



Distributed Systems



Giuseppe Anastasi
g.anastasi@iet.unipi.it
Pervasive Computing & Networking Lab. (PerLab)
Dept. of Information Engineering, University of Pisa



Overview

- Introduction
- Network-Based Operating Systems
- Communication Networks
- Communication Protocols
- Distributed Programming



Distributed Systems 2 Operating Systems



Objectives

- Provide a high-level overview of distributed systems and the networks that interconnect them
- Discuss the general structure of operating systems for distributed systems
- Introduce different types of communication networks
- Discuss networking protocols that allow communication in a distributed environment
- Introduce principles of distributed programming

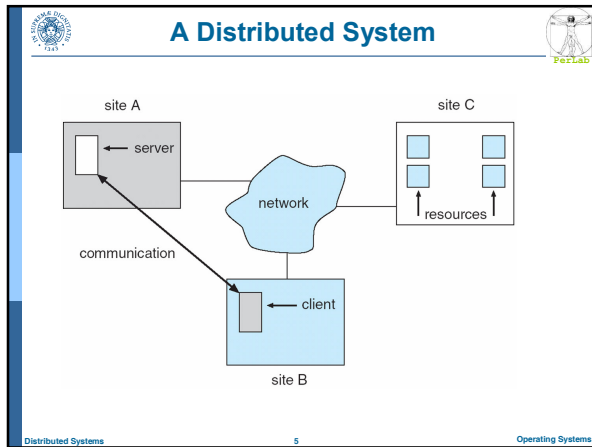
Distributed Systems 3 Operating Systems



Motivation

- **Distributed system** is collection of loosely coupled processors interconnected by a communications network
- Processors variously called *nodes*, *computers*, *machines*, *hosts*
 - Site is location of the processor
- Reasons for distributed systems
 - Resource sharing
 - ▶ sharing and printing files at remote sites
 - ▶ processing information in a distributed database
 - ▶ using remote specialized hardware devices
 - Computation speedup – **load sharing**
 - Reliability – detect and recover from site failure, function transfer, reintegrate failed site
 - Communication – message passing



Distributed Systems 4 Operating Systems



Overview



- Introduction
- **Network-based Operating Systems**
- Communication Networks
- Communication Protocols
- Distributed Programming

Distributed Systems 6 Operating Systems

 **Operating Systems for Distributed Systems** 



- Network Operating Systems
- Distributed Operating Systems

Distributed Systems 7 Operating Systems

 **Network-Operating Systems** 

- Users are aware of multiplicity of machines.
- Access to resources of various machines is done explicitly by:
 - Remote logging into the appropriate remote machine (telnet, ssh)
 - Remote Desktop (Microsoft Windows)
 - Transferring data from remote machines to local machines, e.g., via File Transfer Protocol (FTP) or Web (HTTP)

Distributed Systems 8 Operating Systems

 **Distributed-Operating Systems** 

- Users not aware of multiplicity of machines
 - Access to remote resources similar to access to local resources
- Data Migration
 - transfer data by transferring entire file, or only those portions of the file necessary for the immediate task
- Computation Migration
 - transfer the computation, rather than the data, across the system
- Process Migration
 - execute an entire process, or parts of it, at different sites

Distributed Systems 9 Operating Systems

Distributed-Operating Systems (Cont.)

- **Load balancing**
 - distribute processes across network to even the workload
- **Computation speedup**
 - Sub-processes can run concurrently on different sites
- **Hardware preference**
 - process execution may require specialized processor
- **Software preference**
 - required software may be available at only a particular site
- **Data access**
 - run process remotely, rather than transfer all data locally

Distributed Systems 10 Operating Systems

Overview

- Introduction
- Network-Based Operating Systems
- **Communication Networks**
 - LAN
 - WAN
 - Internet
- Communication Protocols
- Distributed Programming

Distributed Systems 11 Operating Systems

Classification of Communication Networks

- Wide Area Networks (WAN)
- Metropolitan Area Networks (MAN)
- Local Area Networks (LAN)
- Personal Area Networks (PAN)
- Body Area Networks (BAN)

Distributed Systems 12 Operating Systems

Network Topologies

fully connected network

partially connected network

tree-structured network

star network

ring network

Distributed Systems 13 Operating Systems

Service Types

- **Connection Oriented (Stream)**
 - Inspired from the telephone system
 - Connection Setup, Data Transfer, Connection Tear down
 - Messages tend to follow the same path from source to destination
- **Connectionless (Datagram)**
 - Inspired from the mail system
 - Each message includes the destination address
 - Different messages follows different paths
 - No guarantee on message ordering

Distributed Systems 14 Operating Systems

Layered Architecture

OSI	TCP/IP
Application ← LAYER 7	Application ← LAYER 5
Presentation ← LAYER 6	Transport ← LAYER 4
Session ← LAYER 5	Internet ← LAYER 3
Transport ← LAYER 4	Network Interface ← LAYER 2
Network ← LAYER 3	Physical ← LAYER 1
Data Link ← LAYER 2	
Physical ← LAYER 1	

Distributed Systems 15 Operating Systems

Physical Addressing Scheme

How to implement unicast communication?

↓

Each node in the LAN is assigned with a unique address (physical address or hw address or MAC address)

The sending node inserts the destination address in each packet it sends

The destination node accepts the received packet *only if* the destination address corresponds to its own address.

Distributed Systems 19 Operating Systems

Ethernet LANs

- Ethernet a 10 Mbps
 - 10BaseT
 - 10Base5
 - 10Base2
- Fast Ethernet (100 BaseT, 100 Mbps)
- Gigabit Ethernet (1, 10 Gbps)
 - ▶ IEEE 802.3z (1 Gbps)
 - ▶ IEEE 802.3ae (10 Gbps)

Distributed Systems 20 Operating Systems

An Example of Ethernet LAN

Figure 5.29 ♦ An institutional network using a combination of hubs, Ethernet switches, and a router

Distributed Systems 21 Operating Systems

Frame Ethernet

- Preamble
 - 64-bit string used for clock synchronization
- Destination Address
 - Broadcast address: 11111.....1
- Source Address
- Type
 - Specify the data type (e.g., 0800: IPv4)
- Payload
- CRC (Cyclic Redundancy Check)
 - Error detection

Distributed Systems 22 Operating Systems

Ethernet MAC Protocol

- **CSMA/CD** - Carrier sense with multiple access (CSMA); collision detection (CD)
 - A node determines whether another packet is currently being transmitted over that link (*carrier sense*).
 - If the link is sensed as free the packet transmission is started. The node continues listening while transmitting
 - If two or more nodes begin transmitting at exactly the same time, then they will register a *collision* and will stop transmitting
 - A collided packet is re-tried after a random backoff interval

Distributed Systems 23 Operating Systems

Wide Area Networks (WANs)

- Links geographically separated sites
- Point-to-point connections over long-haul lines (often leased from a phone company)
- Speed \approx 1.544 – 45 Mbps

Distributed Systems 24 Operating Systems

Point-to-Point Physical Links

- Collegamenti via modem
- Linee dedicate
- Linee ISDN
- Linee DSL
 - ▶ ADSL (Asynchronous Digital Subscriber Line)
 - ▶ SDSL (Synchronous Digital Subscriber Line)
 - ▶ HDSL (High-Rate Digital Subscriber Line)
 - ▶ VDSL (Very high-rate Digital Subscriber Line)

Distributed Systems 25 Operating Systems

Point-to-Point Virtual Links

- Frame Relay
- SMDS
- ATM

CP: usually mainframes
Broadcast usually requires multiple messages

Distributed Systems 26 Operating Systems

Limitazioni delle LAN

- Le LAN permettono di interconnettere calcolatori in ambito locale
 - Ambienti SOHO
 - Edifici
 - Campus

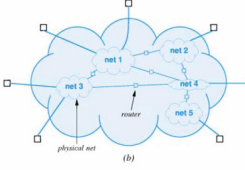
Come far comunicare fra loro calcolatori collegati a reti di tipo diverso e in posti diversi?

Distributed Systems 27 Operating Systems

Quello che vorremmo ...

Due calcolatori devono poter comunicare indipendentemente dalla rete fisica a cui sono collegati

- Aumento di produttività
- Problemi da risolvere
 - Diversi segnali elettrici,
 - Diverso formato dei pacchetti
 - Diverso schema di indirizzamento

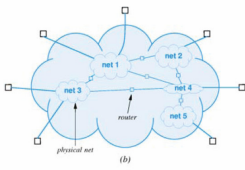


Distributed Systems 28 Operating Systems

Inter-rete

Occorre

- Raccordare fisicamente le reti fisiche mediante opportuni dispositivi fisici (router)
- Aggiungere uno strato software sopra l'hardware di rete
 - Fa apparire l'insieme di reti eterogenee come un unico sistema




Il sistema che si realizza è detto inter-rete o internet

Distributed Systems 29 Operating Systems

Router

- Calcolatore specializzato dedicato alla interconnessione
 - CPU, memoria, ...
 - una interfaccia per ciascuna rete a cui è collegato
- Si collega come un qualsiasi altro calcolatore
 - Collegato contemporaneamente a (almeno) due reti fisiche
- Può interconnettere reti di tipo qualsiasi

Inter-rete = reti di calcolatori collegate da router



Distributed Systems 30 Operating Systems

Protocolli software

- Simulano una rete virtuale
 - si può collegare un calcolatore come si farebbe con una rete singola
- Nascondono i dettagli delle reti fisiche sottostanti
 - L'utente si può disinteressare del tipo di reti fisiche sottostanti, della presenza o meno di router, ecc.
- Realizzano un servizio universale
 - Ogni calcolatore è individuato tramite un **indirizzo software**
 - Ogni calcolatore può scambiare messaggi con altri calcolatori collegati alla inter-rete

Distributed Systems 31 Operating Systems

Protocolli TCP/IP

TCP = Transmission Control Protocol
IP = Internet Protocol



- Famiglia di protocolli usati in **Internet**
- Usati anche per la realizzazione di inter-reti private (Intranet)
- Progettati verso gli inizi degli anni '70 su iniziativa del Pentagono
 - Agenzia ARPA → Arpanet
 - Arpanet → Internet

Distributed Systems 32 Operating Systems

The TCP/IP Protocol Layers



ISO	TCP/IP
application	HTTP, DNS, Telnet SMTP, FTP
presentation	not defined
session	not defined
transport	TCP-UDP
network	IP
data link	not defined
physical	not defined

Distributed Systems 33 Operating Systems

 **Overview** 



- Introduction
- Network-Based Operating Systems
- Communication Networks
- **Communication Protocols**
 - TCP/UDP
 - IP
- Distributed Programming

Distributed Systems 34 Operating Systems

 **Internet Layer (Protocollo IP)** 

- Formato dei pacchetti (datagram)
- Formato degli indirizzi software (indirizzi IP)
- Intradamento dei datagrammi
- Servizio best-effort
 - Connectionless
 - ▶ Possibili fuori sequenza
 - Non affidabile
 - ▶ Perdite e/o alterazioni dei datagram
 - ▶ Nessuna garanzia di QoS (ritardo, jitter, throughput)

Distributed Systems 35 Operating Systems

 **Transport Layer (Protocollo UDP)** 

- Demultiplexing dei datagram
 - Riceve un flusso indistinto di datagrammi IP
 - Recapita i datagram ai processi applicativi a cui sono destinati
- Nessun incremento al servizio offerto da IP
 - Servizio connectionless e non affidabile

Distributed Systems 36 Operating Systems

Transport Layer (Protocollo TCP)

- Flusso di byte (stream)
 - ma la comunicazione è sempre a pacchetti (segmenti)
- Trattamento di fuori-sequenza e duplicati
- Rilevazione dei segmenti alterati o persi
- Recupero dei segmenti alterati, persi, ritardati
- Controllo del flusso
- Controllo della congestione
- Servizio **Connection-oriented e affidabile**
 - **Tutti** i segmenti vengono consegnati in sequenza
 - Assenza di duplicati
 - Nessuna garanzia sul ritardo, sul jitter e sul throughput

Distributed Systems 37 Operating Systems

Host, Router e Protocolli

Host = qualsiasi calcolatore collegato alla inter-rete

Distributed Systems 38 Operating Systems

Indirizzi IP

- Indirizzo a 32 bit assegnato a ogni host
- Struttura Gerarchica
 - Indirizzo di rete (prefisso) + Indirizzo di host (suffisso)
- Indirizzo di rete (**network number**)
 - Identifica una rete fisica
 - Assegnato da una autorità centrale che garantisce l'univocità

Indirizzo di host (**host number**)

- Identifica un particolare host all'interno della rete fisica
- Assegnato localmente dall'amministratore

Distributed Systems 39 Operating Systems

Classi di indirizzi IP

bits 0 1 2 3 4 8 16 24 31

Class A 0 prefix suffix

Class B 10 prefix suffix

Class C 110 prefix suffix

Class D 11110 multicast address

Class E 11111 reserved for future use

Distributed Systems 40 Operating Systems

Notazione decimale puntata

- I 4 byte sono interpretati come numeri decimali
 - ↳ compresi fra 0 e 255
- Indirizzo letto come 4 numeri decimali separati da punti

32-bit Binary Number	Equivalent Dotted Decimal
10000001 00110100 00000110 00000000	129 . 52 . 6 . 0
11000000 00000101 00110000 00000011	192 . 5 . 48 . 3
00001010 00000010 00000000 00100101	10 . 2 . 0 . 37
10000000 00001010 00000010 00000011	128 . 10 . 2 . 3
10000000 10000000 11111111 00000000	128 . 128 . 255 . 0

Distributed Systems 41 Operating Systems

Classi e notazione puntata

Classe	Range di Valori
A	da 0 a 127
B	da 128 a 191
C	da 192 a 223
D	da 224 a 239
E	da 240 a 255

Distributed Systems 42 Operating Systems

Indirizzi IP riservati

- Indirizzo per l'intera sottorete
 - Suffisso di tutti zeri (131.114.0.0)
- Indirizzo di trasmissione broadcast orientata
 - Suffisso di tutti 1 (131.114.255.255)
- Indirizzo di trasmissione broadcast ristretta
 - Costituito da tutti 1 (255.255.255.255)
- Indirizzo di "questo calcolatore"
 - Indirizzi di tutti 0
 - Usato all'avvio
- Indirizzo loopback
 - 127.0.0.1
 - Usato nello fase di sviluppo di applicazione di rete

Distributed Systems 43 Operating Systems

Datagram IP

0	4	8	16	19	24	31
VERS		H. LEN		SERVICE TYPE		TOTAL LENGTH
IDENTIFICATION			FLAGS	FRAGMENT OFFSET		
TIME TO LIVE		TYPE		HEADER CHECKSUM		
SOURCE IP ADDRESS						
DESTINATION IP ADDRESS						
IP OPTIONS (MAY BE OMITTED)				PADDING		
BEGINNING OF DATA						
⋮						

Distributed Systems 44 Operating Systems

Livello di trasporto

- Estende il servizio di trasporto host-to-host in un servizio di comunicazione fra processi
 - Esegue il demultiplexing delle informazioni
 - Basato sul concetto di porta
- I processi vengono individuati mediante la coppia

<Host IP Address, Port Number>
- Il processo mittente deve specificare sia l'indirizzo IP che il numero di porta
- Il SO realizza la porta come coda di messaggi

Distributed Systems 45 Operating Systems

Protocollo UDP

The diagram illustrates the flow of data through the UDP and IP protocols. At the top, three processes labeled A, B, and C are shown. Each process sends data packets (represented by blue rectangles) to a central layer labeled 'Protocollo UDP'. This layer then forwards the packets to a lower layer labeled 'Protocollo IP'. The packets are shown in a queue, with B, C, and A in order from top to bottom. The slide includes logos for the University of Pisa and PerLab, and footer text: 'Distributed Systems', '46', and 'Operating Systems'.

Introduzione al TCP

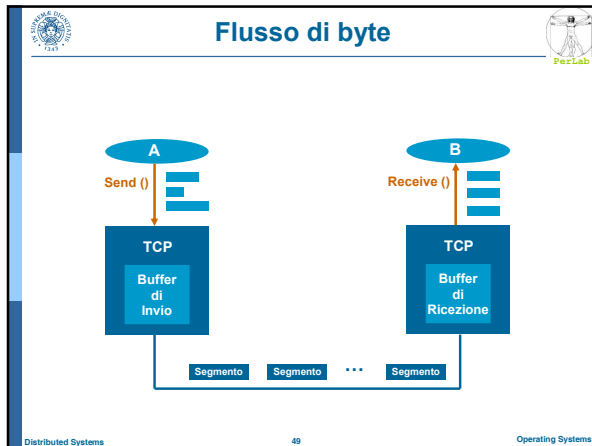
- Servizio di comunicazione fra processi orientato alla connessione
 - Meccanismo di de-multiplexing basato sulle porte
 - Apertura connessione, trasferimento dati, chiusura connessione
- Punto-Punto
 - Ogni connessione TCP collega esattamente due processi
- Affidabile
 - Consegna affidabile, senza duplicati e in sequenza di un flusso di byte
- Full duplex
 - Un flusso di byte in ciascuna direzione

The slide includes logos for the University of Pisa and PerLab, and footer text: 'Distributed Systems', '47', and 'Operating Systems'.

Introduzione al TCP

- Controllo e gestione degli errori
 - Checksum e ACK
 - Timeout e ritrasmissione
- Controllo del flusso
 - Evita che il mittente invii più dati di quanti il ricevitore possa gestire
- Controllo della congestione
 - Evita che il mittente invii più dati di quanti la rete sia in grado di trasportare

The slide includes logos for the University of Pisa and PerLab, and footer text: 'Distributed Systems', '48', and 'Operating Systems'.



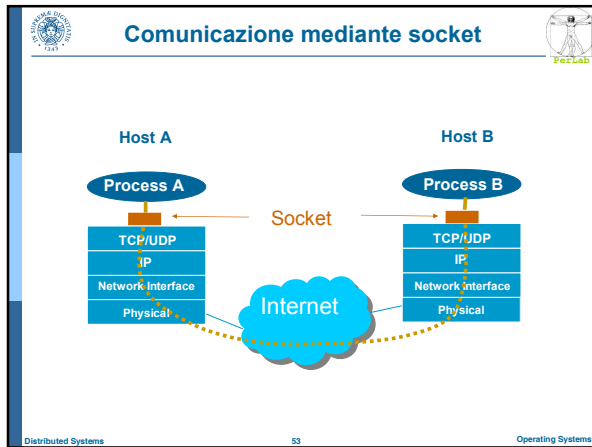
- ### Overview
- Introduction
 - Network-Based Operating Systems
 - Communication Networks
 - Communication Protocols
 - Distributed Programming
 - Socket interface
 - Client-server Model
 - Peer-to-Peer Model
- Distributed Systems 50 Operating Systems

- ### Interfaccia Socket
- Standard *de facto* per la comunicazione fra processi in ambiente distribuito
 - Si può usare anche per la comunicazione fra processi sulla stessa macchina
 - Interfaccia unica per operare con i vari protocolli di rete a disposizione
 - Nasconde tutti i meccanismi di comunicazione di livello inferiore (TCP/UDP, IP, ...)
- Distributed Systems 51 Operating Systems

Socket

- Estremità di canale di comunicazione identificata da un indirizzo
 - ▶ Socket: presa telefonica
 - ▶ Indirizzo: numero di telefono
- Indirizzo
 - ▶ Indirizzo dell'Host (*Indirizzo IP*)
 - ▶ Indirizzo del processo (*Port Number*)
- La comunicazione avviene tramite una coppia di socket

Distributed Systems 52 Operating Systems



Supporto del SO

- Il SO implementa l'astrazione di socket
- System call per
 - Creare un socket
 - Associare indirizzo IP e porta al socket
 - Mettere in ascolto un processo su un socket (server)
 - Accettare una richiesta di servizio su un socket (server)
 - Aprire una connessione verso un socket remoto (client)
 - Inviare un messaggio verso un socket remoto
 - Ricevere un messaggio da un socket
 -

Distributed Systems 54 Operating Systems

Modello Client-Server

- Paradigma basato su scambio di messaggi
 - Paradigma generale
 - Ma usato principalmente in ambito distribuito
- Scambio di msg per
 - Richiesta di servizio
 - Invio dei risultati

Distributed Systems 55 Operating Systems

Client-Server in Sistemi Distribuiti

Distributed Systems 56 Operating Systems

Modello P2P

Distributed Systems 57 Operating Systems

Client-Server Communication

host X
(146.86.5.20)

socket
(146.86.5.20:1625)

web server
(161.25.19.8)

socket
(161.25.19.8:80)

Distributed Systems 58 Operating Systems

Remote Procedure Calls

- Remote procedure call (RPC) abstracts procedure calls between processes on networked systems
- **Stubs** – client-side proxy for the actual procedure on the server
- The client-side stub locates the server and *marshalls* the parameters
- The server-side stub receives this message, unpacks the marshalled parameters, and performs the procedure on the server

Distributed Systems 59 Operating Systems

RPC: general issues

- Data Representation
 - Little endian vs. big endian
 - A system-independent representation is used (e.g., XDR: eXternal Data Representation)
- Exactly-once semantic
 - Acks and retransmissions for avoid message losses (at least once)
 - Timestamps for avoiding multiple execution (at most once)

Distributed Systems 60 Operating Systems

RPC: general issues

- Client-server communication
 - How to locate the RPC port on the server?
- Predefined ports
 - RPC are associated at compile time with fixed port numbers
 - The server cannot change the port number of the required service
- Rendez-Vous
 - The server-side OS provides a rendez-vous daemon (*matchmaker*)
 - The client requires the port number to the matchmaker
 - The daemon replies with the port number
 - The client send the RPC request to the appropriate port number

Distributed Systems 61 Operating Systems

Execution of RPC

Distributed Systems 62 Operating Systems

An RPC Application

- Distributed File System (DFS)
 - Set of daemons and RPC clients
 - Messages containing file-system operations
 - read, write, rename, delete, status
 - The client sends a message to the server
 - The command is executed on the server
 - A reply message is sent to the client

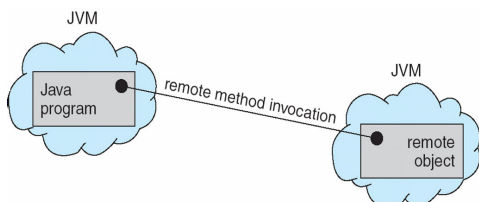
Distributed Systems 63 Operating Systems



Remote Method Invocation



- Remote Method Invocation (RMI) is a Java mechanism similar to RPCs
- RMI allows a Java program on one machine to invoke a method on a remote object





Questions?