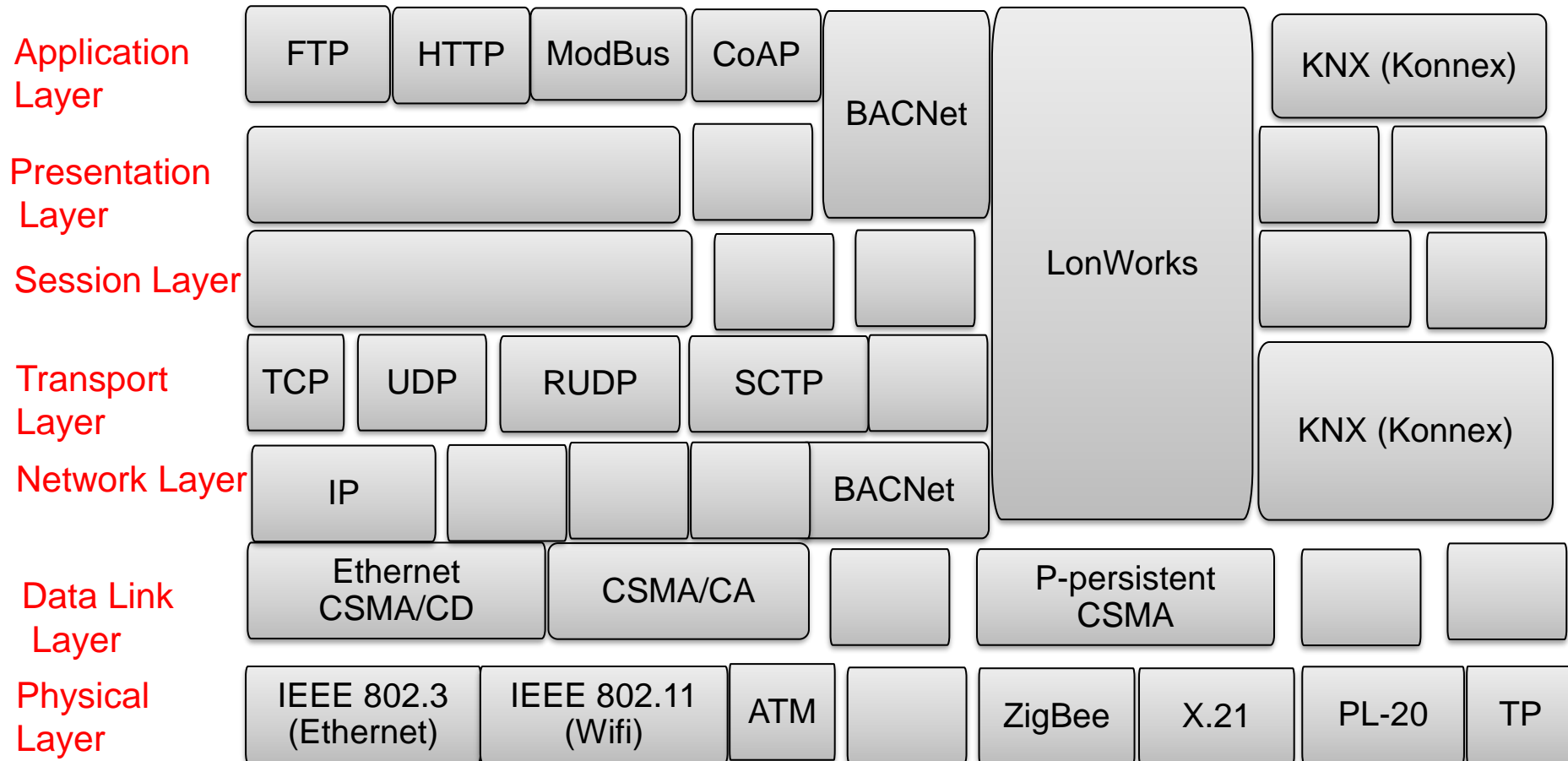Distributed Systems, Bachelor Study          32

# How layered protocols work – message exchange

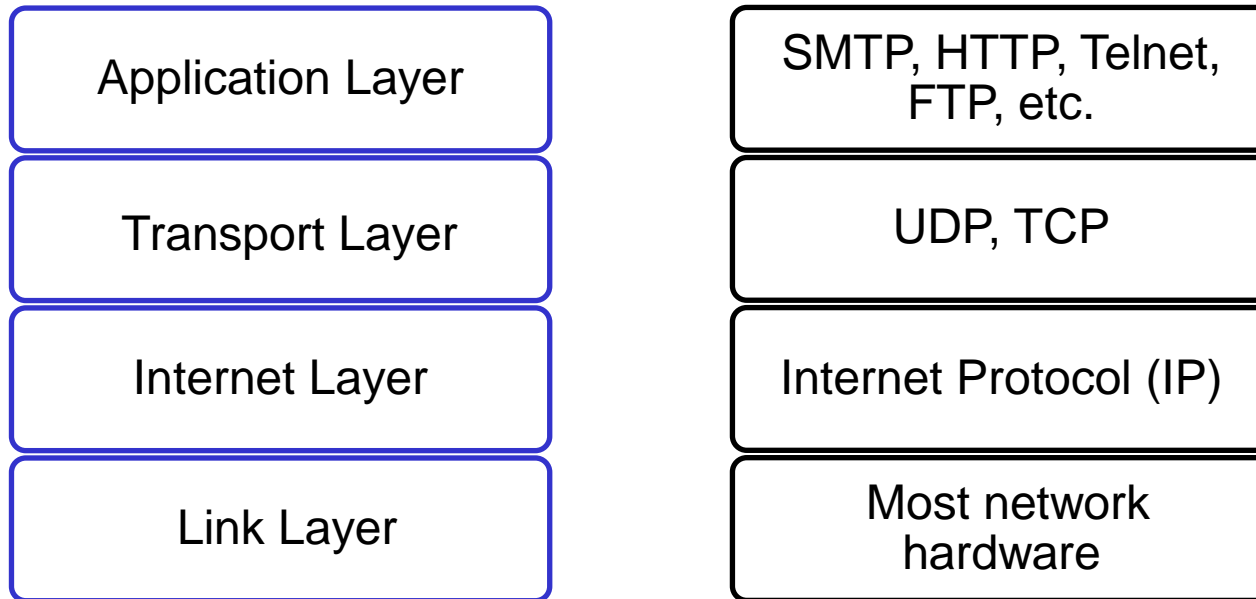- Principles of constructing messages/data encapsulation



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

# Examples of Layered Protocols



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| **Application Layer** | FTP | HTTP | ModBus | CoAP | BACNet | LonWorks | KNX (Konnex) |
| **Presentation Layer** | | | | | BACNet | LonWorks | |
| **Session Layer** | | | | | | LonWorks | |
| **Transport Layer** | TCP | UDP | RUDP | SCTP | | LonWorks | KNX (Konnex) |
| **Network Layer** | IP | | | BACNet | | LonWorks | |
| **Data Link Layer** | Ethernet CSMA/CD | CSMA/CA | | P-persistent CSMA | | | |
| **Physical Layer** | IEEE 802.3 (Ethernet) | IEEE 802.11 (Wifi) | ATM | | ZigBee | X.21 | PL-20 | TP |

# TCP/IP

- The most popular protocol suite used in the Internet

- Four layers

Protocol suite

| Application Layer | SMTP, HTTP, Telnet, FTP, etc. |
| Transport Layer | UDP, TCP |
| Internet Layer | Internet Protocol (IP) |
| Link Layer | Most network hardware |

http://tools.ietf.org/html/rfc1122

# Internet Protocol (IP)

- Defines the datagram as the basic data unit

- Defines the Internet address scheme

- Transmits data between the Network Access Layer and Transport Layer

- Routes datagrams to destinations
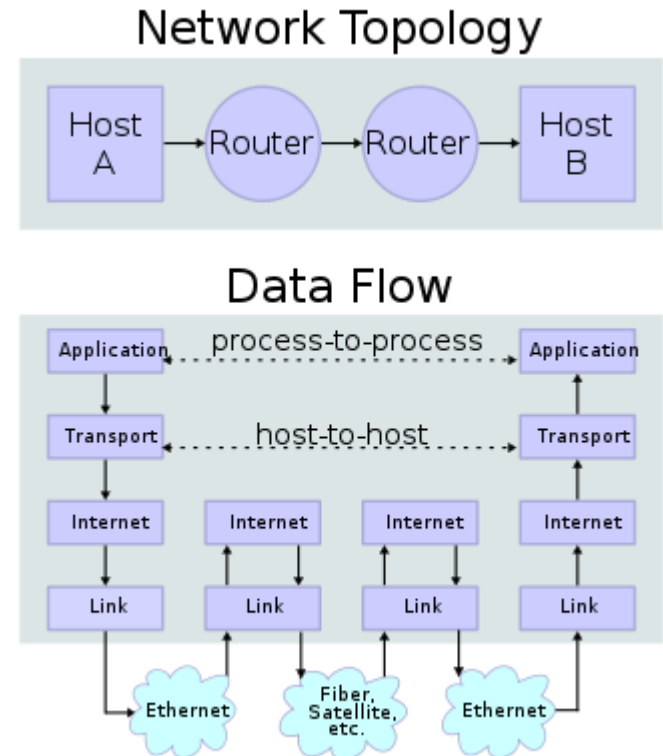
- Divides and assembles datagrams



Figure source:
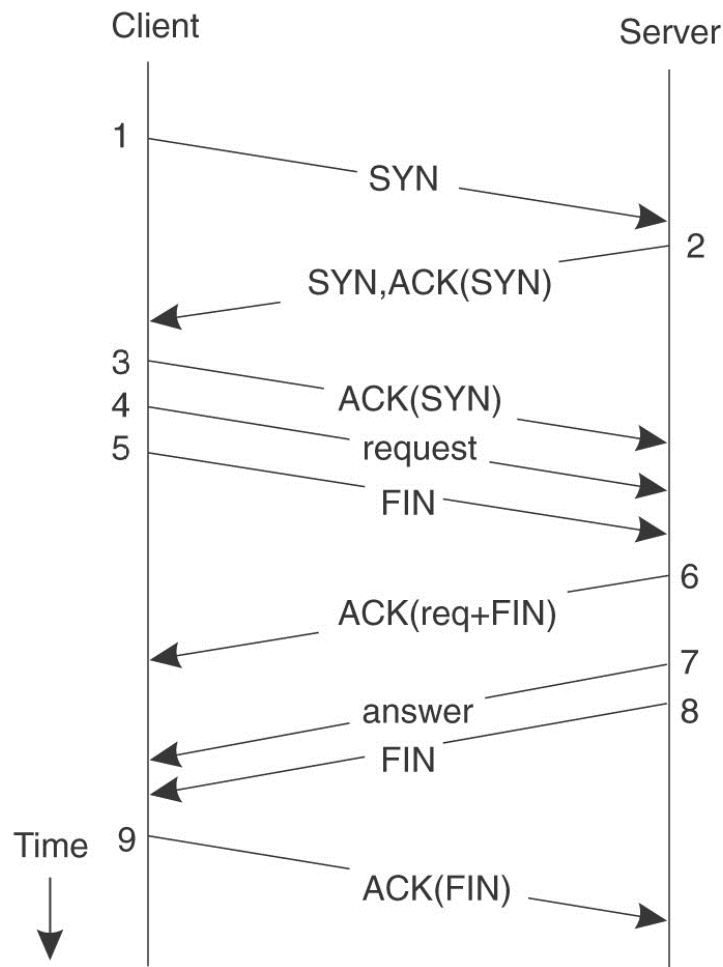http://en.wikipedia.org/wiki/Internet_protocol_suite

# TCP/IP – Transport Layer

- Host-to-host transport features
- Two main protocols: TCP (Transmission Control Protocol) and UDP (User Datagram Protocol)

| Layer\Protocol | TCP | UDP |
| --- | --- | --- |
| Application layer | Data sent via Streams | Data sent in Messages |
| Transport Layer | Segment | Packet |
| Internet Layer | Datagram | Datagram |
| Link Layer | Frame | Frame |

Note: pay attention with the terms „packet/datagram" in TCP/IP versus that in the OSI model

Distributed Systems, Bachelor Study          37

# TCP operations

$sudo nast -d -T iptest >ip.out

$wget www.tuwien.ac.at

Client | Server

1
SYN → 2
SYN,ACK(SYN)
3
4   ACK(SYN)
5   request
FIN
6
ACK(req+FIN)
7
answer   8
FIN
Time   9
ACK(FIN)

Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2002, Prentice-Hall, Inc.

```
---[ TCP ]--------------------------------------------------
192.168.1.7:46023(unknown) -> 128.130.35.76:80(http)
TTL: 64   Window: 14600     Version: 4          Length: 60
FLAGS: -S-----       SEQ: 3308581872 - ACK: 0
Packet Number: 16

---[ TCP ]--------------------------------------------------
128.130.35.76:80(http) -> 192.168.1.7:46023(unknown)
TTL: 54   Window: 14480     Version: 4          Length: 60
FLAGS: -S--A--       SEQ: 3467332359 - ACK: 3308581873
Packet Number: 17

---[ TCP ]--------------------------------------------------
192.168.1.7:46023(unknown) -> 128.130.35.76:80(http)
TTL: 64   Window: 115        Version: 4          Length: 52
FLAGS: ----A--       SEQ: 3308581873 - ACK: 3467332360
Packet Number: 18

---[ TCP ]--------------------------------------------------
192.168.1.7:46023(unknown) -> 128.130.35.76:80(http)
TTL: 64   Window: 115        Version: 4          Length: 166
FLAGS: ---PA--       SEQ: 3308581873 - ACK: 3467332360
Packet Number: 19

---[ TCP Data ]---------------------------------------------

GET / HTTP/1.1


---[ TCP ]--------------------------------------------------
128.130.35.76:80(http) -> 192.168.1.7:46023(unknown)
TTL: 54   Window: 114        Version: 4          Length: 52
FLAGS: ----A--       SEQ: 3467332360 - ACK: 3308581987
Packet Number: 20

---[ TCP ]--------------------------------------------------
128.130.35.76:80(http) -> 192.168.1.7:46023(unknown)
TTL: 54   Window: 114        Version: 4          Length: 1500
FLAGS: ----A--       SEQ: 3467332360 - ACK: 3308581987
Packet Number: 21

---[ TCP Data ]---------------------------------------------

HTTP/1.1 200 OK
```
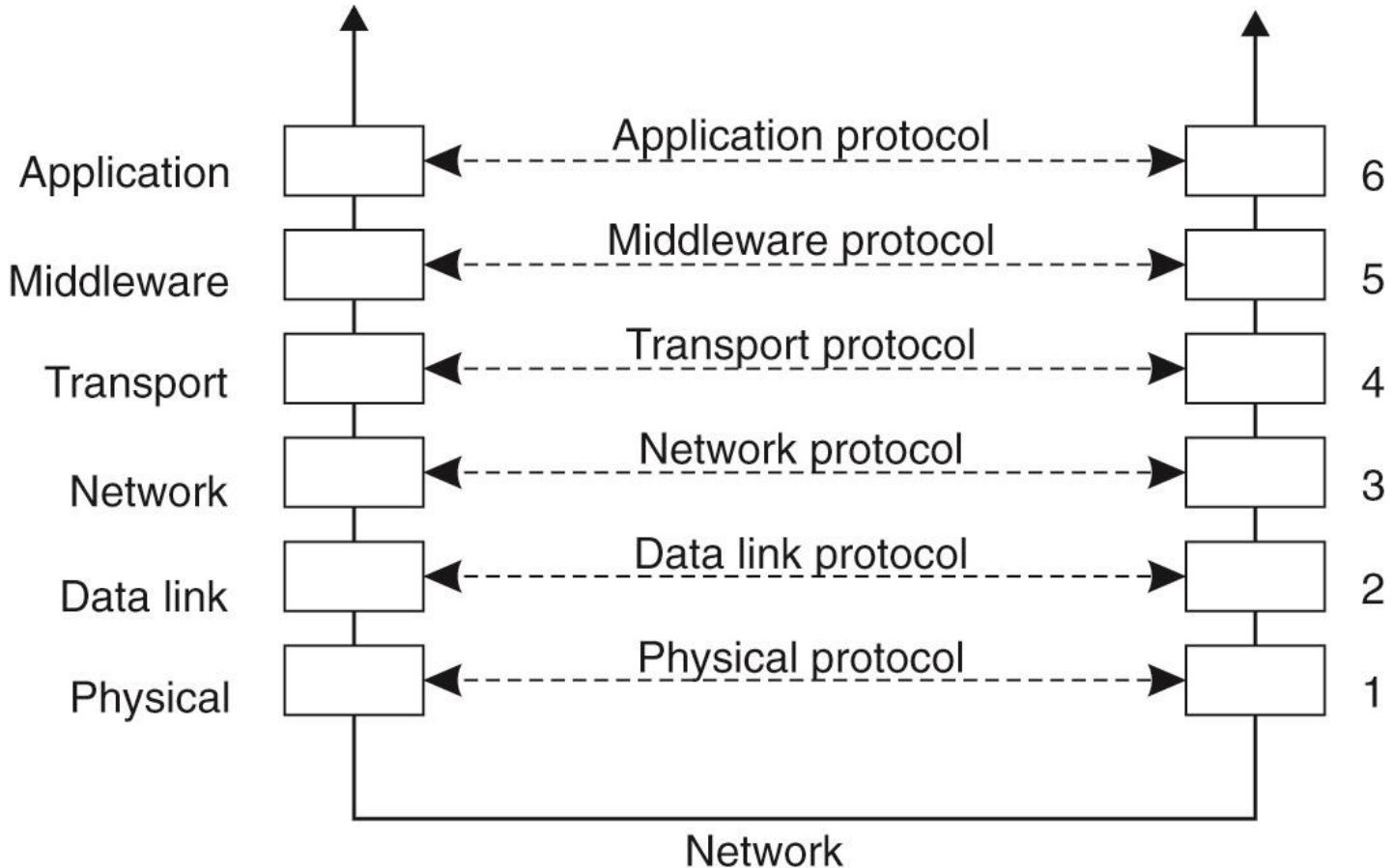
Distributed Systems, Bachelor Study

# Communication protocols are not enough

- We need more than just communication protocols
    - E.g., resolving names, electing a communication coordinator, locking resources, and synchronizing time

- Middleware
    - Including a set of general-purpose but application-specific protocols, middleware communication protocols, and other specific services.

# Middleware Protocols



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

# HANDLING COMMUNICATION MESSAGES/REQUESTS

# Where communication tasks take place?

- Message passing – send/receive
  - Processes send and receive messages
    - Sending process versus receiving process
    - Communication is done by using a set of functions for communication implementing protocols
- Remote method/procedure calls
  - A process calls/invokes a (remote) procedure in another process
    - Local versus remote procedure call, but in the same manner
- Remote object calls
  - A process calls/invokes a (remote) object in another process

Distributed Systems, Bachelor Study                42

# Basic send/receive communication

```
# Echo client program
import socket

HOST = 'daring.cwi.nl'    # The remote host
PORT = 50007          # The same port as
used by the server
s = socket.socket(socket.AF_INET,
socket.SOCK_STREAM)
s.connect((HOST, PORT))
s.send('Hello, world')
data = s.recv(1024)
s.close()
print 'Received', repr(data)
```

Network

```
# Echo server program
import socket

HOST = ''              # Symbolic name meaning the
local host
PORT = 50007          # Arbitrary non-privileged
port
s = socket.socket(socket.AF_INET,
socket.SOCK_STREAM)
s.bind((HOST, PORT))
s.listen(1)
conn, addr = s.accept()
print 'Connected by', addr
while 1:
    data = conn.recv(1024)
    if not data: break
    conn.send(data)
conn.close()
```

Python source: http://docs.python.org/release/2.5.2/lib/socket-example.html

# Remote procedure calls

```
void
hello_prog_1(char *host)
{
        CLIENT *clnt;
        char * *result_1;
        char *hello_1_arg;

#ifndef    DEBUG
        clnt = clnt_create (host, HELLO_PROG, HELLO_VERS, "udp");
        if (clnt == NULL) {
                clnt_pcreateerror (host);
                exit (1);
        }
#endif    /* DEBUG */

        result_1 = hello_1((void*)&hello_1_arg, clnt);
        if (result_1 == (char **) NULL) {
                clnt_perror (clnt, "call failed");
        }
#ifndef    DEBUG
        clnt_destroy (clnt);
#endif     /* DEBUG */
     printf("result is: %s\n",(*result_1));
}


int
main (int argc, char *argv[])
{
        char *host;

        if (argc < 2) {
                printf ("usage: %s server_host\n", argv[0]);
                exit (1);
        }
        host = argv[1];
        hello_prog_1 (host);
exit (0);
}
```

Network

## Procedure in a remote server

```
char **
hello_1_svc(void *argp, struct svc_req *rqstp)
{
        static char * result ="Hello";

        /*
         * insert server code here
         */

        return &result;
}
```

# Remote object calls

```java
public class ComputePi {
   public static void main(String args[]) {
      if (System.getSecurityManager() == null) {
         System.setSecurityManager(new SecurityManager());
      }
      try {
         String name = "Compute";
         Registry registry = LocateRegistry.getRegistry(args[0]);
         Compute comp = (Compute) registry.lookup(name);
         Pi task = new Pi(Integer.parseInt(args[1]));
         BigDecimal pi = comp.executeTask(task);
         System.out.println(pi);
      } catch (Exception e) {
         System.err.println("ComputePi exception:");
         e.printStackTrace();
      }
   }
}
```
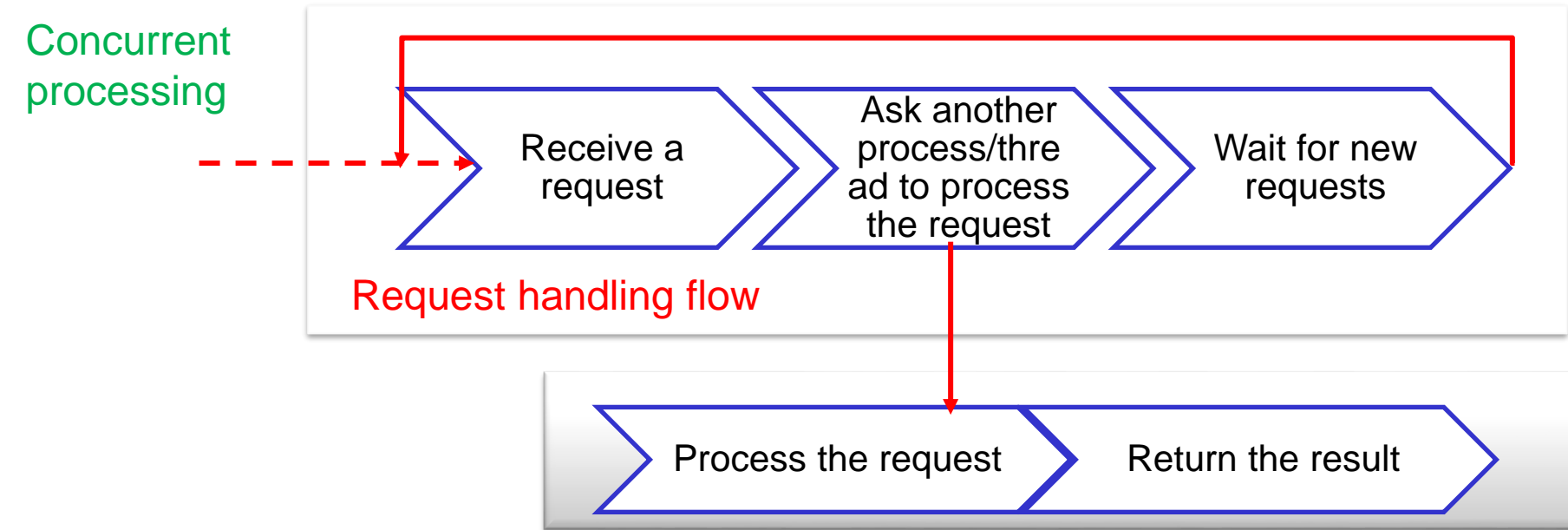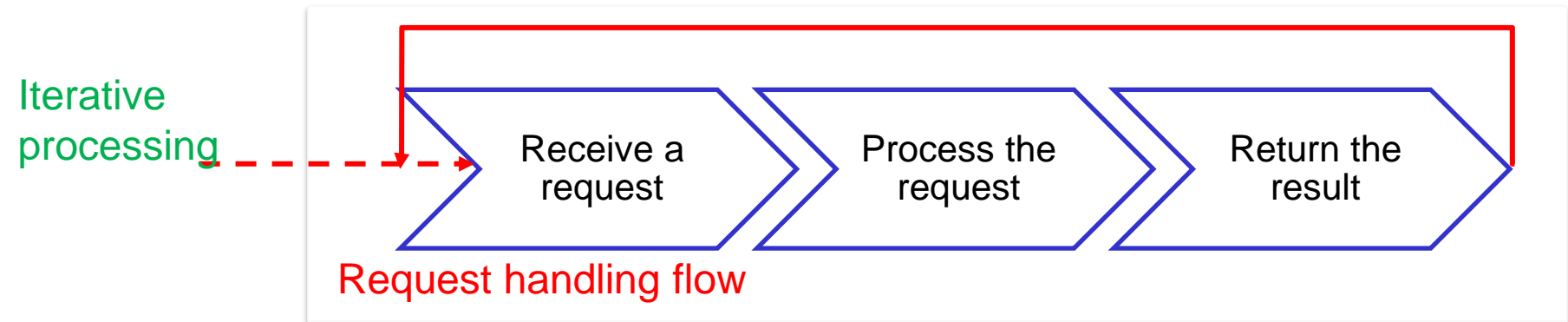
```java
public interface Compute extends Remote {
   <T> T executeTask(Task<T> t) throws RemoteException;
}
....
public class ComputeEngine implements Compute {

   public ComputeEngine() {
      super();
   }

   public <T> T executeTask(Task<T> t) {
      return t.execute();
   }

   public static void main(String[] args) {
      if (System.getSecurityManager() == null) {
         System.setSecurityManager(new SecurityManager());
      }
      try {
         String name = "Compute";
         Compute engine = new ComputeEngine();
         Compute stub =
            (Compute) UnicastRemoteObject.exportObject(engine, 0);
         Registry registry = LocateRegistry.getRegistry();
         registry.rebind(name, stub);
         System.out.println("ComputeEngine bound");
      } catch (Exception e) {
         System.err.println("ComputeEngine exception:");
         e.printStackTrace();
      }
   }
}
```

Java  source:
http://docs.oracle.com/javase/tutorial/rmi/overview.html

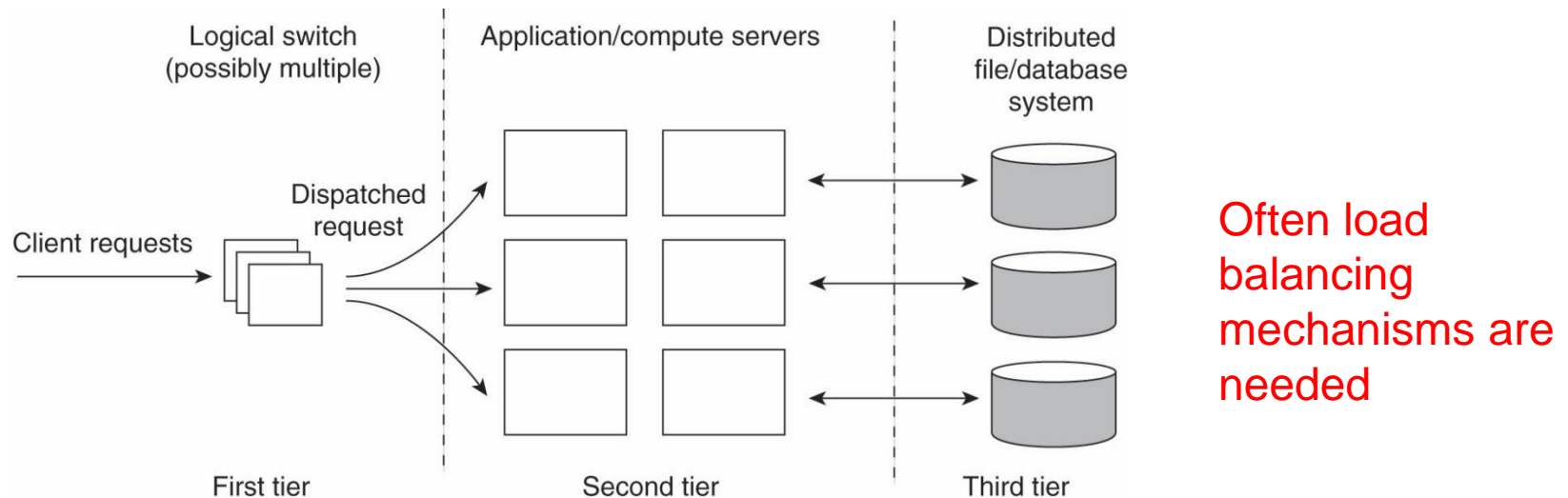Distributed Systems, Bachelor Study

45

# Processing multiple requests

- How to deal with multiple, concurrent messages received?

- Problems:

    - Different roles: clients versus servers/services

        - A large <span style="color:red">number of clients</span> interact with <span style="color:red">a small number of servers/services</span>

        - A single process might receive a lot of messages at the same time

- Impacts

    - performance, reliability, cost, etc.

# Iterative versus concurrent processing

**Iterative processing**

Receive a request → Process the request → Return the result

Request handling flow

---

**Concurrent processing**

Receive a request → Ask another process/thread to process the request → Wait for new requests

Request handling flow

Process the request → Return the result

# Using replicated processes



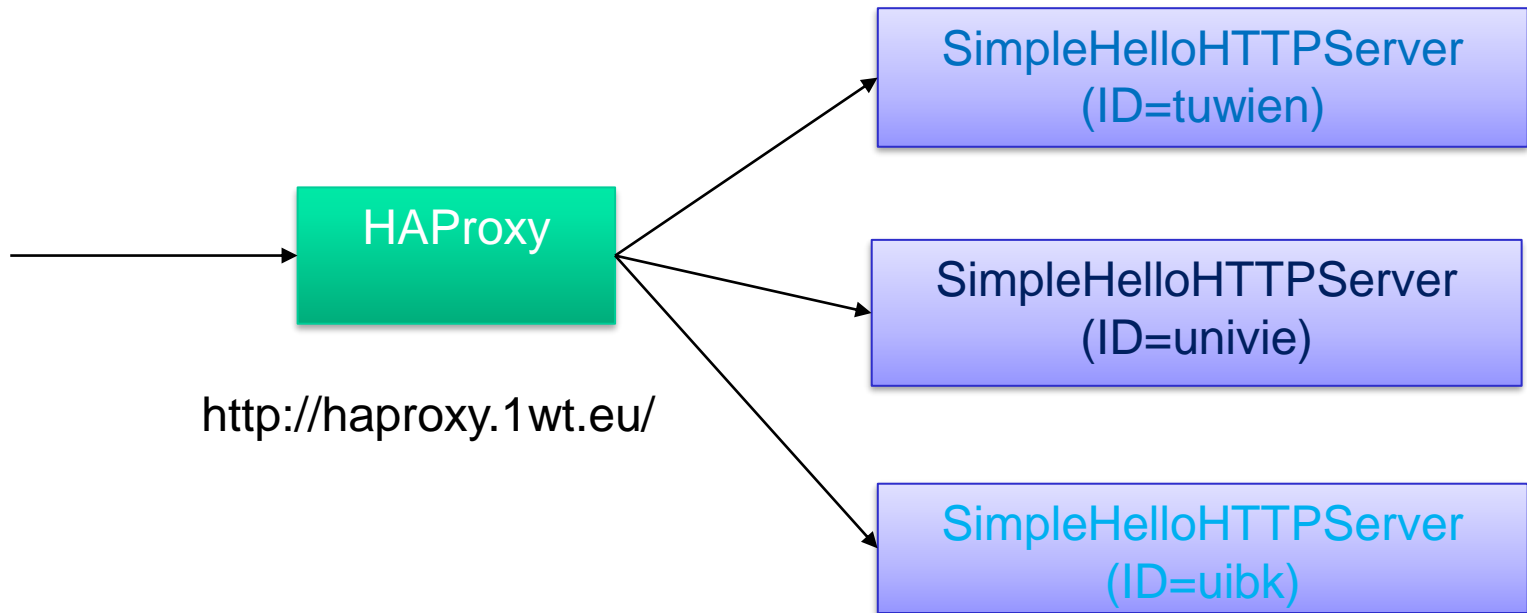Often load balancing mechanisms are needed

Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall

Q: How does this model help to improve performance and fault-tolerance? What would be a possible mechanism to reduce costs based on the number of client requests?
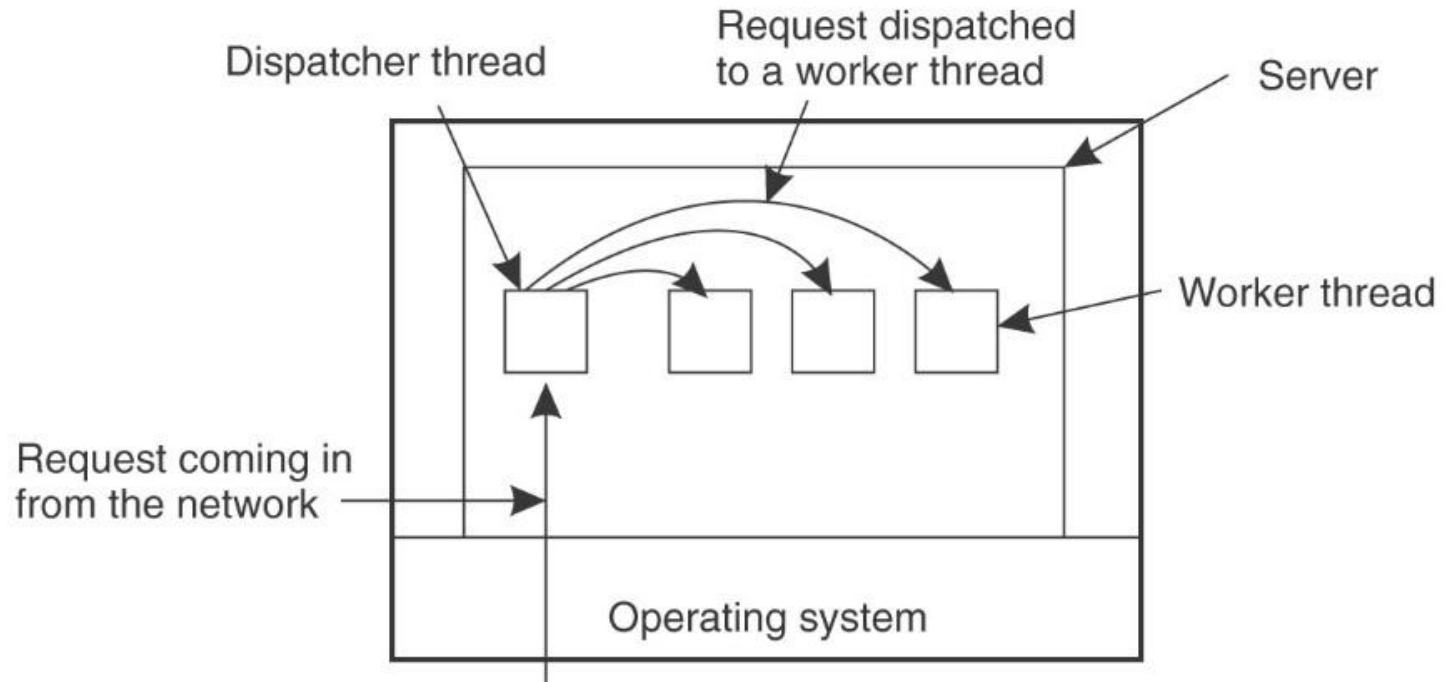
# Example



- Get a small test
  ▪Download haproxy, e.g.
  $sudo apt-get install haproxy
  - Download SimpleHelloHTTPServer.java and haproxy configuration
    - http://bit.ly/19xFDRC
  - Run 1 haproxy instance and 3 http servers
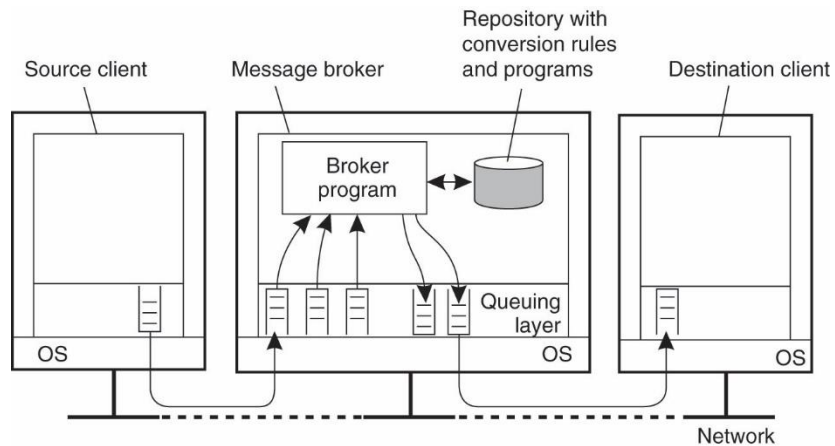    - Modify configuration and parameters if needed
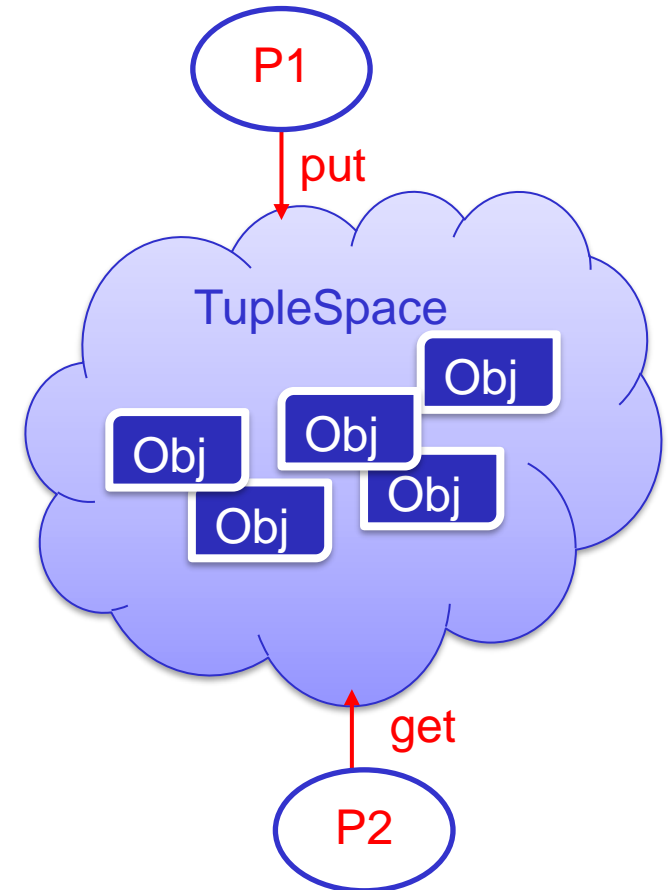  - Run a test client

# Using multiple threads



Dispatcher thread

Request dispatched to a worker thread

Server

Worker thread

Request coming in from the network

Operating system

Q: Compare this architectural model with the super-server model?

# Using message brokers/space repository



Source: Andrew S. Tanenbaum and Maarten van Steen, Distributed Systems – Principles and Paradigms, 2nd Edition, 2007, Prentice-Hall



P1

put

TupleSpace

Obj
Obj
Obj
Obj
Obj

get

P2

# Example

- Get a free instance of RabbitMQ from cloudamqp.com
- Get code from: https://github.com/cloudamqp/java-amqp-example
- First run the test sender, then run the receiver

Test sender → RabbitMQ ⇢ Test receiver

Messaging that just works

cloudamqp.com

```
channel.queueDeclare(QUEUE_NAME, false, false, false, null);
    for (int i=0; i<100; i++) {
        String message = "Hello distributed systems guys:  "+i;
        channel.basicPublish("", QUEUE_NAME, null, message.getBytes());
        System.out.println(" [x] Sent '" + message + "'");
        new Thread().sleep(5000);
    }
```

```
while (true) {
    QueueingConsumer.Delivery delivery = consumer.nextDelivery();
    String message = new String(delivery.getBody());
    System.out.println(" [x] Received '" + message + "'");
  }
```

Note: i modified the code a bit

# Summary

- Complex and diverse communication patterns, protocols and processing models

- Choices are based on communication requirements and underlying networks

  - Understand their pros/cons

  - Understand pros and cons of their technological implementations

- Dont forget to play with some simple examples to understand existing concepts

# Thanks for your attention