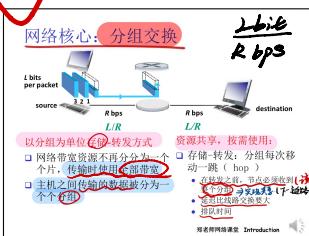


△ Network Core

~~packet switching~~

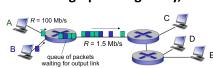
Lecture 1 Part



store and forward

⇒ queuing delay, loss

Packet Switching: queuing delay, loss



- * If arrival rate (in bits/s) exceeds transmission rate of link for a period of time:
- * packets will wait to be transmitted on link
- * packets can be dropped (lost) if memory/buffer fills up

丢包情况
retransmitted by previous
source and system
not at all

packet switching allows more users to use network.

~~circuit switching~~

P55

Packet Switching特点

优点:

- resource sharing 共享网络资源
- simpler, no call setup 想连就连
- 缺点:
- excessive congestion possible: delay and loss
- protocols needed for reliable data transfer, congestion control

还需下层提供可靠数据传输和拥塞控制等服务

How to provide circuit-like behavior PS?

• Bandwidth guarantees 对带宽要求很高 带宽要求

• New methods should be developed

~~circuit switching~~

P55

Circuit switching

线路交换的优点:

1. Dedicated resources: no sharing

2. circuit-like (guaranteed) performance

- Circuit segment is idle if not used by call (no sharing)
- Commonly used in traditional telephone networks

应用方式:

FDM: 按频率划分给多条线路中每一条各占一个频率区间

TDM: 按时间块划分给多条线路中的一条使用, 在终点可根据顺序复原合成完整数据

计算举例

在一条电路交换网络上, 从主机A到主机B发送一个40,000比特的文件需要多长时间?

- 所有的链路速率是1.5Mbps
- 每条链路使用TDM为4/4的TDM
- 建立端-to-end的电连接500 ms

每条链路的速率 (一个时间片): 1.5Mbps/24=64kbps

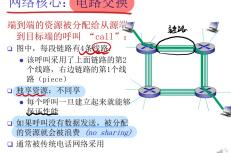
传输时间: 640kb / 64kbps = 10s

共用时间: 传输时间 + 建立链路时间 = 10s + 500ms = 10.5s

把所有链路都加起来

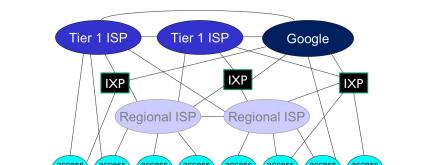
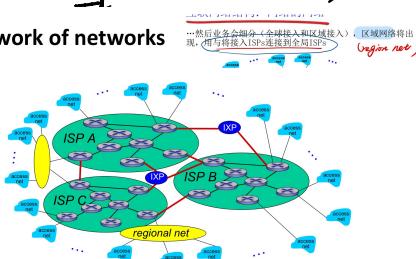
消耗的时间为10.5s

网络核心: (电路交换)



端系统通过 ISP 连入 Internet → connect each access ISP to one global transit ISP

Network of networks



- Small view of well-connected large networks
- "tier-1" commercial ISPs (e.g., ChinaTel, Sprint, AT&T, NTT), national & international coverage
- content provider network (e.g., Google): private network that connects its data centers to Internet, often bypassing tier-1, regional ISPs

四种延迟.

Lecture 1 P68

Delay in Packet-Switched Networks.

延迟通常由四部分组成：

..Transmission delay 传输延迟，指将整个分组packet

完整的从路由器打出的时间

? Propagation delay 传播延迟，指在链路上的客观物理时间

? Queuing delay 排队延迟，存在于过多用户占用带宽

超过了链路带宽上限的情况，涉及到流量强度的问题

i. Nodal processing delay 路由处理延迟，通常不会问

可以当没有

吞吐量.

Lecture 1 P72

Throughput 吞吐量

- Throughput: rate (bits/time unit) at which bits transferred between sender/receiver 由多段链路中最小的带宽来决定

R19. 假定主机 A 要向主机 B 发送一个大文件。从主机 A 到主机 B 的路线上有 3 段链路，其速率分别为

$R_1 = 500 \text{ kbps}$, $R_2 = 2 \text{ Mbps}$, $R_3 = 1 \text{ Mbps}$.

a. 假定该网络中没有其他流量，该文件传送的吞吐量是多少？

b. 假定该文件为 4MB，用吞吐量除以文件长度，将该文件传输到主机 B 大致需要多长时间？

c. 重复 (a) 和 (b)，只是这时 R_3 减少到 100kbps。

由此实例可知，吞吐量由最小的R1决定，故文件传输时只能以最小的带宽来传输，所以是4MB/500kbps。

如果R2带宽减小到100kbps比R1带宽还小，那就只能按R2的带宽来传了

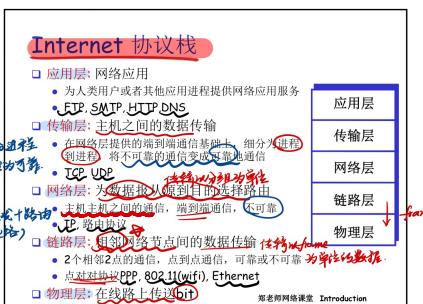
$$1 \text{ kbps} = 10^3 \text{ bps}$$

$$1 \text{ Mbps} = 10^6 \text{ bps}$$

分层的好处

Why layering? 分层的好处

- Divide complex systems to simple components
- Easy for maintenance
- Flexible for updating



△ Application Layer

client - server architectures.

Lecture 2 P1

弱一致性 poor scalability.

P2P architecture

self-scalability. 难以管理. Lecture 2 Pro.

P2P architecture

- No always-on server is needed
- End systems directly exchange data
- Peers request service from other peers, provide service in return to other peers
- Self scalability** – new peers bring new service capacity, as well as new service demands
- Peers are intermittently connected
- Dynamic IP addresses
- Q: Did you use Thunder Downloader(迅雷)? Why does it download so fast?



进程通信

进程通信

进程：在主机上运行的应用程序

□ 在同一个主机内，使用 进程间通信机制 通信（操作系统定义）

□ 不同主机，通过交换 Message 来通信

□ 使用 OS 提供的通信服务

□ 按照应用协议交换报文
□ 借助传输层提供的服务

clients, servers

客户端进程：发起通信的进程

服务器进程：等待连接的进程

Application Layer 3-11

问题1：对进程进行编址 (addressing)

□ 进程为了接收报文，必须有一个标识

即：SAP (发送也需要标示)

○ 主机：唯一的 32位IP地址 (区别于主机)

● 仅有有IP地址不能够唯一标示一个进程；在一台端系统上有很多应用进程在运行

○ 所采用的传输层协议：TCP or UDP

○ 端口号 (Port Numbers)

□ 一些知名端口号的例子：

● HTTP: TCP 80 Mail: TCP25 ftp:TCP 2

□ 一个进程：用IP+port标示端节点 (end point)

□ 本质上，一对主机进程之间的通信由2个端节点构成

Application Layer 3-13

一个进程可以有多少 network identifiers.

Addressing processes 对进程进行编址。

① 每台设备有唯一的 32-bit IPv4 和 / 或 128-bit IPv6

② Process network identifier:

• IPv4 port 192.168.1.100:80

• IPv6 port [2469:3a1:4c01:69d0:f40c:4269:74a2:7ea3]:80

• Can a process have multiple network identifiers?

✓ 3P + port number = end point

一对主和进程之间的连接由2个端点构成

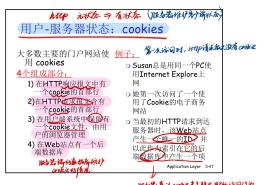


IP + port 构成一个 process.

Cookies

Lecture 2 Pg3

使得 HTTP: stateless → stateful.



产生问题

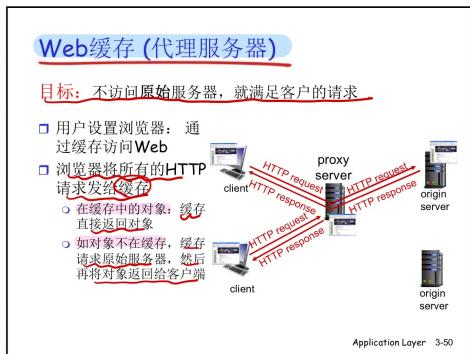
Web Caches

Lecture 2 Pg7.

Web 缓存的好处 (作用)

Why Web caching?

- reduce response time for client request
 - reduce traffic on an institution's access link
 - Internet dense with caches: enables "poor" content providers to effectively deliver content (so too does P2P file sharing)
- 减少机构内部网络
与 Internet 相连链路上的流量
多款网关都采用了缓存
加速数据的传输并能有效提供
内容 (65)



减少了通过链路的节点，Web 访问更快。(速度更快)
(increased access link speed)

应用层协议: HTTP, SMTP, FTP

DNS

C/S model

Lecture 3 P5

实现主机名 - IP 地址的转换。

(a distributed, hierarchical database)

卷

不集中管理:

集中 DNS

Why not centralize DNS?

- Single point of failure
- Traffic volume
- Distant centralized database
- Maintenance

VDP

Recursive query

P9. Authoritative DNS servers: 为被询问的 DNS 提供

TLD name server

local name server → root name server → TLD
→ authoritative DNS server → local

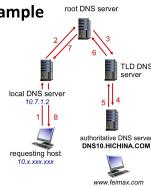
P11

Iterated query.

DNS name resolution example

• Recursive query: 递归查询

- Puts burden of name resolution on contacted name server
- Heavy load at upper levels of hierarchy

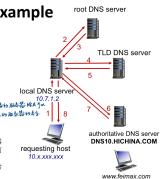


DNS name resolution example

• Host at XJTLU wants IP address for www.felixmax.com

• Iterated query: 迭代查询

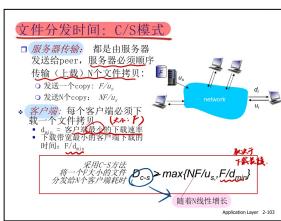
- contacted server replies with name of server to contact
- "I don't know the name, but ask this server"



P2P

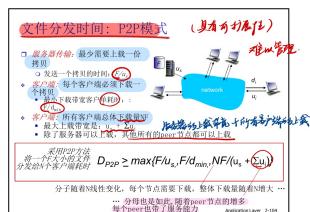
P29.

△ 文件分发



老师

文件量大，加速度上载需要



文件量大，

加速度上载需要

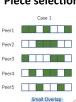
BitTorrent

Lecture 3 P36

Lecture 3 P60

△ Socket programming

Piece selection - Macro view



⇒ selection policy

Random First Piece

- Initially, a peer has nothing to trade
- Want to get a complete piece ASAP
- Select a random piece of the file and download it

First Fit

- Determine the peers that are **first fit** among your peers, and download these first.
- This ensures that the most commonly available pieces are left till the end to download.

Endgame Mode

• Near the end, missing pieces are requested from every peer containing them.

• This results in a download bottleneck due to a single peer with a slow transfer rate.

- Some bandwidth is wasted, but in practice, this is not too much.

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

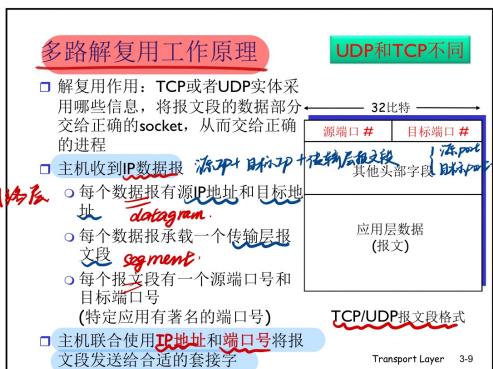
↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓ ↓

↓ ↓ ↓ ↓

多路复用 / 解复用

Lecture 4 P11



具备相同**目标IP地址**和**目标端口号**，即使不是原IP地址且源端口号的IP数据报，将会被传到相同的目標UDP套接字上

UDP

TCP

因为每条连接有各自的
不同的进程
不同的socket

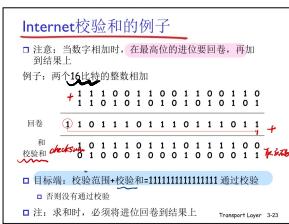
源IP地址
源端口号
目的IP地址
目的端口号

网络层

UDP checksum

Lecture 4 P25

→ detect "errors" in transmitted segment (e.g. bit errors)
addition of pseudo header, UDP header, UDP data



UDP checksum

- Goal: detect "errors" in transmitted segment

Sender:

- Treat segment contents, including header fields, as sequence of 16-bit integers.
- Checksum addition (one's complement sum) of pseudo header, UDP header and UDP data
- Sender puts checksum value into UDP checksum field.

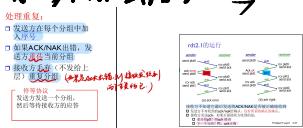
Receiver:

- Compute checksum of received segment \oplus 1111111111111111
- Check if computed checksum equals checksum field value:
 - NO - error detected 小概率发生差错
 - YES - no error detected. But maybe errors nonetheless?

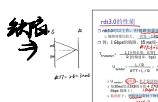
Reliable data transfer (RDT)

Lecture 4 P26

- rdt 2.0 → 若ACK/NAK出错？ \Rightarrow 每个分组中加入序号。
- rdt 2.1 →



- rdt 2.2 → 若NAK.
- rdt 3 → 超时重传机制 (可能会丢分组)



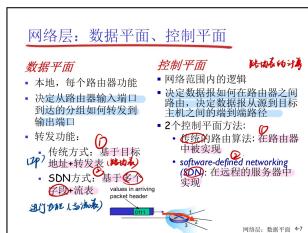
NAK free

Lecture 5, 6 还没有

超时重传机制

一次只能发一个 \Rightarrow Pipelined protocols.

△ Network layer



Lecture 6 p18

Two key network-layer functions

• 转发和路由

Network-layer functions:

- **Forwarding:** move packets from router's input to appropriate router output
- **Routing:** determine route taken by packets from source to destination

Analogy: taking a trip

- forwarding: process of getting through single interchange
- routing: process of planning trip from source to destination

路由器结构

~~最长前缀匹配~~

Longest prefix matching

When looking for **forwarding table entry** for given destination address, use **longest address prefix** that matches destination address.

Destination Address Range	Link Interface
11010100 00001111 0001***.....	0
11010100 00001111 00011000	1
11010100 00001111 000111**	2
otherwise	3

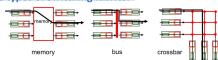
Example:
DA: 100000 00001111 .0001100000 which interface?
DA: 10010000 00010111 .0001100000 which interface?

• switching fabrics 支持结构

P29

Switching fabrics 支持结构

- Transfer packets from input buffer to appropriate output buffer
- Switching rate: rate at which packets can be transferred from inputs to outputs
 - Often measured as multiple of input/output line rate
 - N inputs: switching rate $\leq N$ times line rate desirable
- Three types of switching fabrics:



• output ports

P30

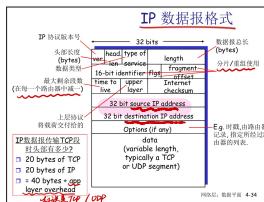
→ 优先级调度

P31

{ FIFO
 priority
 RR
 Weighted Fair Queuing

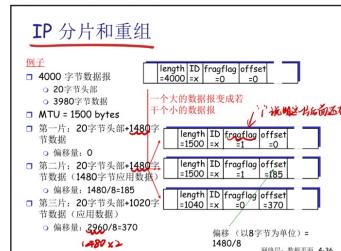
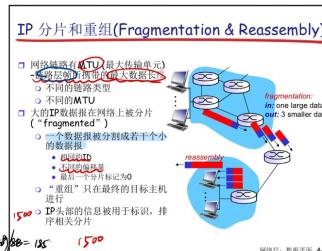
数据报文

• IP datagram.



~~IP 分片和重组 (Fragmentation & Reassembly)~~
~~MTU 讲义~~

MTU 讲义

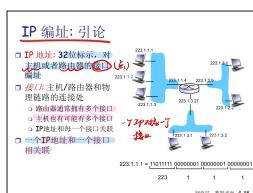


目标主机署名没有花时间内收到所有分片

⇒ 将收到的所有分片也丢弃

• IP addressing

Lecture 7 Pg



IP addressing: introduction

Q: how are interfaces directly connected?

A: by switch, we'll learn it in the future lecture



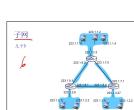
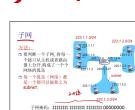
Does switch contain IP addresses? Why?

地址块与网段可以是任意长短

⇒ a.b.c.d /x
↑
subnet mask.

Subnets

判断子网掩码



CIDR.

IP addressing: CIDR

CIDR: Classless Inter Domain Routing

- Subnet portion of address can have arbitrary length
- Address format: a.b.c.d/x, where x is # bits in subnet portion of address



Network name with address format: a.b.c.d/x
Subnet mask: /x
Network prefix: x

192.168.2.1

• 主机怎么获得 IP 地址. \Rightarrow DHCP Lecture 7 P12

DHCP: Dynamic Host Configuration Protocol

Goal: allow host to dynamically obtain its IP address from network server when it joins network

- Can renew its lease on address in use
- Allows reuse of addresses (only hold address while connected/“on”)
- Support for mobile users who want to join network (more shortly)

DHCP overview:

- host broadcasts “DHCP discover” msg [optional] $\xrightarrow{\text{discover}}$
- DHCP server responds with “DHCP offer” msg [optional] $\xrightarrow{\text{offer}}$
- host requests IP address: “DHCP request” msg $\xrightarrow{\text{request}}$
- DHCP server sends address: “DHCP ack” msg

12

(UDP)

DHCP: more than IP addresses

DHCP can return more than just allocated IP address on subnet:

- Address of first-hop router for client $\xrightarrow{\text{first-hop router}}$
- Name and IP address of DNS sever
- Network mask (indicating network versus host portion of address)

• Subnet 什么获得 IP. ?

IP addresses: how to get one?

一个子网如何分配 IP 地址?

Q: how does network get subnet part of IP addr?

A: gets allocated portion of its provider ISP's address space

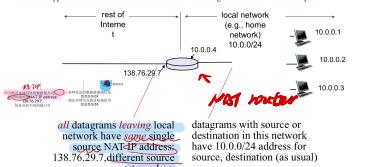
ISP's block	11001000_00010111_00010000_00000000	200.23.16.0/20
Organization 0	11001000_00010111_00010000_00000000	200.23.16.0/23
Organization 1	11001000_00010111_00010010_00000000	200.23.18.0/23
Organization 2	11001000_00010111_00010100_00000000	200.23.20.0/23
Organization 7	11001000_00010111_00011110_00000000	200.23.30.0/23

• NAT.

NAT: network address translation

Why NAT? If the subnet grows bigger, what if the ISP had already allocated the contiguous portions of address range?

And what typical network admin wants to know is how to manage IP addresses in the first place?



NAT: network address translation

Motivation: local network uses just one IP address as far as outside world is concerned:

- Range of addresses not needed from ISP: just one IP address for all devices
- Can change addresses of devices in local network without notifying outside world
- Can change ISP without changing addresses of devices in local network
- Devices inside local net not explicitly addressable, visible by outside world (a security plus)

NAT: Network Address Translation

实现: NAT 路由器必须:

- 从外数据包: 替换源地址和端口号为 NAT IP 地址和端口号, 目标 IP 地址不变
- 远端的 C/S 将会用 NAT IP 地址, 新端口作为目标地址

- 记住每个转换替换对 (在 NAT 转换表中)
 - 源 IP, 端口 vs NAT IP, 新端口

- 进入数据包: 替换目标 IP 地址和端口号, 采用存储在 NAT 表中的 mapping 表项, 用 (源 IP, 端口)

FIGURE: 图片 4-43

NAT: Network Address Translation



FIGURE: 图片 4-44

• IPv6

(IPv4的32bit) 不可用

- 头部格式改变帮助QoS

32bit \rightarrow 16bit

IPv6 数据报格式:

- 固定的40字节头部
- 数据报传输过程中，不允许分片

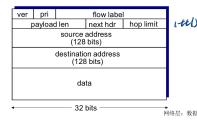
分块大小 \rightarrow 支持 \rightarrow 发送 ICMP
错误报告

网络层：数据平面 4-70

IPv6 头部 (Cont)

Priority: 标示流中数据报的优先级
Flow Label: 标识数据报在一个“Flow”
(*“flow”的概念没有被严格地定义)

Next header: 标示上层协议

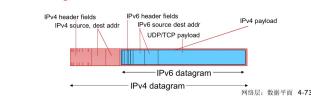


从IPv4到IPv6的平移

- 不是所有的路由器都能够同时升级的
 - 没有一个标记叫“flag days”
 - 在IPv4和IPv6路由器混合时，网络如何运转？

疑惑: 在IPv4路由器之间传输的IPv4数据报中携

带IPv6数据报

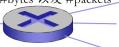


网络层：数据平面 4-73

• SDN

OpenFlow 数据平面抽象

- 流: 由分组 (帧) 头部字段所定义
- 通用转发: 简单的分组处理规则
 - 模式: 将分组头部字段和流表进行匹配
 - 行为: 对于匹配上的分组, 可以是丢弃、转发、修改、修改将匹配的分组发送给控制器
 - 优先级Priority: 几个模式匹配了, 优先采用哪个, 消除歧义
 - 计数器Counters: #bytes 以及 #packets



路由器中的流表定义了路由器的匹配+行动规则
(流表由控制器计算并下发)

网络层：数据平面 4-96

Two approaches to structuring network control plane:

- per-router control (traditional)
- logically centralized control (software defined networking) SDN

Destination-based layer 2 (switch)

Switch Port	MAC Src	MAC dst	Eth Type	VLAN ID	IP Src	IP Dst	IP Prot	TCP sport	TCP dport	Action
*	22:A7:23	*	*	*	*	*	*	*	*	port3

layer 4: 任播层

layer 3: 网络层

layer 2: 链路层

layer 2 frames from MAC address
22:A7:23:11:E1:02 should be forwarded to output port 3

port 3

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

+

• Scalable routing

lecture 8 P.14

Internet approach to scalable routing

将路由器集合为“AS”

Aggregate routers into regions known as “autonomous systems” (AS) (a.k.a. “domains”)

Intra-AS routing

- Routing among hosts, routers in same AS (“network”)
- All routers in AS must run some intra-domain protocol
- Routers in different AS can run different intra-domain routing protocol
- Gateway router: at “edge” of its own AS, has link(s) to router(s) in other ASes

Inter-AS routing

- Routing among ASes
- Gateways perform inter-domain routing (as well as intra-domain routing)

13

层次路由

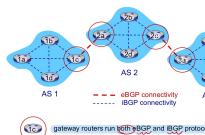
- 路由变成了：2个层次路由
 - AS内路由器：在同一个AS内路由器运行相同的路由协议
 - “intra-AS routing protocol”: 内部网关协议
 - 不同的AS可能运行着不同的内部网关协议
 - 能有效解决规模和管理问题
 - 如: RIP, OSPF, IGRP
 - 或“AS internal gateway”
 - AS间运行“AS间路由协议”
 - “inter-AS routing protocol”
 - 或“AS external gateway”
 - 或AS之间的互联私有链路

阅读材料：幻灯片第 5-7 页

Internet inter-AS routing: BGP 外部网关协议

- BGP (Border Gateway Protocol): the de facto inter-domain routing protocol**
 - “glue that holds the Internet together”
- BGP provides each AS router a means to:**
 - eBGP: obtain subnet prefix reachability information from neighboring ASes. Allows subnet to advertise its existence to rest of Internet: *I am here.*
 - IBGP: propagate reachability information to all AS-internal routers, and determine “best” routes to other networks based on reachability information and *policy*

eBGP, iBGP connections



eBGP: gateway routers run BGP and iBGP protocols
iBGP: internal routers run BGP and iBGP protocols

BGP 基础

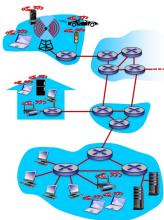
- BGP 会话：两个 BGP 路由器 (“peers”) 在一个半永久的 TCP 连接上交换 BGP 路由信息

- 通告自己目标子网前缀的“路径” (BGP 是一个“路径为驱动”的协议)
 - AS 3 到 AS 1 的路径: AS3|AS2|AS1 (从 AS3 到 AS1 的网关路由器沿 2 条路径: AS3|X 或 X|AS3)
 - 34 个 AS 的出路点，如何从 AS3 找到 X 地址
 - 语义上: AS3 和 AS2 邻居，它可以向子网 X 广播前缀
 - X 是 Z 从 Y 的下一跳 (next hop)

幻灯片第 5-7 页

Link Layer.

What's Link Layer



- Data-Link Layer has responsibility of**
 - transferring **datagram**
 - from one **node** to **physically adjacent node**
 - over a **link**
- Node:**
 - Hosts and routers
- Link:**
 - communication channel
 - connection adjacent nodes
- Layer-2 packet: frame**
 - Encapsulates datagram

~~OSI Model~~

Detecting errors.

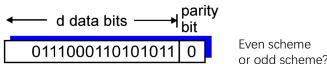
▷ parity checking

Parity checking – Single bit parity

• Even/Odd parity scheme:

- Add an additional bit
- Total number of 1 (D+1) is even/odd

• Detect single bit errors



▷ checksum

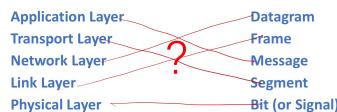
Checksum

- In the **TCP and UDP** protocols, the Internet checksum is computed over **all fields (header and data fields included)**.

- In **IP**, the checksum is computed over the **IP header** (since the UDP or TCP segment has its own checksum).

- 发送方首先根据数据算出来一个校验和，然后接收方再自行计算，看看一不一样，使用addition (one's complement sum) of segment contents

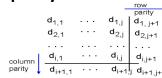
Data Unit



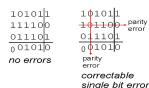
2

Parity checking – 2D parity

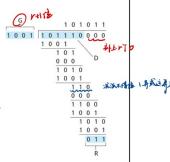
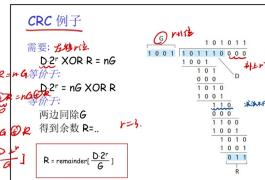
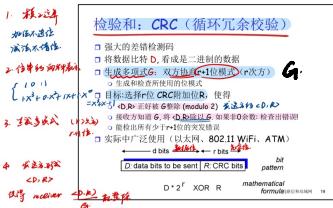
• Detect and correct single bit errors



• Two-dimensional even parity



▷ CRC



Cyclic redundancy check (CRC)

• view Data bits, D , as a binary number

- choose an $r+1$ bit pattern (generator), G
 - both sender and receiver must agree on what generator
 - generator finds CRC remainder, R , such that
 - a is bin pattern G ; b is every divisible by G (modulo-2)
 - receiver knows G , divides D by G .
 - if non-zero remainder: error detected!
 - can detect all burst errors less than $r+1$ bits



G given for both sender and receiver. The most significant (leftmost) bit of G is a 1
needs to be computed by the sender

R needs to be computed by the sender

Cyclic redundancy check (CRC)

- all calculations are done in modulo-2 arithmetic without carries in addition or borrows in subtraction, including Addition, Subtraction, Bitwise exclusive-or (XOR) (异或运算, 同得0, 异得1)
- 题目会给你用去除的多项式. 假设给你的是 $X^8 + X^7 + X^6 + X^5 + X^4 + X^3 + X^2 + X + 1$, 展开是 100000111
- 在开始算CRC之前还要在原数据之后加上多项式最高阶的0
- 算完之后那部分就是加在原数据和发送方之后
- 校验的点在于接收方在收到数据和发送方算出的原数据的校验和后, 还需自行计算出一个收到数据的校验和来与收到的进行比较, 从而达到检验的效果



• 多点访问

Lecture 9 P21 - 33

~~多点访问~~

• MAC address & ARP

MAC addresses and ARP

• IP address

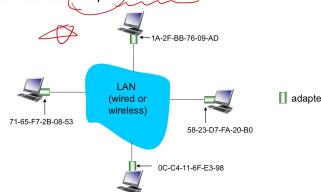
- IPv4 (32 bits) and IPv6 (128 bits)
- Network-layer address
- Layer-3 forwarding

• MAC (LAN) address

- 48 bits (e.g., 12-34-56-78-90-AB)
 - hexadecimal (base 16) notation
(each "numeral" represents 4 bits)
- Burned in NIC ROM
- Used "locally" get frame from one interface to another physically-connected interface (same network, in IP-addressing sense)

MAC addresses and ARP (MAC)

Each adapter on LAN has unique LAN address



ARP: address resolution protocol

• How to determine interface's MAC address, if knowing its IP address?



ARP table: each IP node (host, router) on LAN has the table

- An IP-to-MAC address mappings for LAN nodes
- IP address MAC address: TTL = 1
- TTL (Time To Live): time after which address mapping will be forgotten (typically 20 min)

Note that an ARP table is not necessary to contain an entry for every host and router on the subnet; some may have never been entered into the table, and others may have expired.

ARP协议：在同一个LAN（网络）

- A发送请求包-SIB的IP地址
- B接收到后，如果目的IP地址不是自己的IP地址，则将IP-MAC映射关系存入自己的ARP缓存中
- B将包含自己的IP地址的ARP应答包
- Dest MAC: 56-84-7B-0F-FF-FF
 - LAN: 有线或无线会自动判别机架
- LAN: 将源IP和目的IP映射到自己的MAC地址
- B将接收到的ARP包，回复A自己的IP地址
- A接收到后，更新自己的ARP表缓存
- 无需网络管理员干预
- 比如说A
 - 例如A: 192.168.1.100
 - 用A的MAC地址（单播）

ARP实例

只有跨过子网的传输才需要将帧交给路由器；帧中包含目的地IP和路由器与出发地所处子网接口的MAC地址

ARP报文
源MAC: 00-0C-11-6F-E3-98
源IP: 192.168.1.100
目标MAC: FF-FF-FF-FF-FF-FF
目标IP: 192.168.1.111
TTL: 1
校验和: 0x0000
用MAC地址(单播)

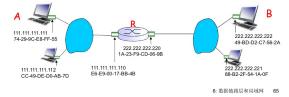
→ ARP消息的生成与发送
请求报文的生成。

△ 不同 LAN

路由到其他 LAN

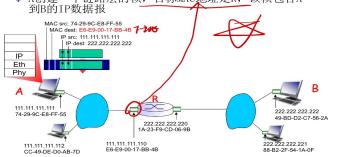
Walkthrough
假设我们知道的IP地址

- 在R上有两个MAC地址，分别对应两个LAN。IP地址是通过~~forwarding~~转发的
- 在源主机的ARP表中，发现目标主机的下一跳时111.111.111.110
- 在源主机的ARP表中，发现其MAC地址是E6-E9-00-17-BB-48。MAC
- etc



编址：路由到其他 LAN

- A创建数据报，源IP地址: A; 目标IP地址: B
- A创建一个链路层的帧，目标MAC地址是R，该帧包含A到B的IP数据报



• Ethernet Frame (protocols)

Ethernet frame structure

- Sending adapter **encapsulates** IP datagram (or other network layer protocol packet) in **Ethernet frame**



preamble: 8-byte

- 7 bytes with pattern 10101010 followed by one byte with pattern 10101011
- used to synchronize receiver, sender clock rates

Ethernet frame structure

- **Address:** 6 byte source, destination **MAC addresses**.

- If adapter receives frame with matching destination address, or with broadcast address (e.g. ARP packet), it passes data in frame to network layer protocol

- Otherwise, adapter discards frame

- **Type:** indicates higher layer protocol (mostly IP but others possible, e.g., Novell IPX, AppleTalk)

- **CRC:** cyclic redundancy check at receiver

- Error detected: frame is dropped



Ethernet: unreliable, connectionless

- **Connectionless:** no handshaking between sending and receiving NICs

- **Unreliable:** receiving NIC doesn't send acks or nacks to sending NIC

- Data is lost if the receiving NIC does not receive frame before sender uses higher layer rdt (e.g., TCP), otherwise dropped data lost

- Ethernet's MAC protocol: unslotted CSMA/CD with binary backoff (using with sensing channel idle)

APP表: IP层上 ~~同一IP~~ 两个APP表.

路由表: 路由器上

> forwarding table

Switch table: 支线加上.

• Switch.

多路同时传输

Switch: multiple simultaneous transmissions

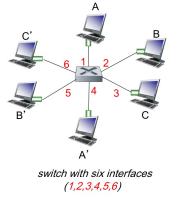
- Hosts have dedicated, direct link to switch
- Switches buffer packets

- Ethernet protocol used on each incoming link, but no collisions; full duplex

- Each link is its own collision domain

- Switching, e.g.; A-to-A' and B-to-B' can transmit simultaneously, without collisions





Self-learning

- Switch learns which hosts can be reached through which interface

- When frame received, switch "learns" location of sender

- Records MAC/interface (sender/location) pair in switch table

- Entries get removed if no frames with that MAC received after timeout (which can be configured)

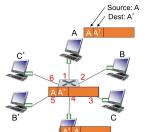
MAC addr	interface	Time
A	1	9:10
A'	4	9:15

switch table (initially empty)

Self-learning

- frame destination, A', location unknown: flood

- destination A location known: selectively send on just one link



MAC addr	interface	Time
A	1	9:10
A'	4	9:15

switch table (initially empty)

Switches vs. routers

Both are store-and-forward:

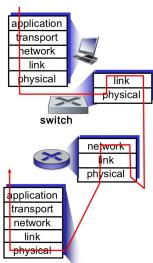
- **Routers:** network-layer devices (examine network-layer headers)

- **Switches:** link-layer devices (examine link-layer headers)

Both have forwarding tables:

- **Routers:** compute tables using routing algorithms, IP addresses

- **Switches:** learn forwarding table using flooding, learning, MAC addresses



10

Lecture 10 Pg5

请求 goole 网址的全过程.

Network Security

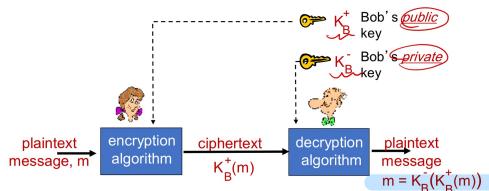
△ Symmetric key cryptography

DES, AES

symmetric key crypto

- requires sender, receiver know shared secret key
- Q: how to agree on key in first place (particularly if never "met")?

△ public key cryptography



RSA

有3種不能解密的
key

Lecture 10 P51

$$K_B^-(K_B^+(m)) = m = K_B^+(K_B^-(m))$$

RSA 特點: secure

缺點: } computationally intensive
DES is faster than RSA

應用: use public key crypto to establish secure connection
then establish second key - symmetric session key - for encrypting data
faster.

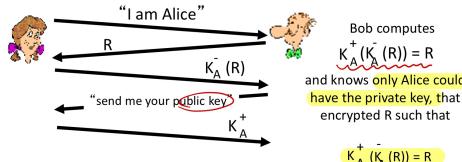
session key, K_s

- Bob and Alice use RSA to exchange a symmetric key K_s
- once both have K_s , they use symmetric key cryptography

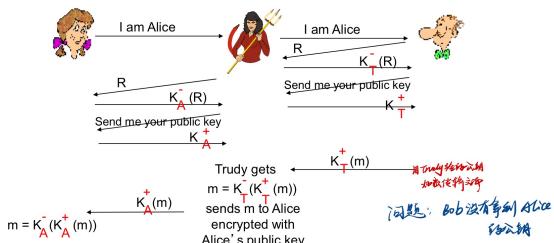
Authentication

- ap4.0 requires shared symmetric key
- can we authenticate using public key techniques?

ap5.0: use nonce, public key cryptography



man (or woman) in the middle attack: Trudy poses as Alice (to Bob) and as Bob (to Alice)



13

密文 Bob 与公钥 k_A : $K_A^-(Alice, K_A^+)$

Message integrity

Digital signatures

- suppose Alice receives msg m , with signature: $m, K_B^-(m)$
- Alice verifies m signed by Bob by applying Bob's public key K_B^+ to $K_B^-(m)$ then checks $K_B^+(K_B^-(m)) = m$.
- If $K_B^+(K_B^-(m)) = m$, whoever signed m must have used Bob's private key.

Alice thus verifies that:

- Bob signed m
- no one else signed m
- Bob signed m and not m'

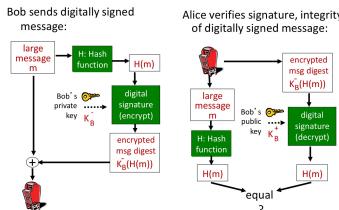
non-repudiation:

Y. Alice can take m and signature $K_B^-(m)$ to court and prove that Bob signed m

computationally expensive to public key- encrypt long messages

$\Rightarrow H(m)$: Hash function.

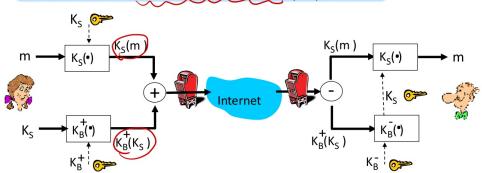
Digital signature = signed message digest



~~Confidential~~

Secure e-mail

Alice wants to send confidential e-mail, m , to Bob.



Alice:

- generates random *symmetric* private key, K_S
- encrypts message with K_S (for efficiency) $K_S(m)$
- also encrypts K_S with Bob's public key $K_B^+(K_S)$
- sends both $K_S(m)$ and $K_B^+(K_S)$ to Bob

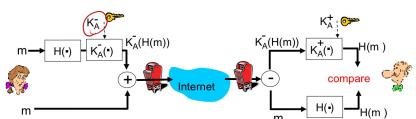
28

认证 + 完整性.

Secure e-mail (continued)

实境认证.

Alice wants to provide sender authentication message integrity



- Alice digitally signs message
- sends both message (in the clear) and digital signature

发送机密 (博主) + 电子签名

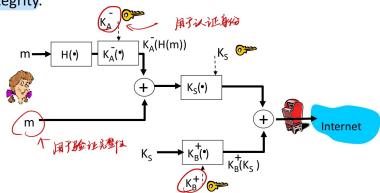
30

认证 + 完整性 + 加密

Secure e-mail (continued)

发送机密信息 + 实境认证.

Alice wants to provide secrecy, sender authentication, message integrity.



Alice uses three keys: her private key, Bob's public key, newly created symmetric key

SSL : TCP 上的 confidentiality / integrity / authentication

IPsec : IP 上的