

Internet Protocols and Networks.

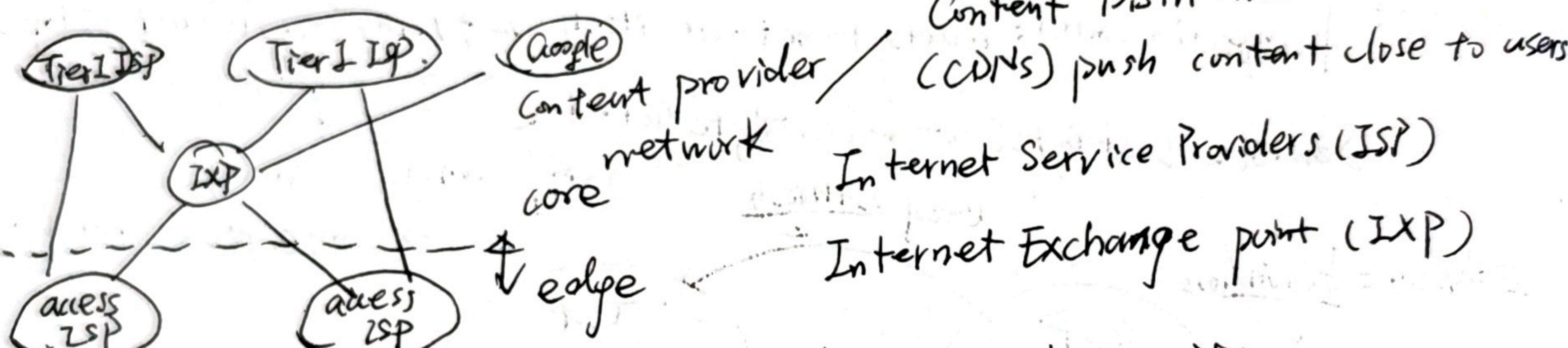
<Block 1>

A. Introduction to IP.

1. The network core (3/1)

store-and-forward $\xrightarrow{\frac{L \text{ bit}}{R \text{ bits}}} \text{end-to-end delay}$ ($1 \text{ m bits} = 10^6 \text{ bits}$)

2. network of network



3. loss \leftarrow packet arrival rate to link exceeds output link capacity

$$d_{node} = d_{processing} + d_{queue} + d_{transmission} + d_{propagation}$$

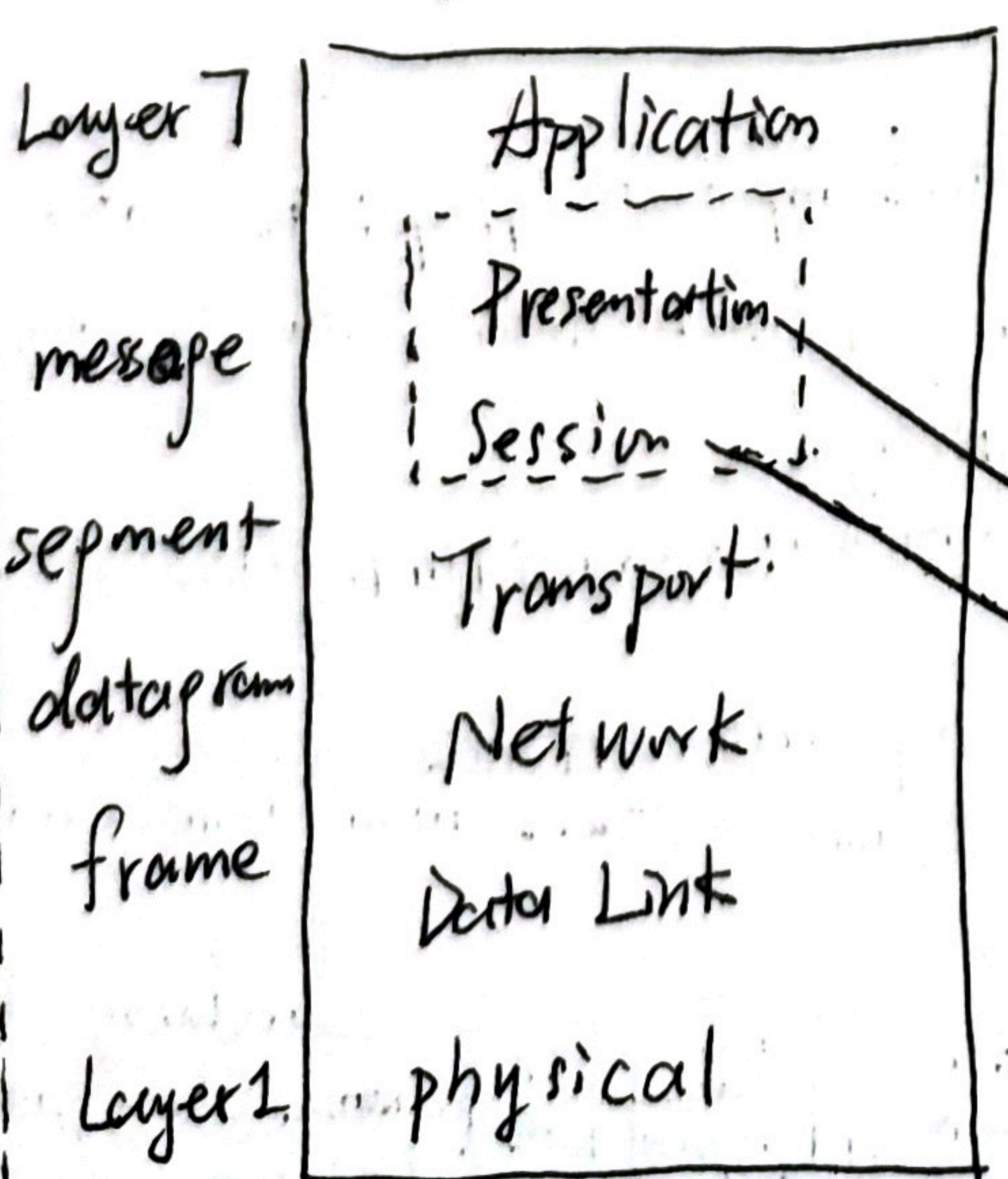
\hookrightarrow determine output link

Delay } limit \rightarrow Throughput } instantaneous rate (bps) at which bits transferred
loss } \rightarrow average between sender/receiver

bottleneck link: link on end-end path that constrains end-end throughput.

网络瓶颈 - 一般高達率連接， bottleneck link \rightarrow 在於 \rightarrow 哪個連接。

4. 協議
Why \leftarrow explicit structure allows identification, relationship of complex system pieces.
modularization eases maintenance, updating of system.



TCP/IP (5层): Transmission Control Protocol / Internet Protocol

ISO/OSI (7层): International Standards Organization / Open Systems Interconnection.

Traslates data for an application (encode)

Take care of a connection between two hosts for the life-time of that connection (Authentication + authorisation)

shift + 按鍵

B. Transport layer (I).

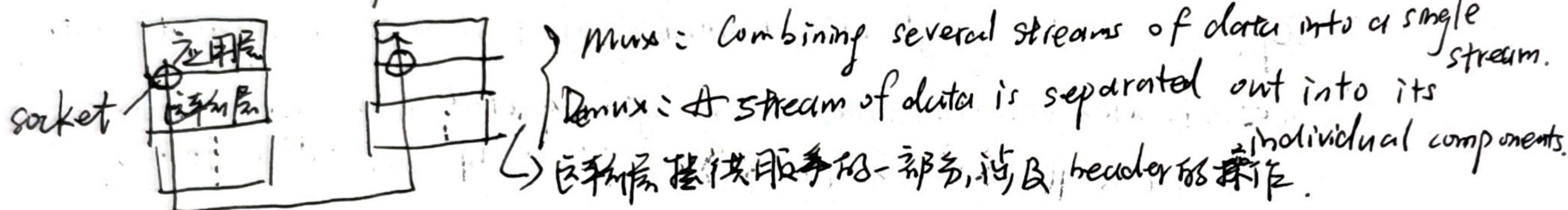
1. Transport-layer services.

Transport layer → logical communication between processes
Network layer → logical communication between hosts.

send. break app messages into segments pass to network layer
receiver ~~ress~~ reassemble segments into messages, pass to app layer

? reliable, in-order TCP, UDP;

2. Multiplexing / demultiplexing



UDP socket = \langle destination IP, destination port # \rangle

TCP socket = \langle destination IP, destination port, other source IP, source port \rangle

3. Connectionless transport: UDP (User Datagram Protocol)

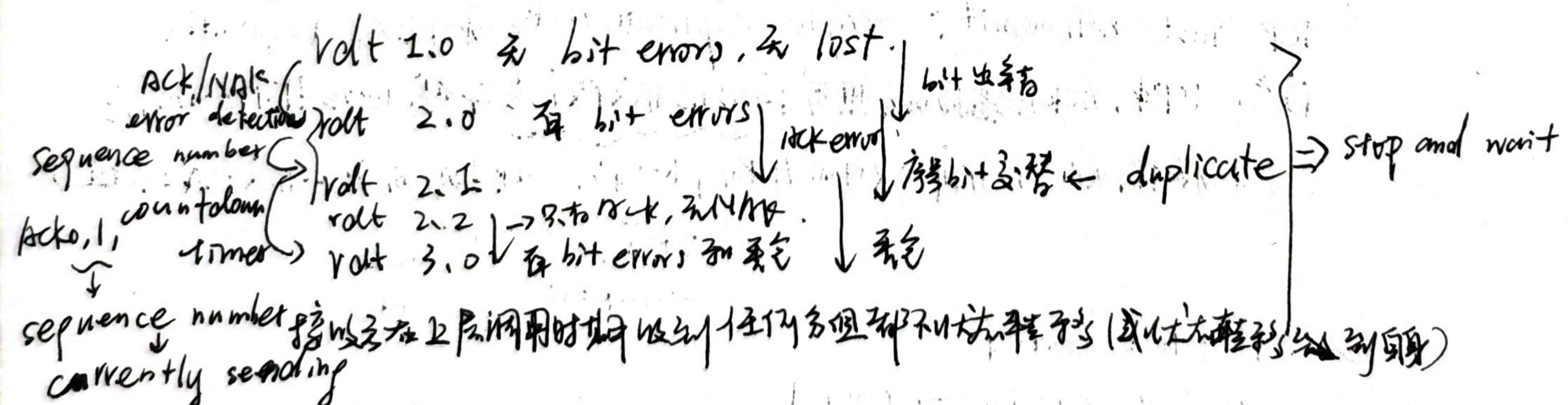
<u>source port</u>	<u>dest port</u>
<u>length</u>	<u>checksum</u>
<u>payload</u>	

64 bit = 8 Byte

detect "errors" in transmitted segment.

4. Principles of reliable data transfer

reliable data transfer protocol → (rdt) 模型



A. Transport layer (II)

$$\lambda \rightarrow \text{avg. utilization} \quad \frac{\text{RTT} + 4R}{\text{RTT}}$$

1. pipelined protocols

Gro-Baick-N
(GrBN)

Go-Back-N
(GBN) } 只能确认 cumulative ACK } 重传机制 } 自
失序报文 drop }

Selective R

女孩破坏人 cu
girl破坏人 cu

只有一千五时哭了五连哭。

向日有時事失，失時有時事失。

其系尔脉，有口能言

自從打 whalar

GRISY 有 sender 有
Grisy 有 sender 有

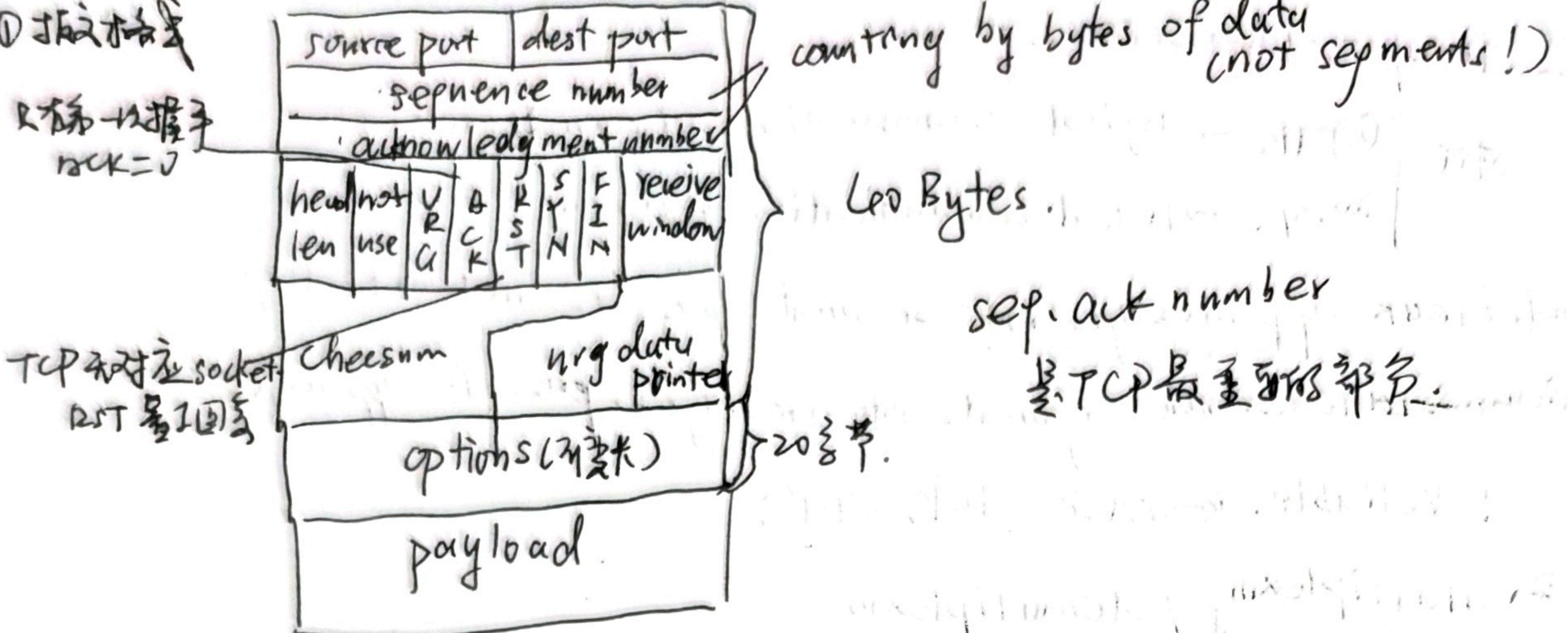
DR. Smith window
不一致，今理

立和落差之间，一列

商貿以 ≤ 133

2. TCP

① 报文格式



sequence number: byte stream "number" of first byte in segment's data

acknowledgement number: next byte expected from other side
(TCP 只确认该流至第一个丢失的数据节为止, TCP cumulative Ack).

收到头序报文段的头部时即应用而解.

② RTT (round trip time) and timeout,

$$\text{EstimatedRTT} = (1-\alpha)\text{EstimatedRTT} + \alpha \cdot \text{SampleRTT} \quad (\alpha=0.125)$$

$$\text{DevRTT} = (1-\beta) \text{DevRTT} + \beta | \text{SampleRTT} - \text{EstimatedRTT} | \quad (\beta=0.25)$$

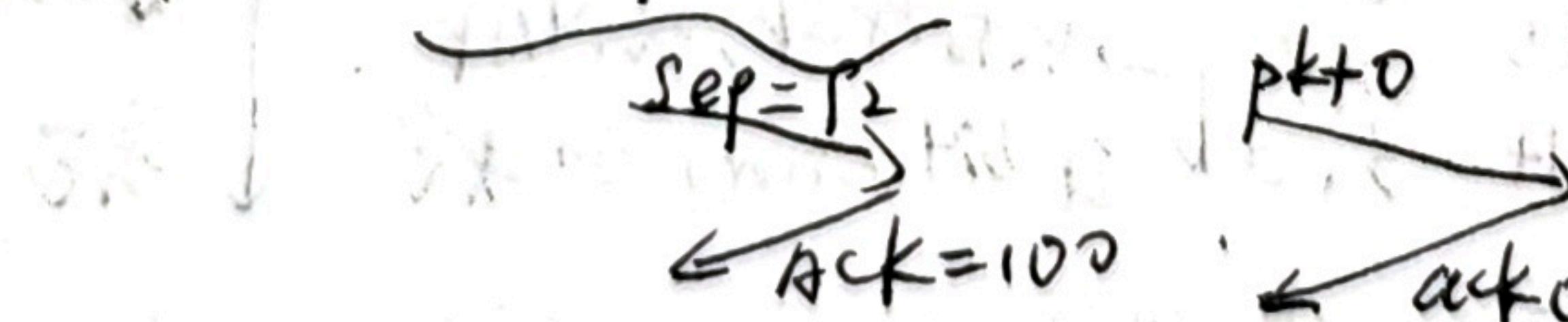
$$\text{TimeoutRTT} = \text{EstimatedRTT} + 4 \text{DevRTT}$$

$\begin{cases} \text{too short: early timeout, unnecessary retransmissions} \\ \text{too long: slow reaction to segment loss.} \end{cases}$

TCP fast retransmit: triple duplicate ACKs (第三个起算 Duplicate)

注意: TCP 中, ACK 序号接收方想要接收的数据段序号而不是回传的接收序号

(ACK 和序号大小在 PPTP 中含义不同, ACK 为下一字段, ACK 1 则响应 1 段数据或 1 个包)



③ TCP flow control

receive window \rightarrow rwnd

receiver $\begin{cases} \text{保证接收缓冲区不 overflow. } \text{LastByteRecd - LastByteRead} \leq \text{RcvBuffer.} \\ \text{同时} \text{rwnd} = \text{RcvBuffer} - [\text{LastByteRecd} - \text{LastByteRead}] \end{cases}$

sender $\begin{cases} \text{LastByteSent - LastByteAcked} \leq \text{rwnd} \end{cases}$

$\text{MSS} = \text{IP header} + \text{TCP header} + \text{MTU}$

maximum segment size, $\text{MSS} \leq \text{Window}$

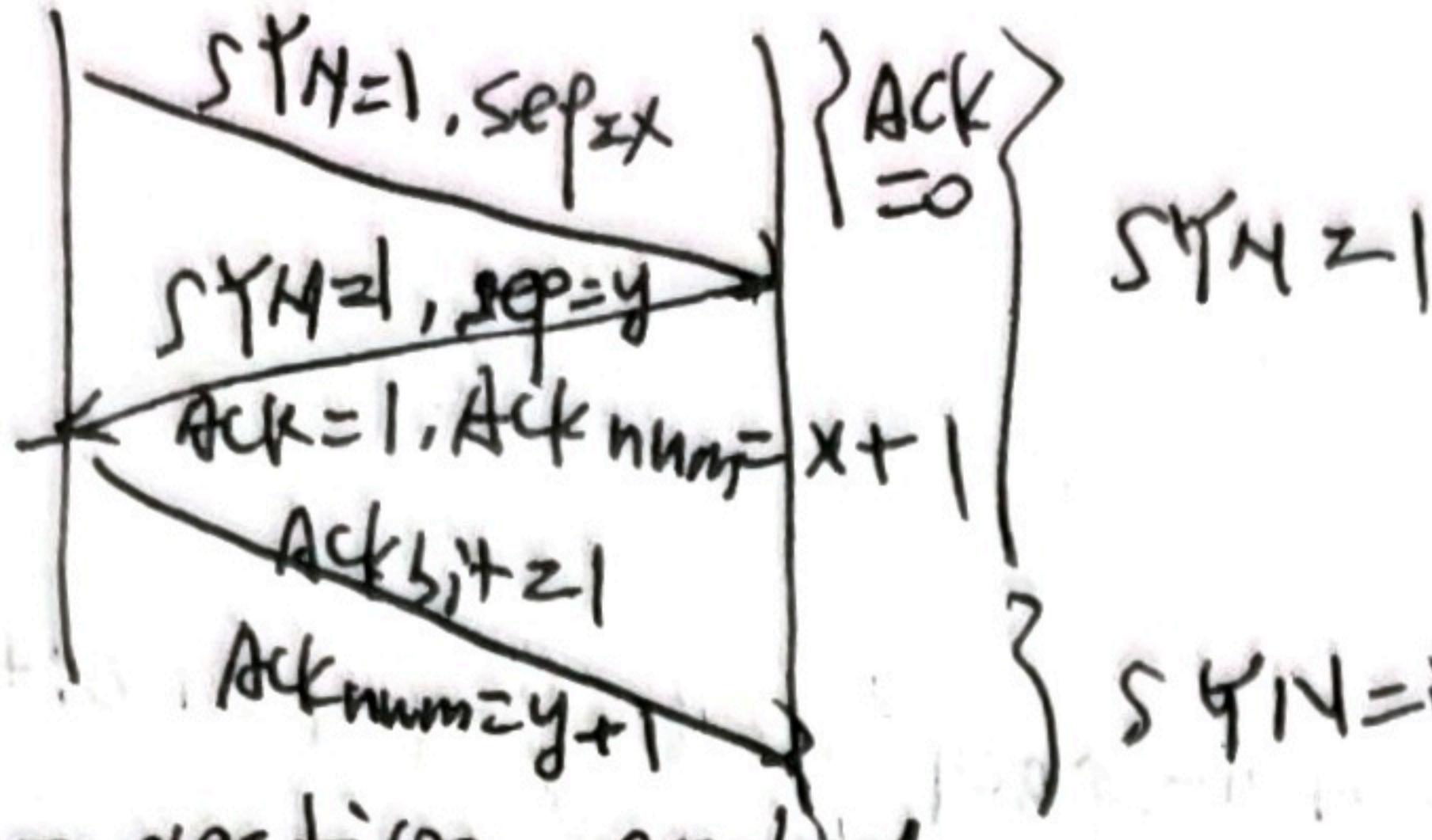
Maxle's Algorithm $\begin{cases} \text{cause: when an application generates data very slowly, Overhead of this is huge} \\ \text{发送方发送一个数据} \\ \text{当窗口缓冲区满时, 直到} \end{cases}$

向 TCP 发送 ACK

窗口缓冲区到 MSS (DataSize)

④ Connection Management

詳註



SYN=1

3. TCP congestion control.

Flow control: controlled by receiver limit window for connection. Stops receiver getting too much traffic

Congestion control: reacts to network itself being overloaded, stops network getting too much traffic

Sender: Last Byte Sent - Last Byte Acked

$ssthresh$

min { cwnd, rwin } } → congestion window.

$$ssthresh = \frac{1}{2} cwnd_{loss}$$

丟失後 cwnd

TCP
Slow Start

TCP

Congestion Avoidance

initially cwnd = 1 MSS

initial rate is slow

but ramps up exponentially fast.

3 duplicate Ack → 快速重傳, $cwnd = \frac{cwnd}{2}$, grow linearly

TCP flavours

TCP Reno: timeout → 快速重傳.

TCP Tahoe: 丢包時降低窗口大小 (設置 cwnd 為 1).

Additive Increase Multiplicative Decrease (AIMD)

loss 發生 $\frac{w}{2}$ $\frac{w}{2}$ $\frac{w}{2}$ $\frac{w}{2}$... avg TCP throughput = $\frac{3}{4} \frac{w}{RTT}$ bps.

TCP 有 Fairness.

B. Network layer. Data plane.

1. Overview.

Forwarding: move packets from router's input to appropriate router output.

Routing: determine route taken by packets from source to destination.

Network service model

guaranteed delivery

guaranteed minimum bandwidth to flow

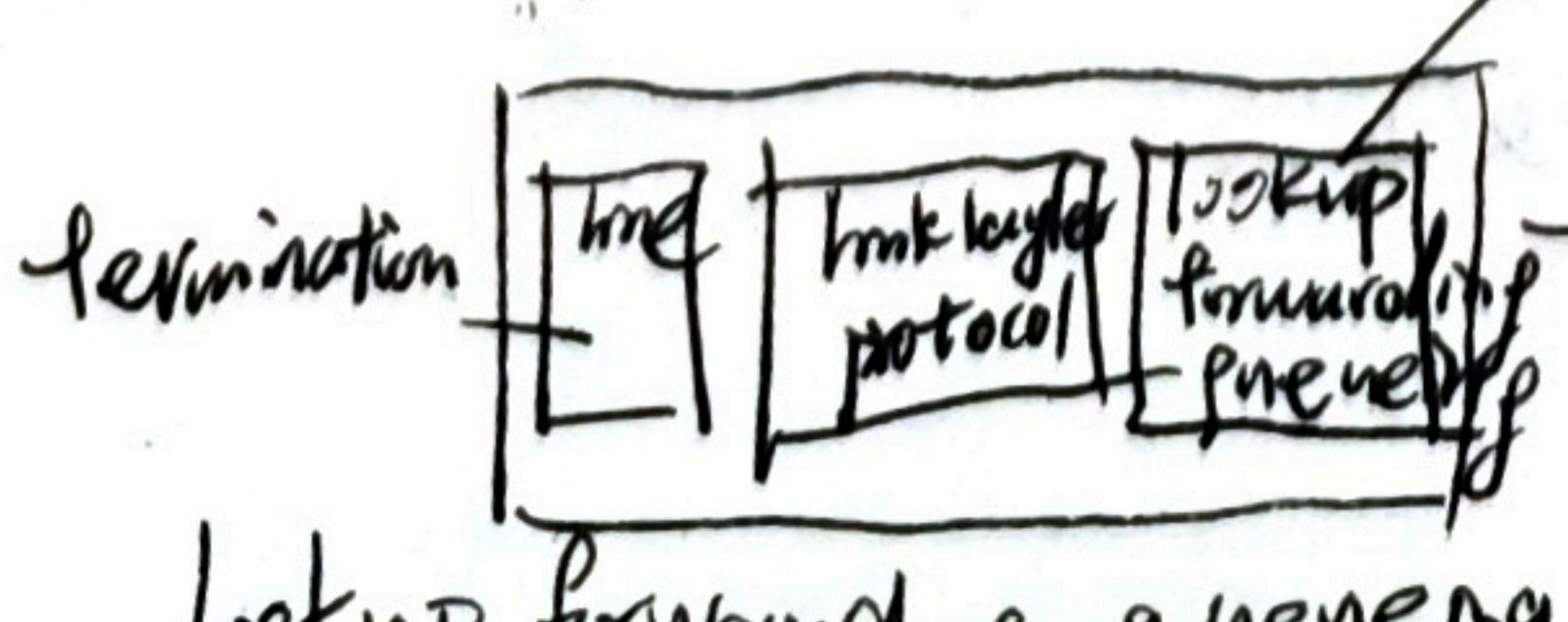
In order

guaranteed delivery with less than time delay.

2. Router.

destination-based forwarding
generalized forwarding

scheduling discipline
chooses among queued
datagrams for transmission.



lookup, forwarding, queuing

switching
fabric

buffer
queueing

link
layer
protocol

Scheduling discipline

3. 路由表 → longest prefix matching

④ IPv4

① IP datagram format.

TOS, QoS 服务质量 IP header

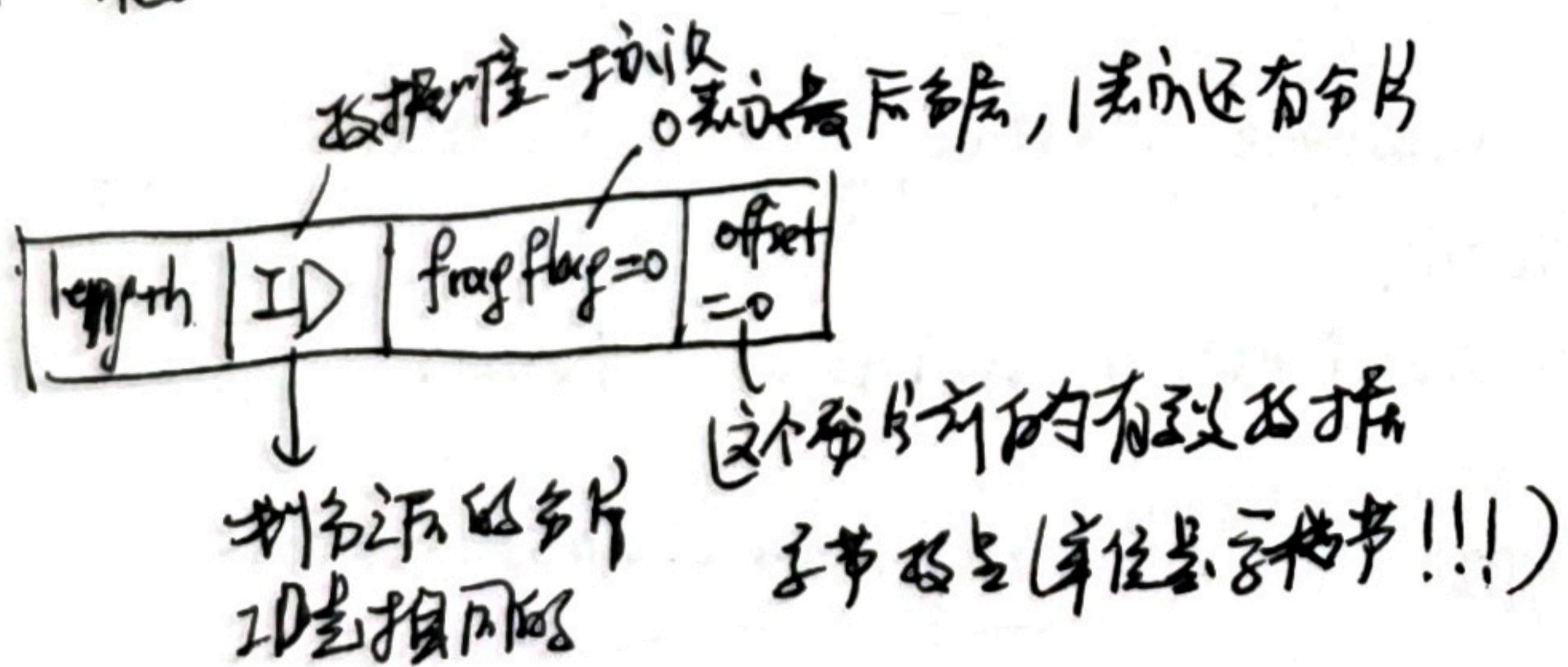
Ver	header length	Type of service	length
16-bit identifier	flag	fragment offset	
time to live	upper layer	header checksum	
			→ 32 bit, 等于四层或五层头部之和
			→ 20 byte.
			32 bit source IP
			32 bit destination IP
			IP options
			payload

逐跳头部
NLPI 头部
协议

② 数据分片、Fragmentation and reassembly

因为 MTU 存在，并且不同链路有不同的 MTU，发送方和路上的路由器都对 my datagram 进行分片，多片数据报只有 final destination 才被 reassembled.

IP header bits used to identify, order reported fragments.



TCP 首部为 20, UDP 首部为 8, 有时用来计算报文
IP 首部一般 20 Byte
→ 报文首部载荷 + 首部 = MTU
一般为 1500 bytes

③ IP

IP address is identifier for host, router interface, associated with each interface.

classful IP addressing. Class A / Class B / Class C / 128 / 112 / 24 网络部分长度, subnet mask

Subnet: can physically reach each other without intervening router

< Subnetting >

What: splitting a block → smaller blocks, each subnet having its own subnet mask.

Why: want divide into individual networks, self, 一个很大的大段 IP, 不同部分利用。
a private site router is used for subnetting

划分子网掩码。

Classless addressing - CIDR →

无类 IP 地址空间划分。

划分 IP 地址 CIDR ← slash

Network	Host
Prefix	suffix
n	32-n

prefix length.

IP → Network address, 除了子网掩码，所有主机号置 0, 在地址上 /n.
Broadcast address, 主机号全部置 1, 不可用为主机 IP.

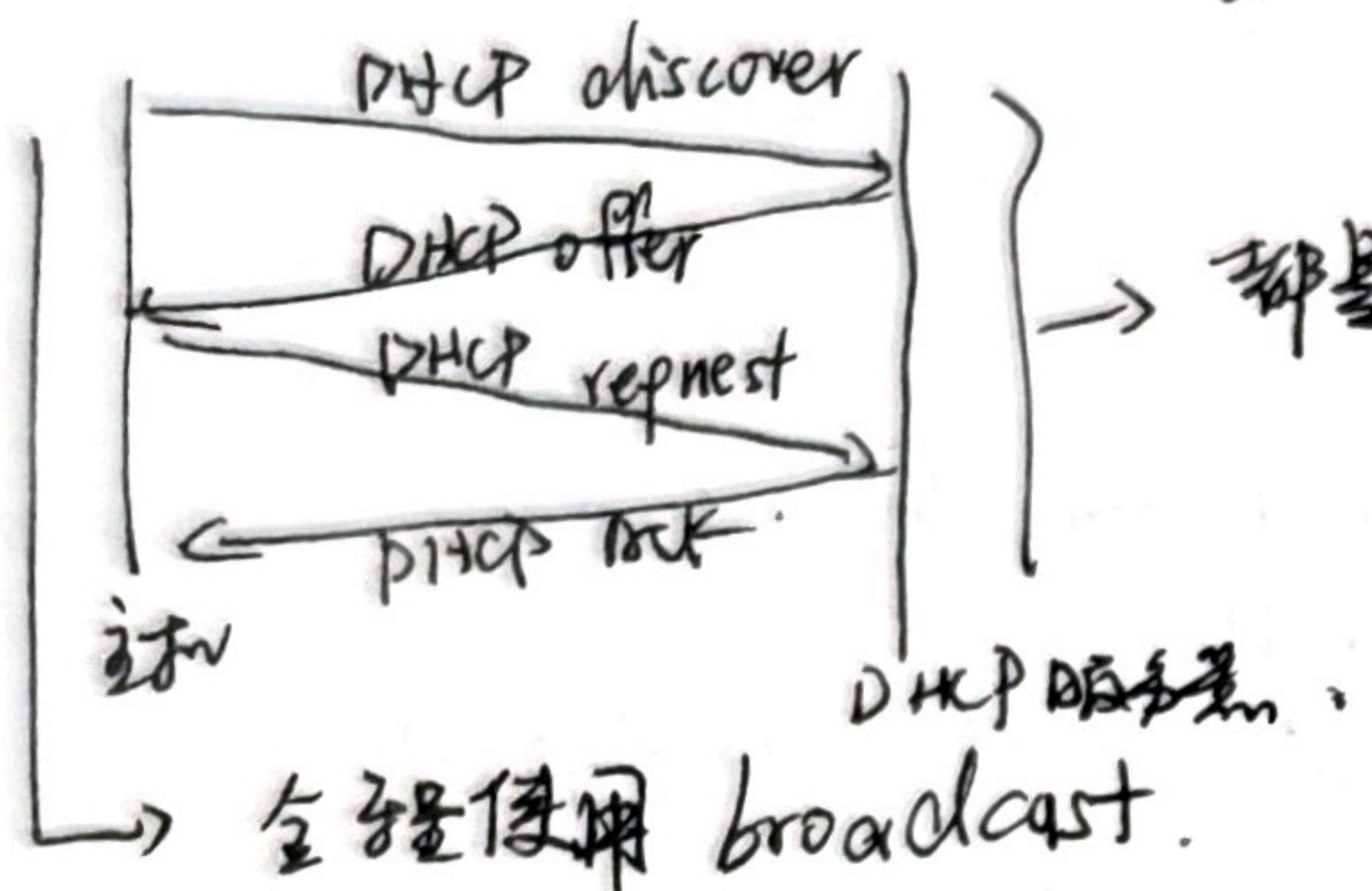
network address 基于第一个子网 IP

这两种划分方法 (适用于不同背景)

- 将一块网络划分为若干子网，使得每个网内主机数相同，且从不属于同一块放入子网后，划分出子网的掩码位相同。
- Variable length subnet mask (VLSM)，每个子网内的主机数量不同，从大到小开始，优先满足用户的需求（位数），然后从小到大设计，一定不能重叠，划分的子网的掩码长度必须不同，因为凡先满足的量主机数位。

5. DHCP = Dynamic Host Configuration Protocol.

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



DHCP不仅可以获得主机IP,

还可以 first-hop router IP.
IP, name for DNS server.
network mask.

TCP为点到点通信，所以DHCP报文
会在UDP协议之上运行

6. 使用逐个网络前缀分配多个网络的网内被称为主址聚合

Hierarchical addressing → route aggregation (地址聚合, 网络聚合).

hierarchical addressing allows efficient advertisement of routing information.

<Block 3>.

A. 7. Network Address Translation (NAT).

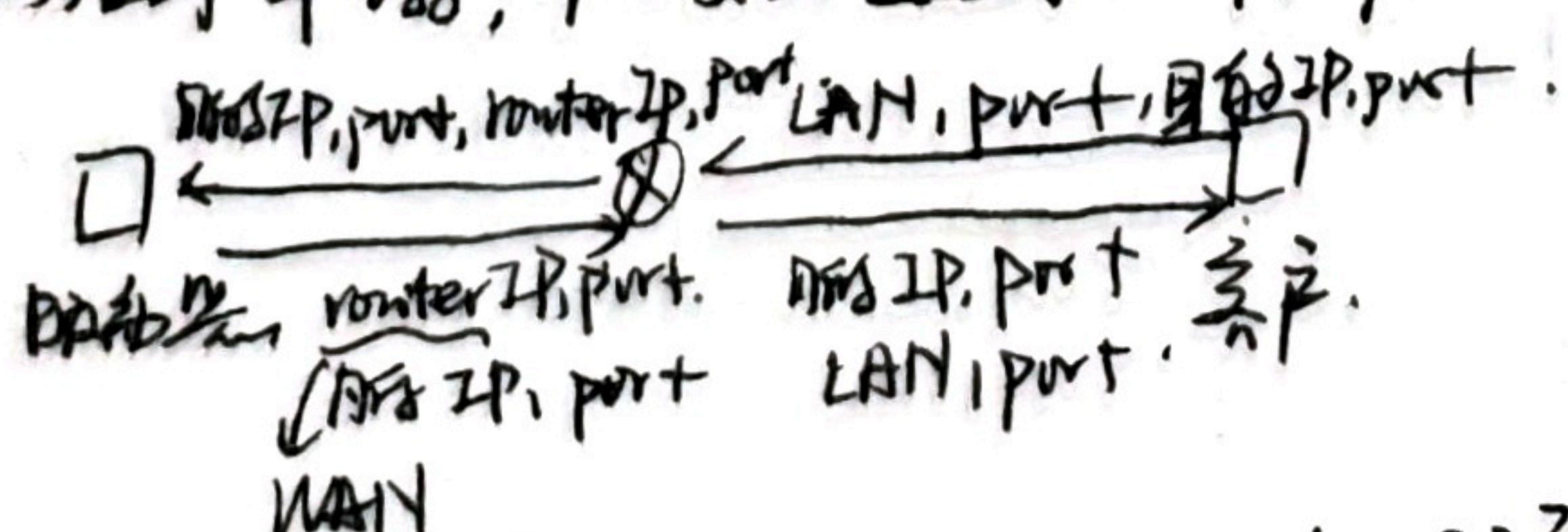
使用 NAT router 隐藏子网中 hosts，使得整个子网从外部看好像只有 1 个 IP 地址。

NAT router 中有一个

NAT translation table.	
WAN side address	LAN side address
198.76.29.7:5001	10.0.0.1:3456

IP address = $\frac{\text{port}}{\text{32 bit.} \quad \text{16 bit.}}$

不需要 NAPT 串接，可以更改主机 IP，可以更改 IP 而主机不变。



路由器对外部都是 WAN 地址，对于内部来说，可以将 NAT router 当作透明的。

路由器对外部都是 WAN 地址，对于内部来说，只知道自己在 router (NAT) 里面。
对于外部来说，它意识不到客户的存在，只知道自己在 router (NAT) 里面。
Controversy 1: router 在第 3 层用共享信道，这个问题已被 Hot Potato, violate end-to-end argument.

8. IPv6.

动机: IPv4 用完了。

additional motivation: header format helps processing / forwarding

less byte 头部，无 fragmentation, attack

① IPv6 datagram format.

version	priority	flow label
payload length	next header	hop limit
source address (128 bits)		
destination address (128 bits)		
TTL		

→ identify priority among datagrams in flow

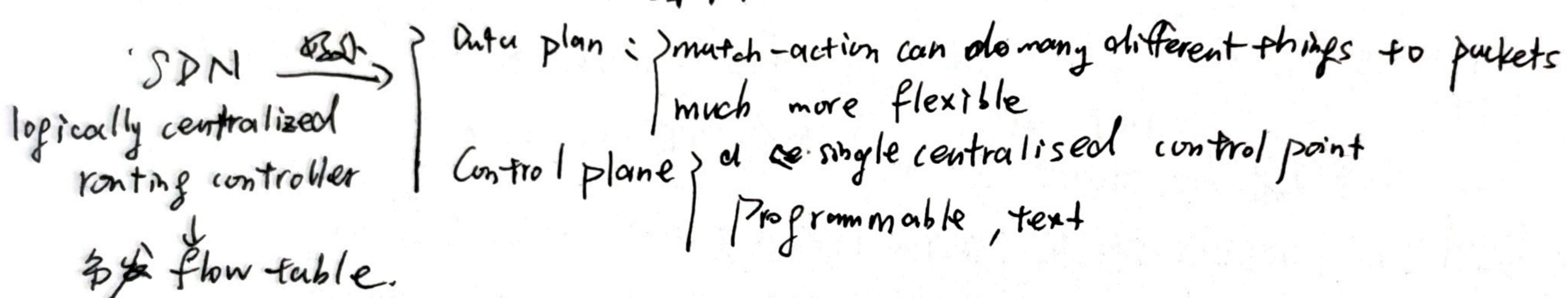
需要对经那个链路层协议。

头部占 40 Bytes.

② IPv4 + IPv6 用途

tunneling: IPv6 datagram carried as payload in IPv4 datagram among IPv4 routers

3. Generalized Forward and SDN.



OpenFlow ← a specific SDN protocol.

Rule	Action	States
Layer 2 Link, Network, Transport	packets + byte counters	
Layer 3, IP protocol	Forward, encapsulate and forward to controller	
	drop, send to normal processing pipeline	
	modify fields	

B. Network layer (Control) < 路由 >

1. { per-router control
logically centralized control

Autonomous System

{ a group of routers under a single administrative authority

{ static
dynamic.

2. Routing algorithm → { global
decentralized

Dijkstra → oscillation possible.

$D_V \rightarrow D_X(y) = \min_V \{ c(x, v) + D_V(y) \}$, count to infinity → poisoned.

3. Intra-AS OSPF.

↳ autonomous system. IGP (Interior gateway protocols) → { RIP
OSPF (LS-IS)
IGRP

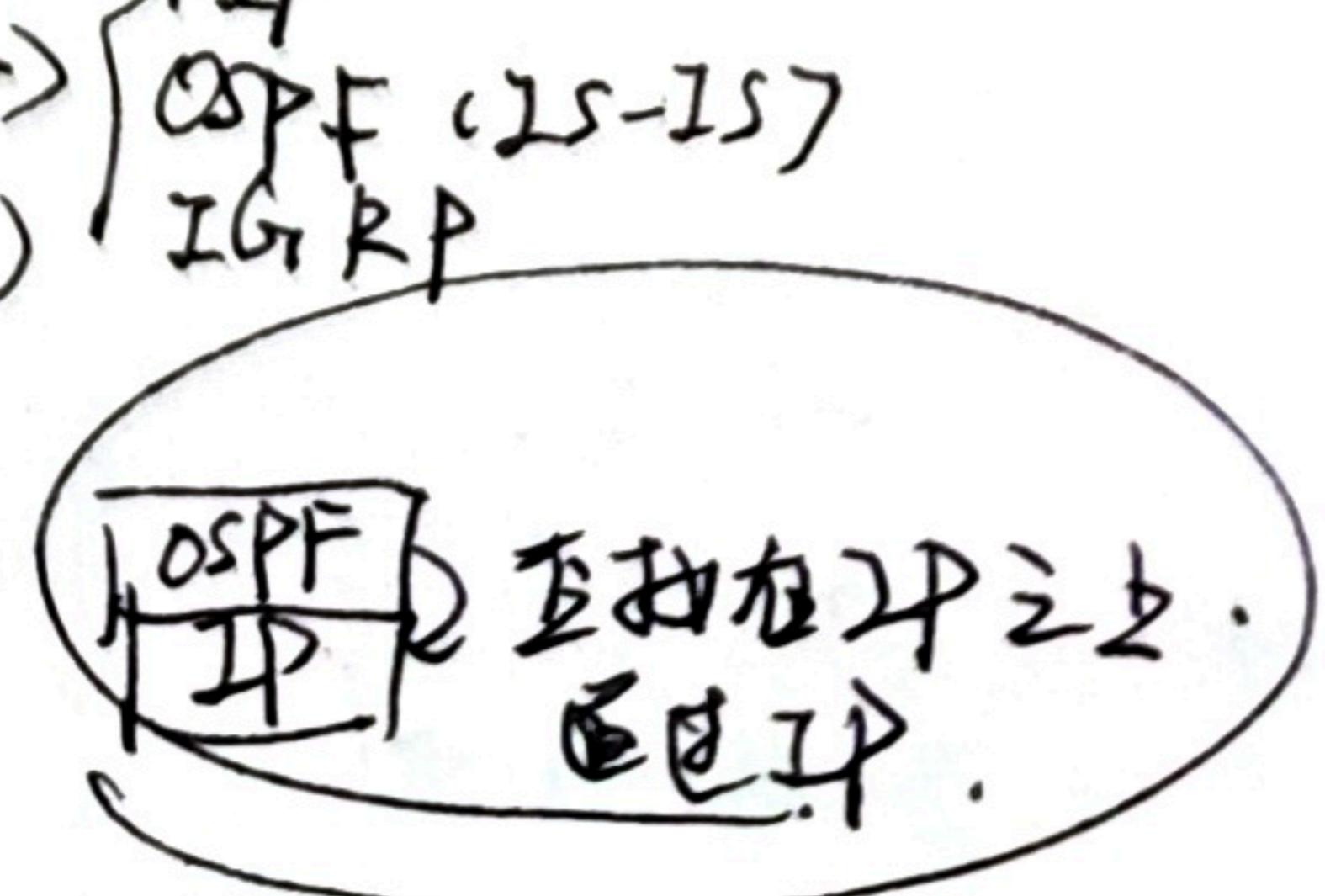
OSPF → { LS algorithm

floods OSPF LS advertisements in entire AS

↳ { security
multiple path

multiple cost-based

hierarchical OSPF.



4. BGP.

Border Gateway Protocol.

基于TCP之上

prefix + attributes = A route

↳ (AS-PATH, NEXT-HOP)

eBGP, iBGP.

BGP route selection } local preference value attribute = policy decision.
 shortest AS-PATH.
 closest NEXT-HOP router = hot potato routing
 additional ~~criterion~~ criterion.

5. SDN control plane.

Why SDN } easier network management
 table-based forwarding allows "programming" routers.
 open implementation of control plane

network-control app } control plane.
 SDN controller

data plane. ↳ OpenFlow → 基于TCP之上

6. ICMP. internet control message protocol. → Traceroute → } UPnP.
 在IP数据报之上 ↳ 可能是ICMP. } unreachable port
 communication network level information. } IP name & address
 Router, RTTs record.

内容: type, code, plus 8 Byte of IP datagram causing error

↳ 错误源→发送方不响应引发该差错的IP datagram.

Q. Network security basics.

1. firewall

isolates organization's internal net from larger Internet, allowing some packets to pass, blocking others.
 例如 → attack, illegal modification / access of internal data.

① Stateless packet filtering

Access Control Lists; table of rules.

② Stateful packet filtering

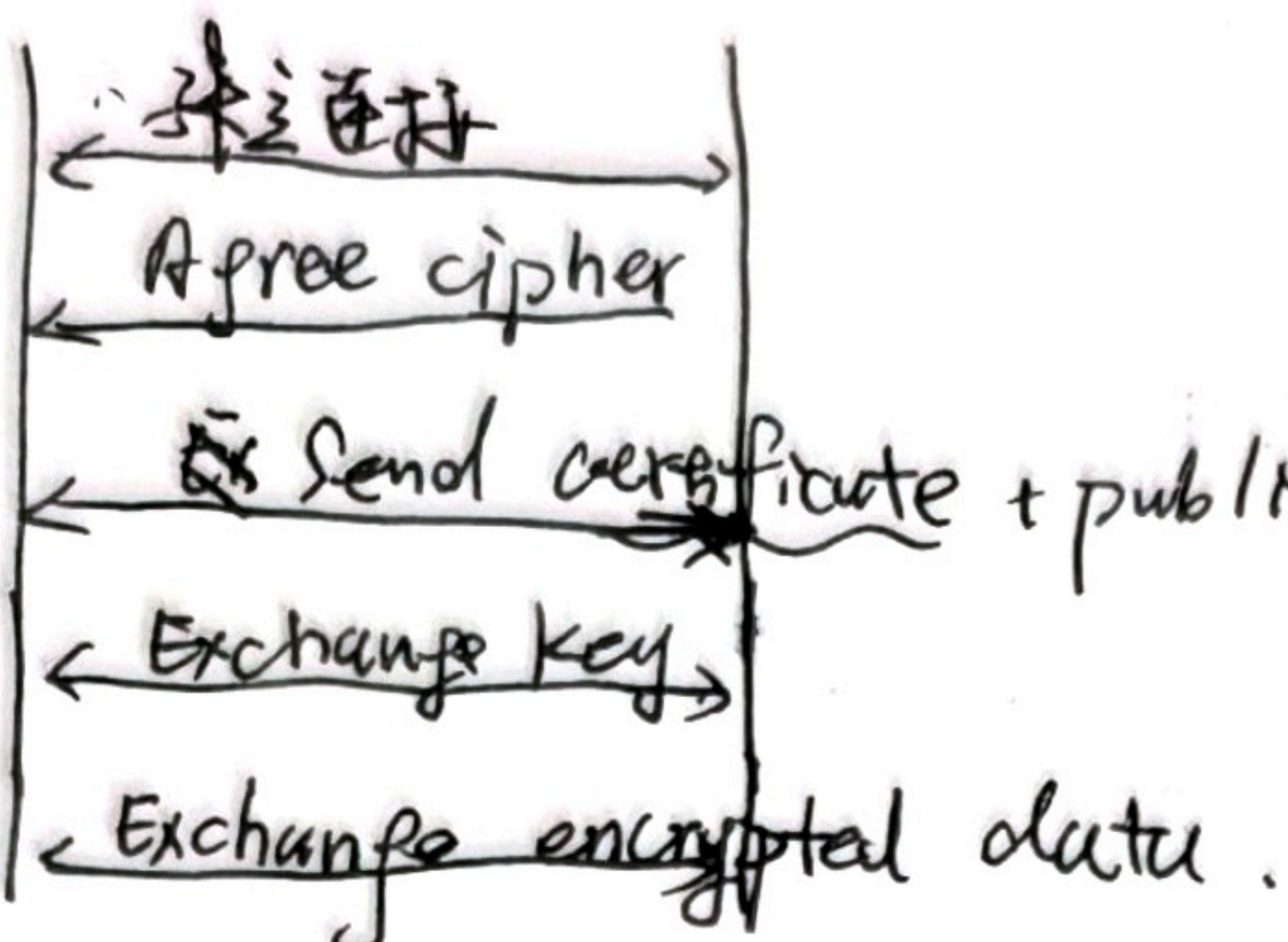
track status of every TCP connection (SYN, FIN, SYN ACK
 timeout inactive connection).

③ Application gateways.

require all telnet users to telnet through gateway, relays data between connections

2. HTTPS security \rightarrow TLS (Transport Layer Security) } correct website

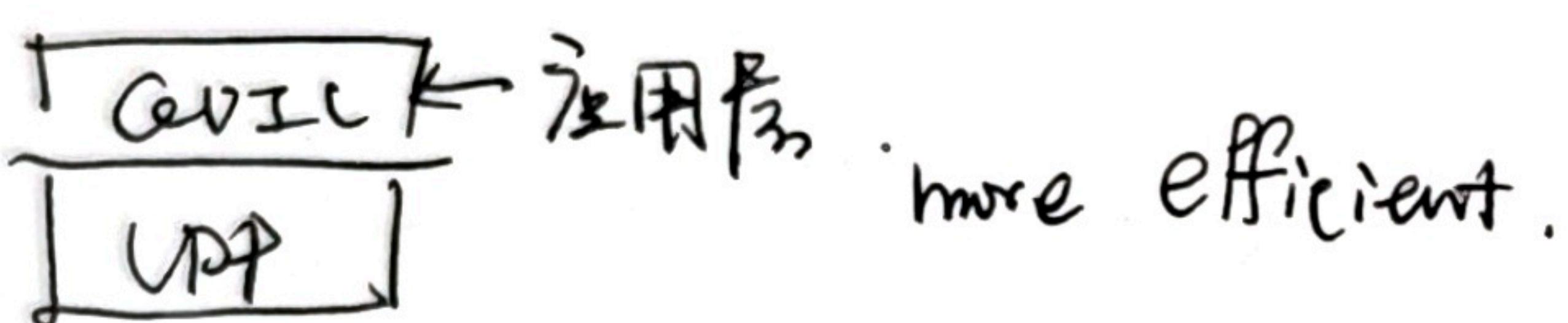
TLS



} secure connection.

send certificate + public key
exchange key
exchange encrypted data.

3. 快速 \rightarrow UDP 与 加密协议 CEVIC.



Asymmetric encryption (use verified
public key for secure exchange)
symmetric encryption (for fast
exchange of data).

4. Network ethics question

informed consent, consent but ~~not~~ not informed consent

some details of uses may be accessed by others.

human subject (harm to others).

agree and understand.

<Block 4>. Data Link layer.

1. Services.

framing, reliable delivery between adjacent nodes.

flow control, error detection, error correction.

<where?> \rightarrow 网络接口控制器 (NIC) \rightarrow network interface card.
硬件实现.

2. Error detection.

奇偶校验 bit.

Parity checking

Internet checksum.

CRC.

3. multiple access links, MAC protocol
 partition channel. } TDMA, time division multiple access.
 } FDMA, frequency

<Random access protocol>

① Slotted ALOHA.

retransmits frame in each subsegment slot with prob p until success.
 full rate, 简单, decentralized.

② Pure ALOHA

无碰撞, 效率极低.

③ CSMA (carrier sense multiple access).

collision detection \rightarrow propagation delay

④ CSMA/CD (Collision detection)

发送时若其他也在发送, 立刻停止, 发送 jamming signal.

停止之后, NIC执行 binary backoff time.

与 m 以互斥, 重发等待时间 t = Random $\{0, 1, 2, \dots, 2^m - 1\}$.

(\geq 争用 (512 bits times))

帧的传播时间大于碰撞时间, 避免产生 early collision.

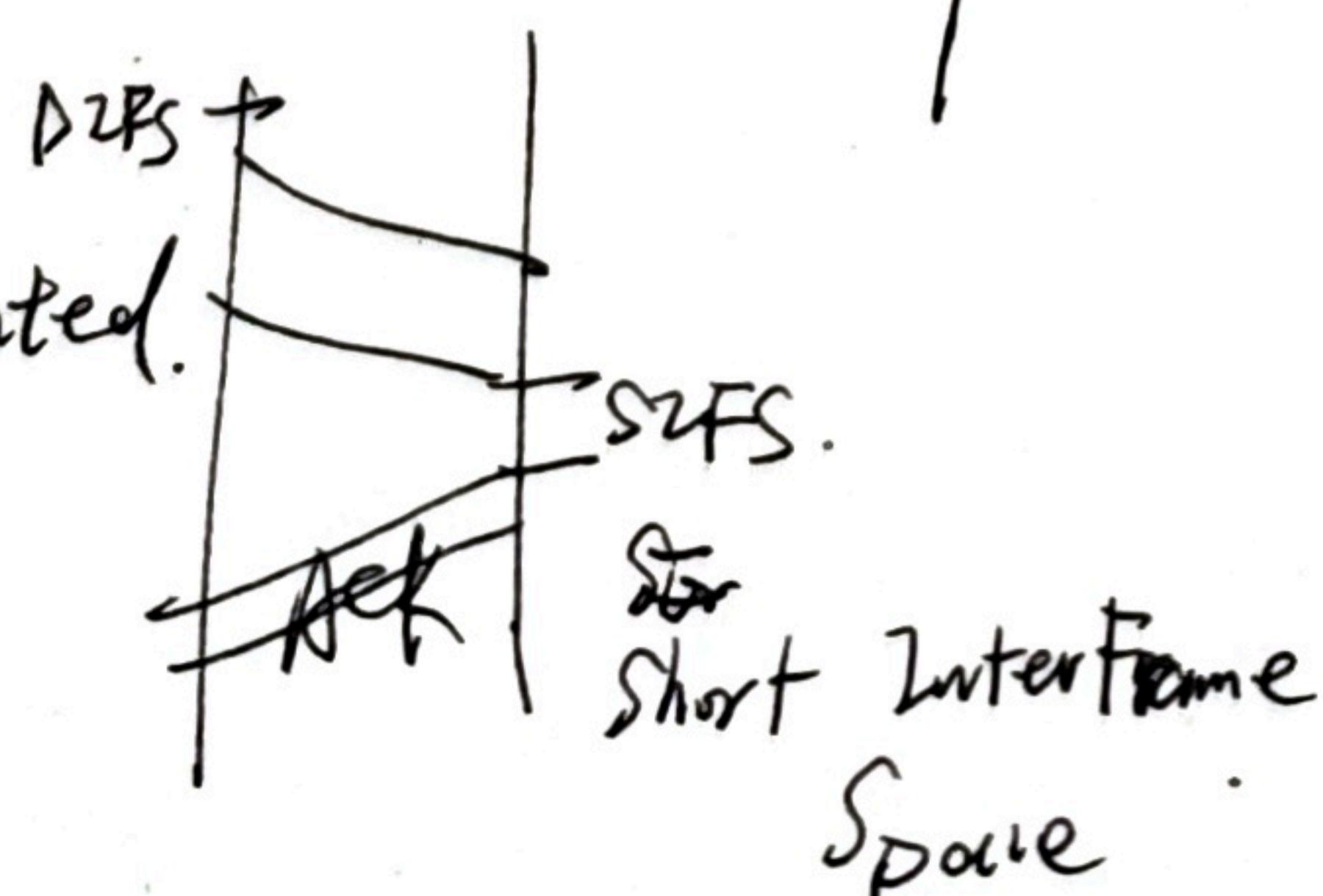
minimum frame size = $\overbrace{\text{RTT} + k}$

$2 \frac{k}{v}$ collision window

⑤ CSMA/CA (Collision avoidance).

无碰撞检测 } hidden station / terminal.

fasting weak ~~recent~~ received signals.



RTS, CTS \rightarrow 解决隐藏节点问题.

Request to Send Clear to Send.

部署于广播信道

4. MAC address.

48 bit, 8 Byte, ~~永久~~。

MAC 地址 - NIC (adapter)

MAC → never changes, flat address, portability.

IP → can changes, hierarchical.

<ARP, Address resolution protocol>

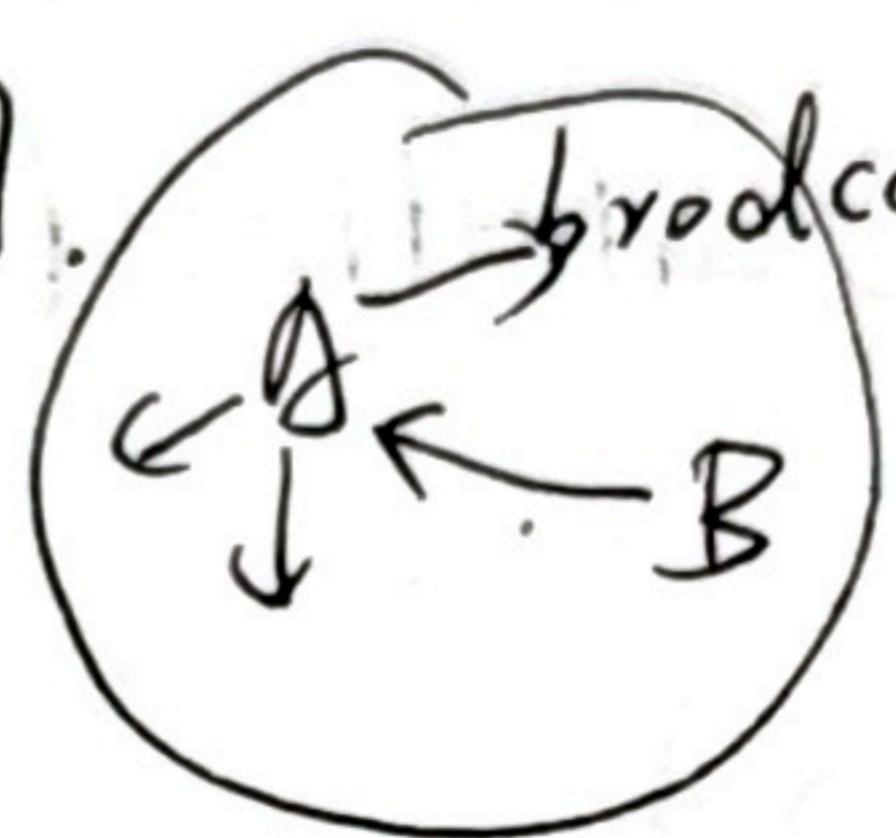
ARP table. $\langle IP, MAC, TTL \rangle$.

局域网 - LAN, ~~局域网中使用~~

~~局域网~~

IP $\xrightarrow{\text{ARP}}$ MAC

query



broadcast ARP query packet (B's IP)

5. Ethernet.

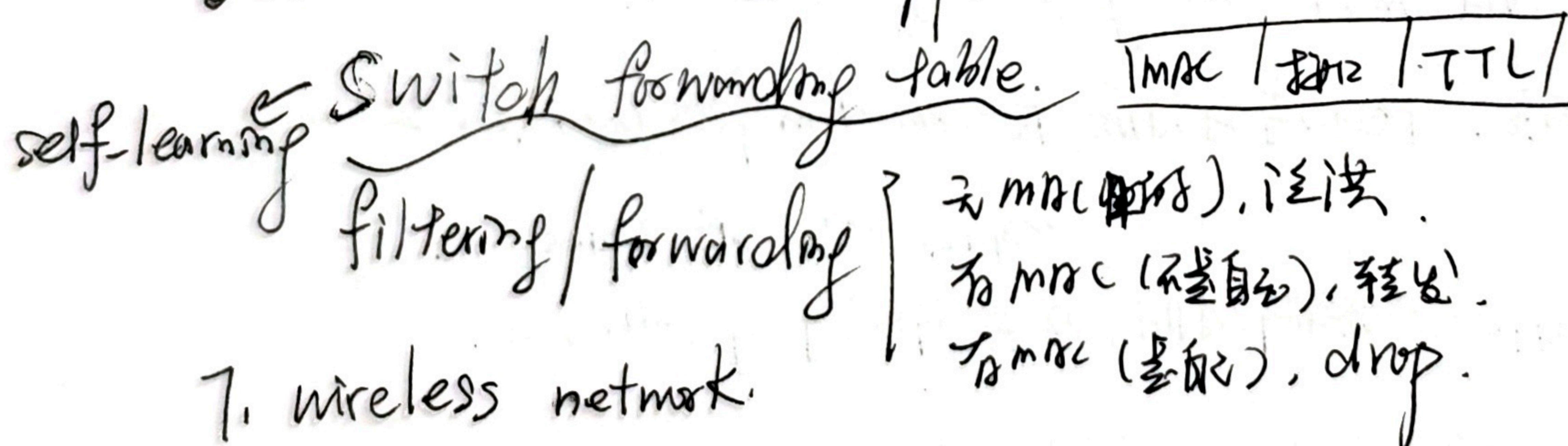
preamble	dst mac	src mac	type	data	crc
----------	------------	------------	------	------	-----

higher layer protocol.

synchronize receiver, sender clock rates.

connectionless, unreliable/unslotted CSMA/CD with binary backoff.

6. Ethernet switch. (Transparent).



7. wireless network.

infrastructure mode / ad-hoc mode.

无线 1, 2, 3 } \rightarrow 802.11 frame.
接续 空中 接收器发送

802.11 和 802.3 相似, 对于 ~~主机~~ 来说, 但是更简单。
路由器

只有 host 才能听到 AP 和有线

8. Virtual Local Area Network

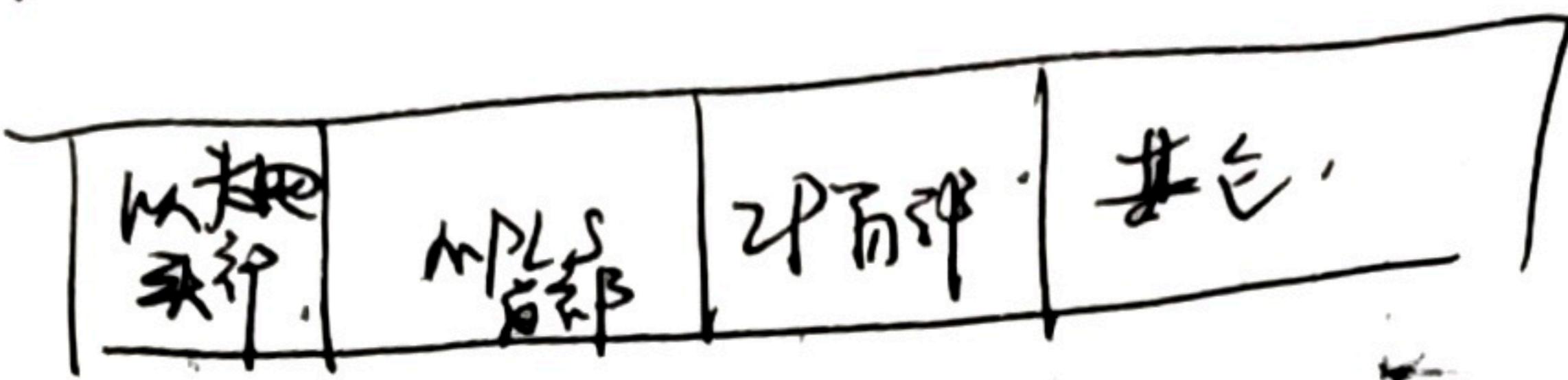
① port-based VLAN 不同端口 属于不同 VLAN

② trunk port. 802.1Q 标记 VLAN.

f. MPLS (Multi-protocol label switching).

使用固定 label 进行 IP forward.

MPLS - IP 平替.



MPLS router

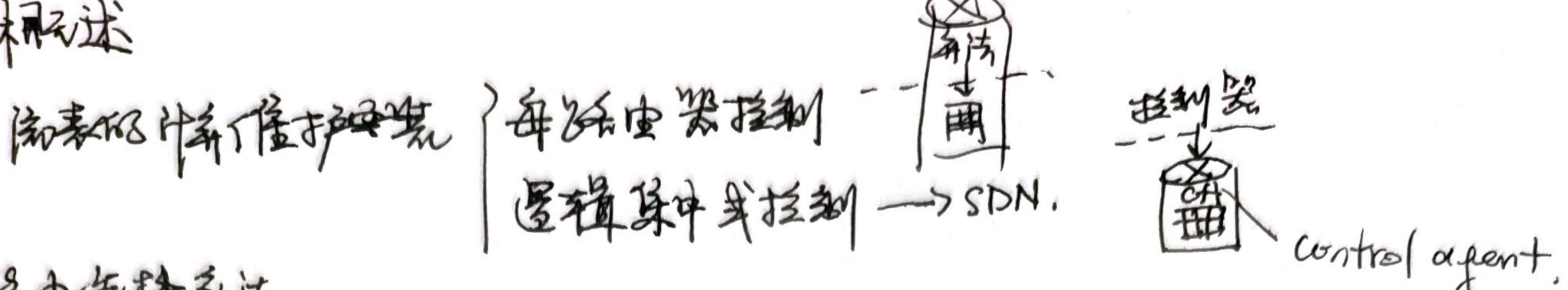
代替 IP，只通过 MPLS 路由器.

MPLS → IP 和同时一起.

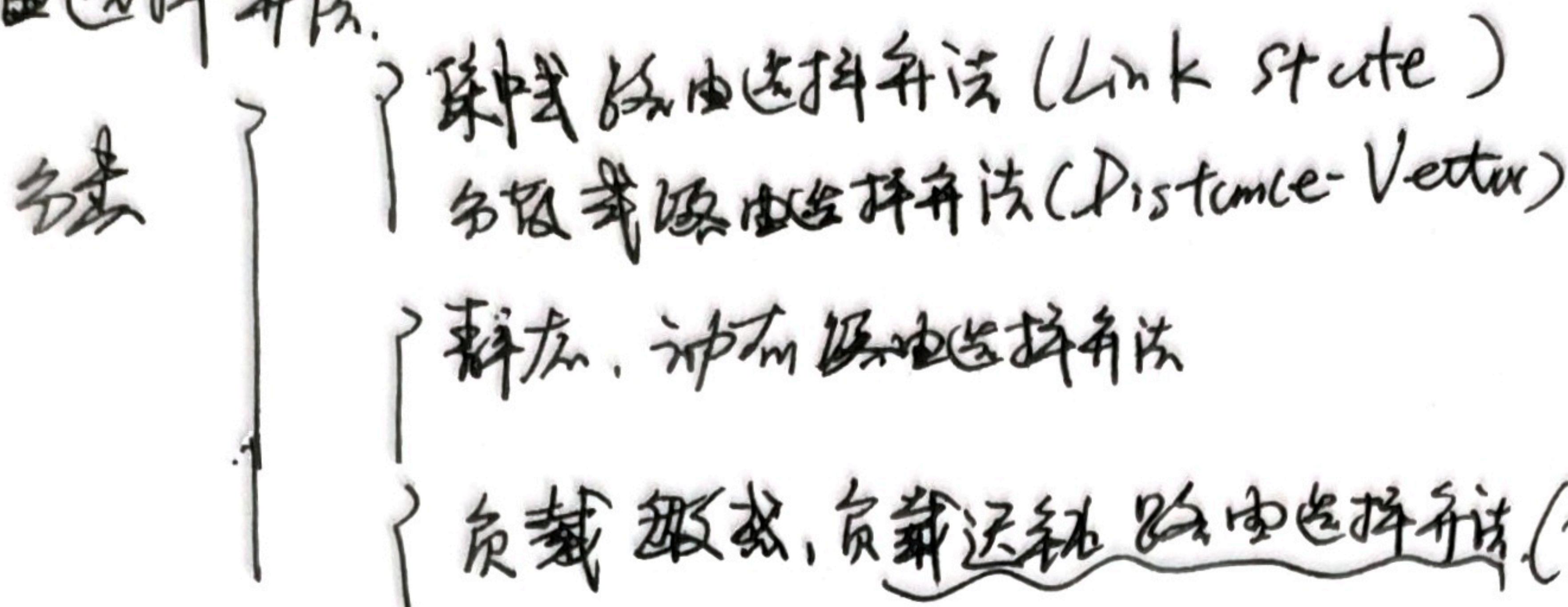
MPLS forwarding table.

<网络层·控制层面>

一、概述



二、路由选择算法



1) 链路状态路由选择算法 (LS).

所有网络状态已知 \rightarrow 每个节点广播链路信息.

Dijkstra 算法， $\forall v \in V$ 为一个初始节点，找到达到所有节点的最近路径.

Input: 初始化邻节点的状态 ($D_{vw} = c(u, v)$, $D_{uw} = u$)

Loop: 加入最小 D_{vw} 的候选集

使用 $D_{vw} = \min\{D_{vw}, D_{wv} + c(w, v)\}$ 更新节点
until: N' 为全集.

时间复杂度 $\frac{n(n+1)}{2}$

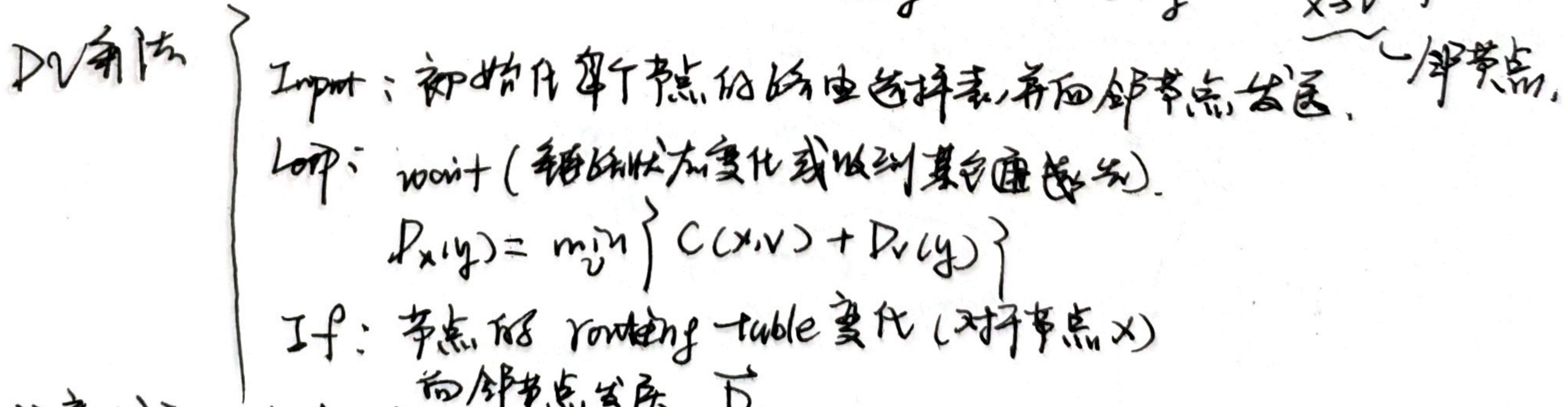
使用完后，将链路状态信息广播发表.

缺点：同步状态下，拥塞敏感振荡 \rightarrow 路径固定时间随拥塞.

2) 距离向量路由选择算法 (DV).

迭代、异步、自修正、分布式.

$$\text{运动方程} \quad \text{Bellman-Ford equation: } D_{x,y} = \min_{v \in V} \left\{ D_{x,v} + c(x, v) \right\}.$$



(注意发送的条件，发送的内环也不算，因为表是一个向量，当然每个节点也只能有一个本节点的反向)

按理说应该开辟了一块有用的存放路径不确定转发表，本节未提及

缺点：链路开销增加并协作，增加侧造成 routing loop，可以使用毒性逆程处理，但是节点多起来时并不实用.

LS 和 DV 比较

1. 简单性及收敛速度: LS $O(N^2) > DV$ (无条件收敛). 2. 可靠性: LS 完整性提供健壮性, DV 由一致错误通过整个网络. 3. 报文复杂性: LS 使用 $O(1/N^2)$ 占报文 (全广播) DV 只需部分点, 支持报文.

三. Internet 中自治系统内部路由选择 (OSPF).

不同网络相接方法使用同一种方法, 因为网是 IP 网络, ISP 有自己的想法.

一般一个 ISP - 一个 AS, 所属自治系统 (Autonomous System, AS).

一个 AS 内运行的路由选择方法 \rightarrow 自治系统内部路由选择协议 (intra AS ...) \Rightarrow OSPF

开放最短路径优先 (OSPF), 链路状态协议 } 泛洪链路状态信息 } IS-IS
...
Dijkstra. } 路由器向 AS 其他 router 广播
... route 选择信息, 未变化也同
期性广播

优势: 主要, $n \rightarrow V$ 多条路径, 单播多播支持, 支持 AS 层级分段. } 使用主干段对
一个 AS 分区.

AS 间路由选择协议 $\xrightarrow{\text{Internet}}$ 边界网关协议, BGP (Border gateway protocol).

$\xrightarrow{\text{Inter-AS}}$ 在 BGP 中, 各组路由器生成一个 CIDR 前缀, 即子网表 (S, I).
 $\xrightarrow{\text{CIRD}}$

1) 通告 BGP 路由信息.

AS 路由器 \rightarrow

网关 (gateway router)	BGP Connection	eBGP
内部 (internal router)		

 并不等于物理
链路对应.

gateway router 通告 eBGP 收到, 再以 iBGP 对等 AS.

2) 确定最佳路由.

路由器 (route) $\xrightarrow{\text{前缀}}$
 $\xrightarrow{\text{BGP attribute}}$ AS-PATH \rightarrow 通告过的 AS 列表
 $\xrightarrow{\text{NEXT-HOP}}$ 逆向传播过程, 跳过当前 AS 最近的 AS 为
 选择办法.
 $\xrightarrow{\text{forwarding}} \text{forwarding} \text{ router 的 IP 地址}$
 $\xrightarrow{\text{指向的端口}}$

a) 土豆土豆路由选择 (hot potato routing)

只关注 AS (之前) 由邻居提供路径最优.

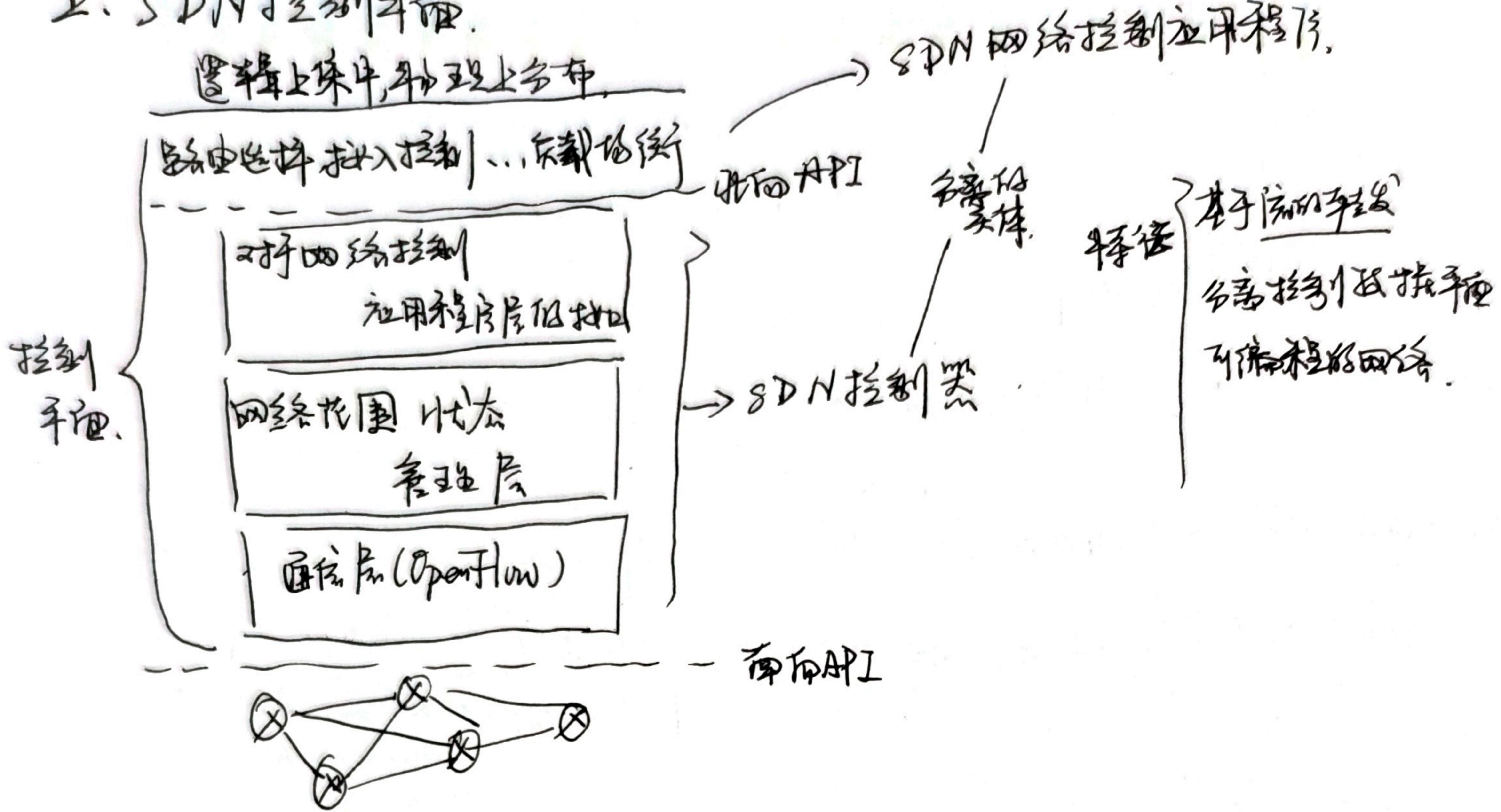
b) 路由器选择办法

本地偏好度量 \rightarrow 路由选择策略, AS-in routing 主导作用.

自下而上
的
选择策略.

最近 AS 路径
跳数
BGP 协议执行.

五、SDN 控制平面



六、ICMP：因特网控制报文协议

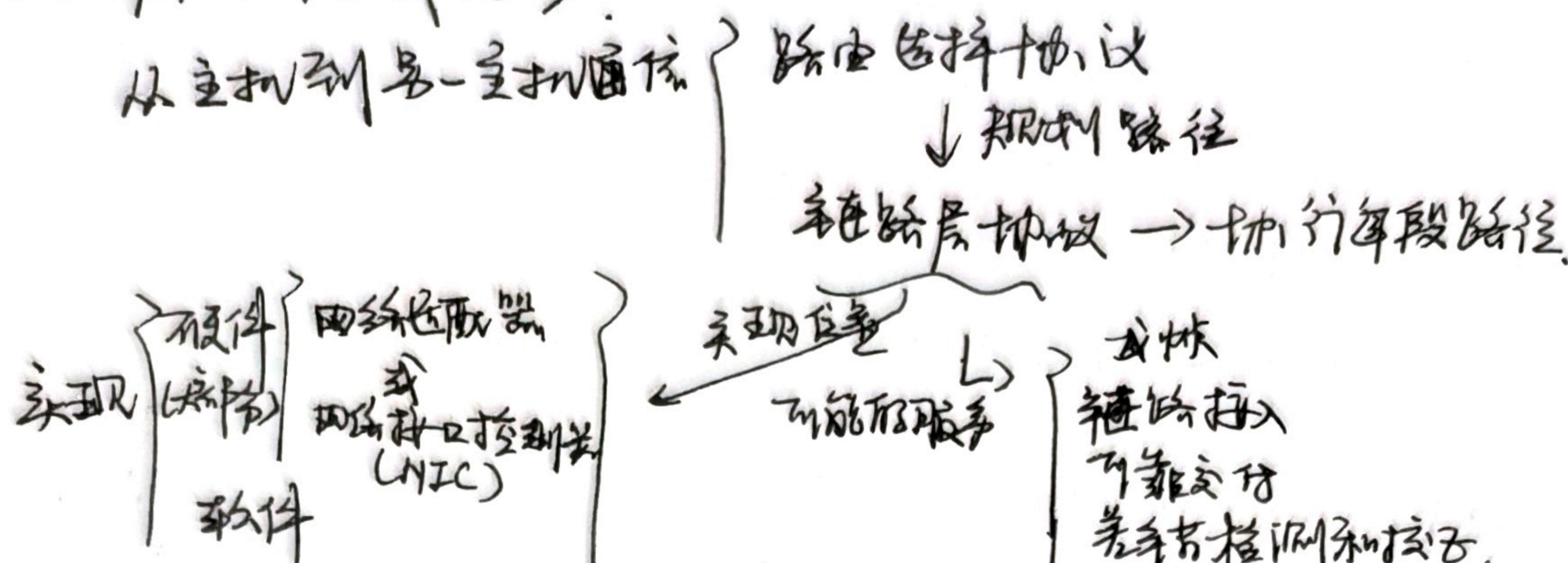
ICMP，被主机和路由器用来彼此向通网络层信息。→ 差错报告。

ICMP 在 IP 的报文载荷中，有类型字段和标志字段

二者匹配即可往复应答。

Traceroute 程序使用 ICMP 实现，通过增大 TTL 字段，获得被丢弃的 ICMP 到达目的主机时，因为空包不回送 ICMP，故得 ICMP。

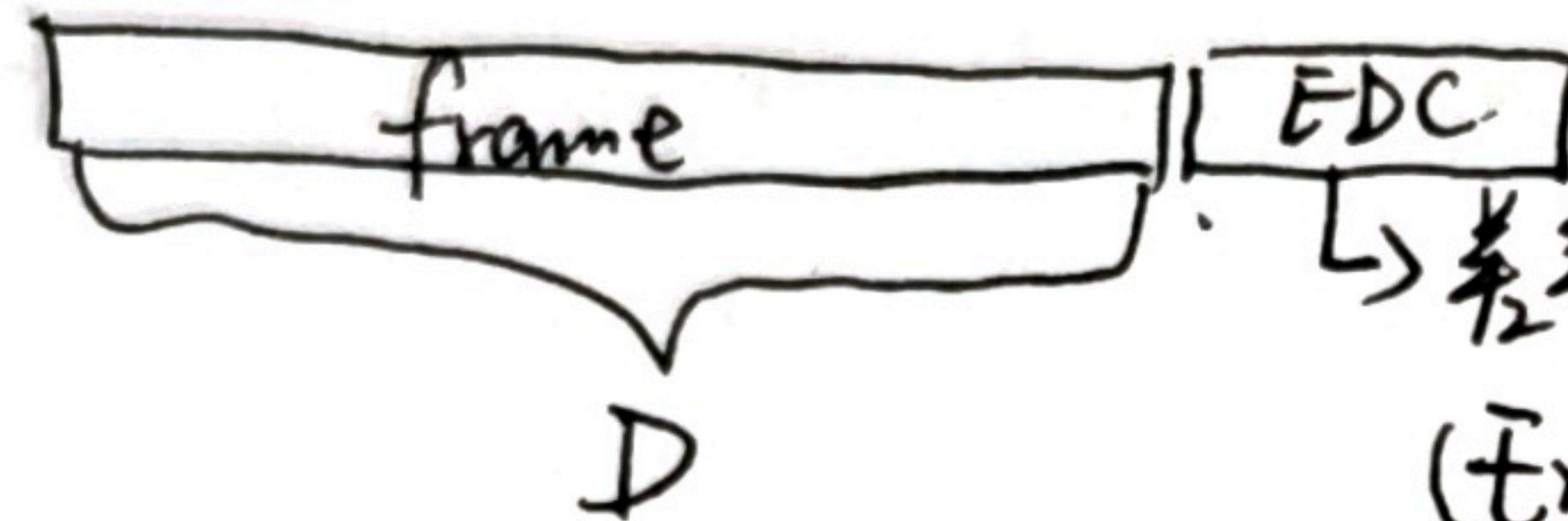
<链路层和局域网>



一、差错检测和纠正技术 (纠错码)

增强数据:

(发送方)



差错检测和纠正

(Error-Detection and -Correction, EDC)

接收方也可能查不到错误，即 未检测出比特差错 (undetected bit error).

其检测和纠正能力为 前向纠错 (Forward Error Correction, FEC)

1) 前向校验.

parity bit → 不对称检测 个别bit error

two-dimensional parity → 单个bit 错, 检测到纠正

→ 行, 列, 行列交错. 两个bit 错, 检测不对纠正.

2) 检验和.

$d \text{ bit} \rightarrow k \text{ bit } (x n) \xrightarrow{\text{全加 (四进制)}} \text{ 反码} \rightarrow \text{Checksum.}$

3) 循环冗余校验 (Cyclic Redundancy Check) → 全为 1, 出现 0 都错误.

即 CRC 循环冗余校验 (CRC) 编码, 也有多项式编码.

发送方设置 CRC

$$R = ((D \cdot z^r) \bmod G)$$

所有位都是 \oplus XOR.

r bit 系数 d bit 左移 r 位, 部分余数 $r+1$ bit

$d+r$ bit

接收方收到 $D|R$, 拆出 DR 的 $d+r$ bit, 使用 G去除.

则 $\frac{d+r}{G}$ 系数 } $D \rightarrow$ 正确
 other → 有错误.

三、多路访问链路和协议 (MAC协议)

broadcast link → 多路访问问题 (multiple access problem).

局域网 → 地理上集中在一个区域的网络.

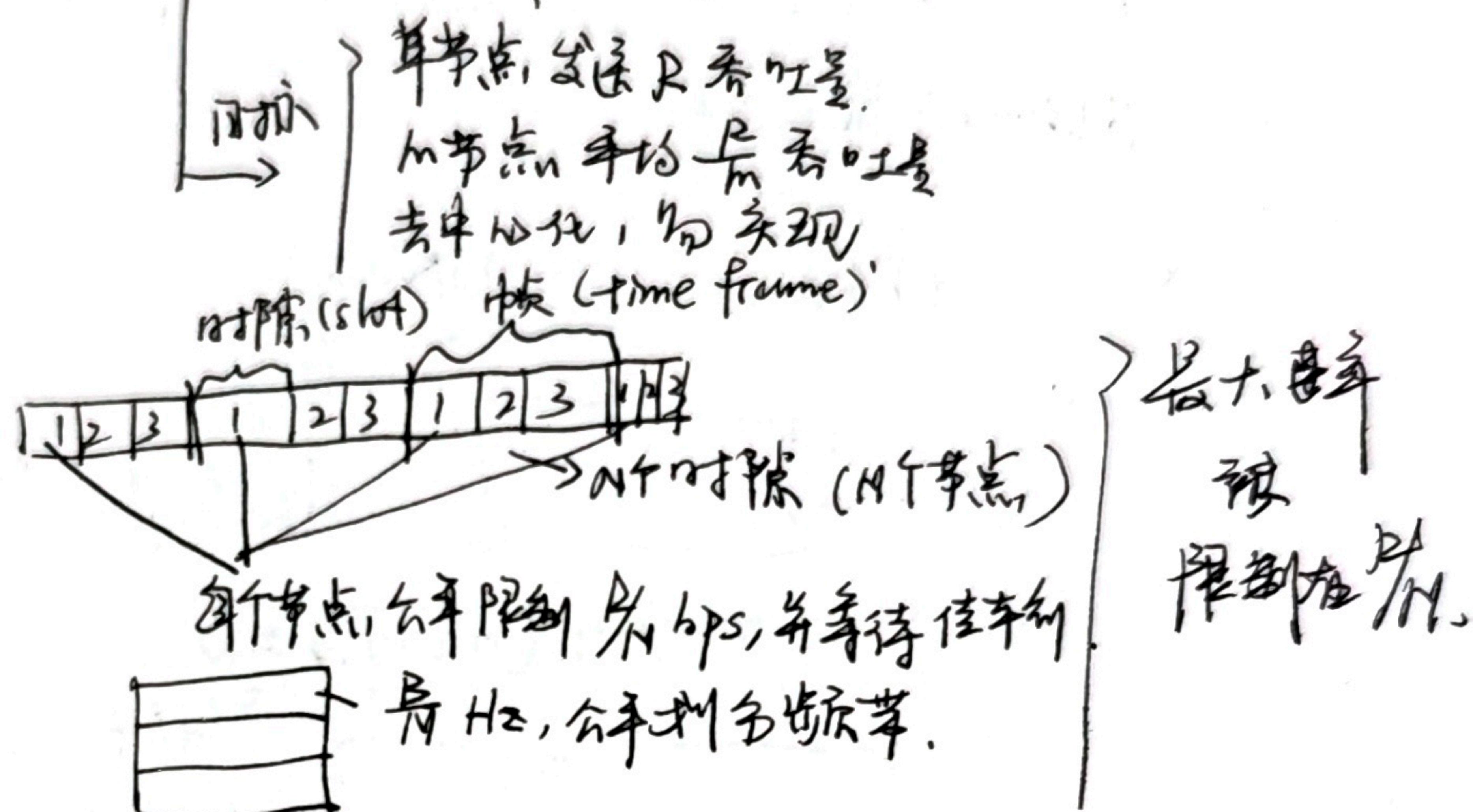
为解决广播碰撞 (collide) → 多路访问协议 $\xrightarrow{\text{方法}}$

channel partitioning
random access
parking-turns.

1) 信道划分协议.

时分多路复用 (TDM)

频分多路复用 (FDM)



码分多址 (Code Division Multiple Access, CDMA)

2) 随机接入协议.

一个节点只会以信道全部速率发送, 碰撞不等待随机时延

a) 对称ALOHA.

发生 collide, 以 P 概率碎片在每个时隙中重传, 成功为止.

$$\text{效率} \quad Np(1-p)^{N-1} \xrightarrow{\text{limit}} 0.37.$$

b) ALOHA

相较于 ALOHA, 取消时隙相隔后, 信号即随时开售. Collide 后, P 报文立刻重传, (1-P) 报文等待一个帧佳车轮时间, 再 P 报文决定.

$$\text{效率} = Np(1-p)^{2(N-1)} \xrightarrow{\text{limit}} \frac{1}{2} \cdot 0.37.$$

c) 截波侦听多路访问 (CSMA).

Carrier Sense Multiple Access → 只有节点未传输 (或尚未收到) 自己才传输.

丢站到端时延 (channel propagation delay) 影响.

d) 具有碰撞检测的载波侦听多路访问 (CSMA/CD)

Collision detection → 发生碰撞立刻停止, 停止一个随机时间量.

二进制指数后退 } 经过 n 次碰撞后, $K = \{ \text{Random} \{ 0, 1, 2, \dots, 2^n - 1 \} \}$ 变量.

$$\text{随机时间量} = K \cdot (512 \text{ bit} + \text{传输时间})$$

$$\text{效率} = \frac{1}{1 + \frac{1}{2} \cdot \frac{K}{2^n} \cdot \text{frames}}$$

3) 轮流协议.

a) 轮询协议 (polling protocol), 主节点轮流向各个节点并指定最大帧

b) 令牌传递协议 (token-passing protocol), 传递令牌, 节点持有令牌为传递指定帧

节点损坏? 令牌丢失? 都是问题.

不超时
令牌最大帧.

三、交换局域网.

1) 线路层寻址和ARP.

a) MAC地址

Medium Access Control, 介质访问控制 地址 → LAN address
physical address
MAC address.

正如IP地址对于接口, MAC地址对应路由器和主机的NIC.

(线路层交换机无MAC地址).

MAC地址为6字节, 用十六进制对, 一共六对表示.

MAC具有扁平结构, 不可变, IP有层次结构, 可变.

对于局域网上任何MAC broadcast address = FF-FF-FF-FF-FF-FF.

b) 地址解析协议.

IP ← Address Resolution Protocol, ARP → MAC.

ARP只在同一个子网上起作用和路由器不能解析IP.

每台主机或router都有自己的ARP table, 有IP, MAC, TTL.

获取table方法 } 发送方 ← 标准映射到IP地址 → | IP | MAC | TTL |



广播范围..

ARP查询多但响应应与请求一致, 会将发送接收IP地址和MAC地址.

ARP报文封装在链路层帧中, 包含MAC, IP. 跨越不同局域网和链路层的协议.

c) 发送数据到子网以外.

主机向另一子网发送数据报, MAC地址并不相同的主机对应的MAC地址, ARP只在同一个子网中使用, 对于该子网, 并没有MAC对应的IP, 会丢弃.

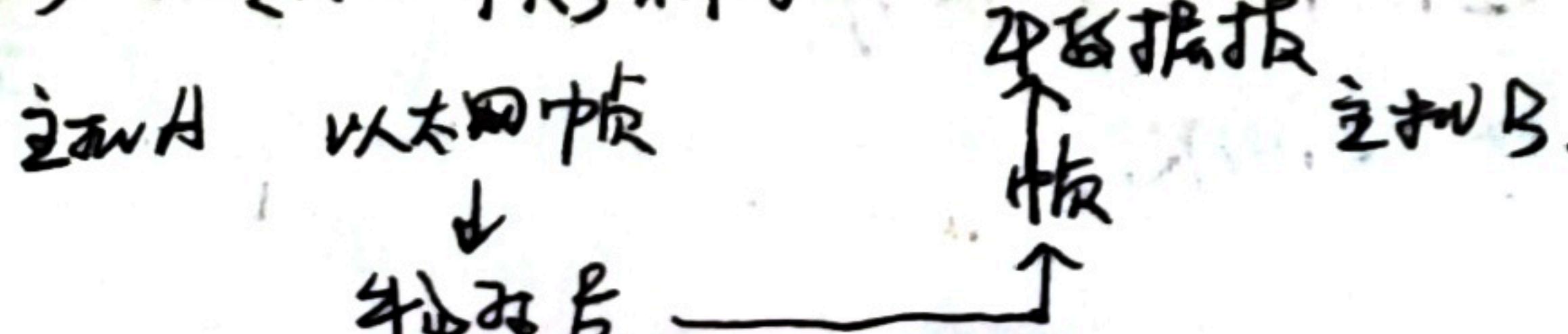
主机 帧 → 路由器接口 → 另一接口 → 帧 → 第二主机

2) 以太网.

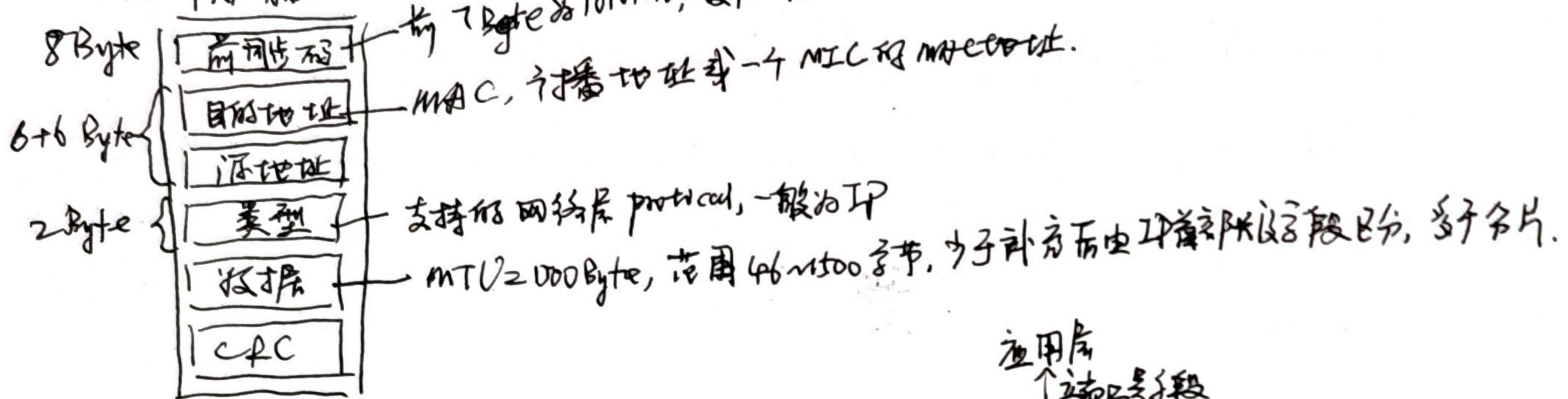
ARP → 路由器接口 MAC.

交换以太网, 无碰撞, 有源转发包过滤交换机.

a) 以太网帧结构.



以太网帧结构：



以太网提供无连接服务(面向网路层).

(b) 现代以太网或许已经需要 MAC 协议, 但其帧结构
(因为交换机可以有转发缓存). 保持不变.

应用层
一个端口号字段

运输层
一个协议字段

网络层
一个类型字段
链路层.

3) 链路层交换机.

对于两个主机和路由器, 交换机是透明的.

a) 转发和过滤

借助于交换机表 (switch table) 来完成交换机表查询为:

地址 (MAC)	通往该地址的接口	表项过期时间
----------	----------	--------

到一个目标 MAC 地址进行交换, 3 种情况:

1. 表中无 MAC 地址, 而该接口之外, 另外有接口缓存广播该帧.
2. 查询到输出接口强输入的这个接口, 并执行过滤行为.
3. 有对应的其它接口, 而其接口仍缓存转发.

b) 学习 (stref-learning).

即插即用设备 (plug-and-play device).

1. 查表

2. 对于每个接口有以下信息: ① 目标 MAC 地址 ② 接口 ③ 时间

3. 一段时间 (老化期 aging time) 后, 支持地没有收到以 MAC 地址为目的地的帧, 则从表中删除这个表项.

交换机的接口能同时发送和接收, 全双工.

c) 性质.

1. 清除碰撞, 交换机不会在同一个网段上转发多个碰撞帧.

2. 异步媒体连接.

3. 分布式互连, 一个 NIC 通常发送帧, 断开并为另一个 NIC 服务.

d) 交换机与路由器.

交换机 > 优: 即插即用、过滤速度快.

缺: 防止广播帧 loop, 限制环形结构.

冲突域大, 处理慢.

对广播风暴不提供保护.

→ 用以太网, 不支持 IP 的配置.

路由器

> 优: 4 层以, 有冗余路径也不循环, 打补丁部署到生成树组

对第二层广播风暴提供防火墙保护

缺: 非即插即用, 处理时间长.

→ 大网关 (两个都用), 灵活性.

4) 虚拟局域网

交换机的缺点:缺乏流量隔离。

无法使用(没有充分利用接口带宽)

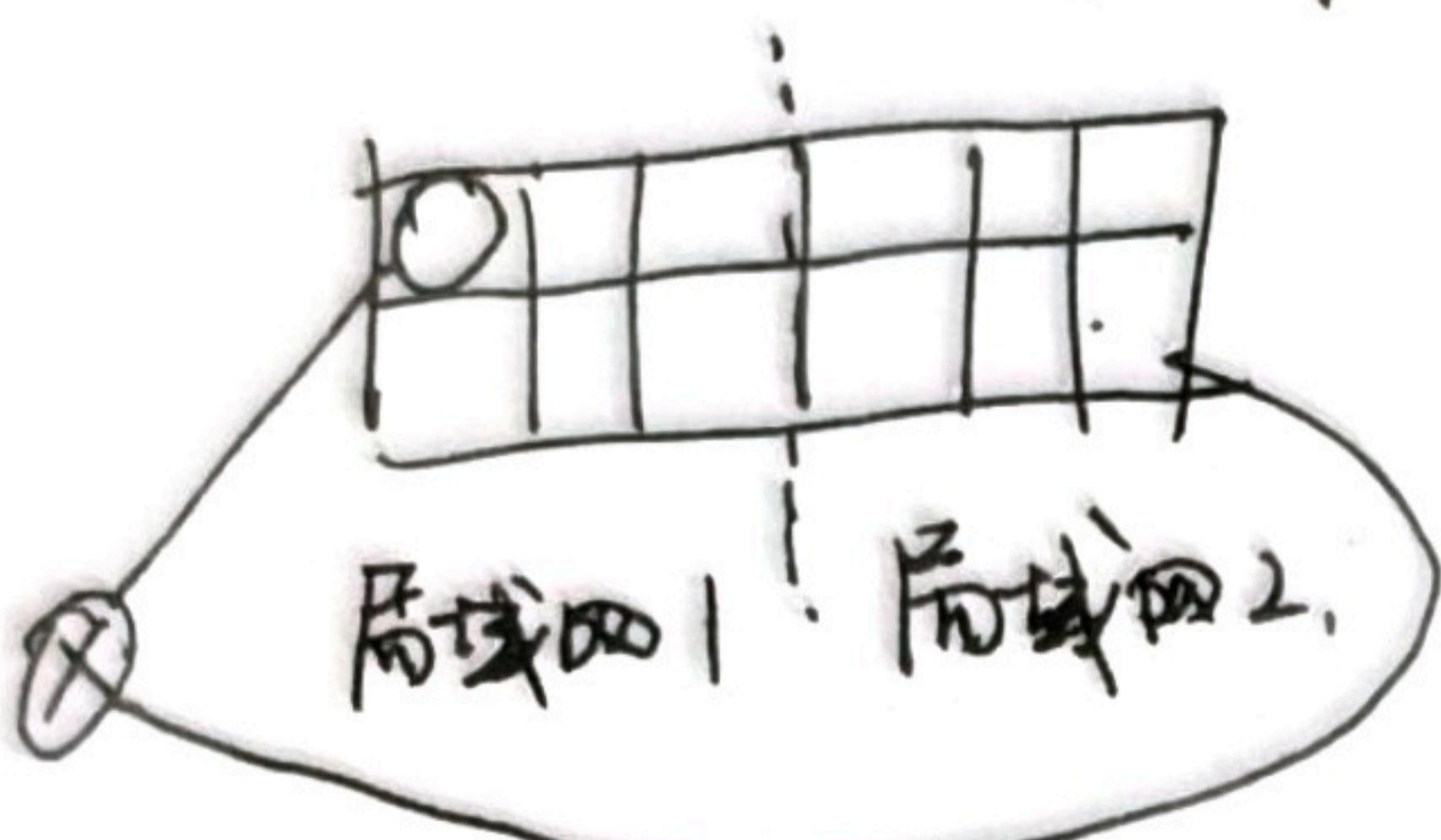
基础用户必须改变布局。

支持虚拟局域网

(Virtual Local Network, VLAN)

仍需交换机来解决。

VLAN交换机允许一个单一的物理局域网基础上分离出多个虚拟局域网。

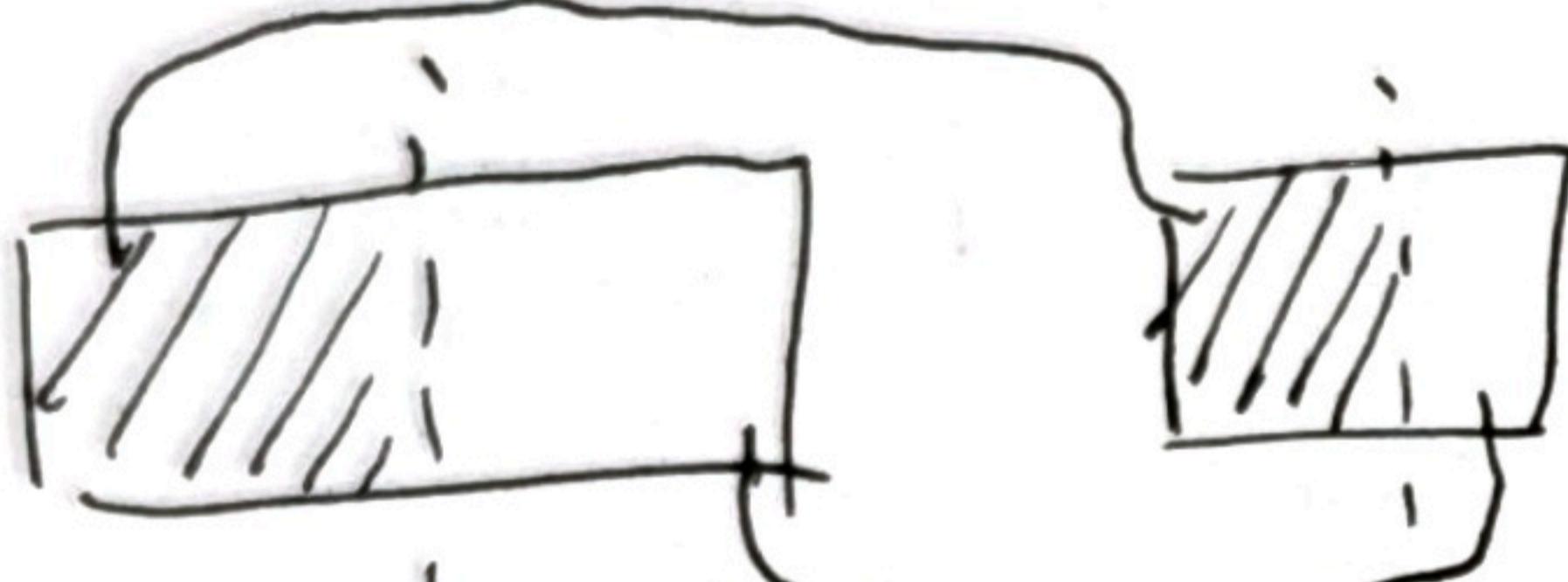


局域网 1, 2 只能有各自的内部链路和连接帧,

跨局域网传输则必须共享于 1, 2 的接口用路由来实现。

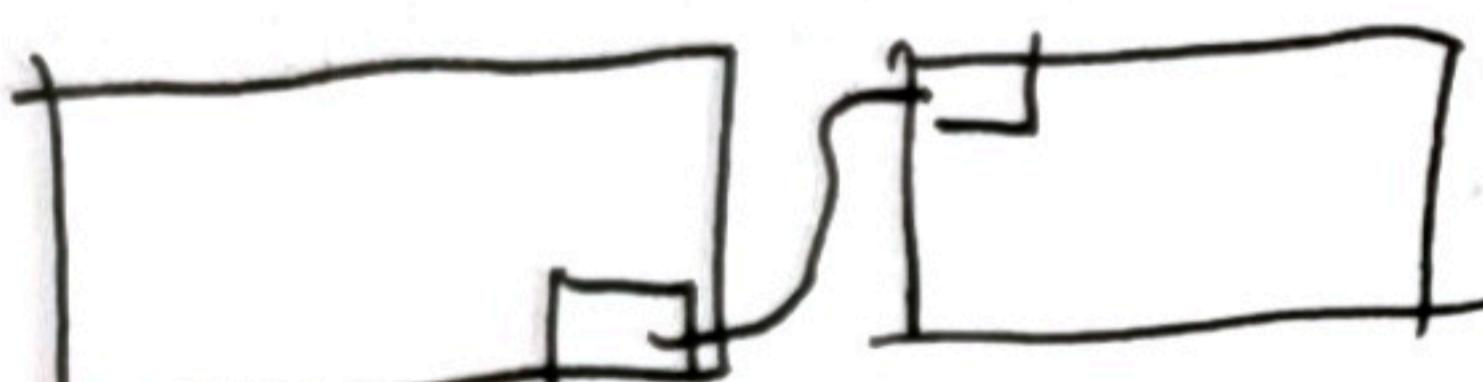
对于需求较多端口, 如何扩展呢?

1.



新增一交换机, 相同局域网互连
不具备扩展性。

2. VLAN trunking (VLAN 干线连接).

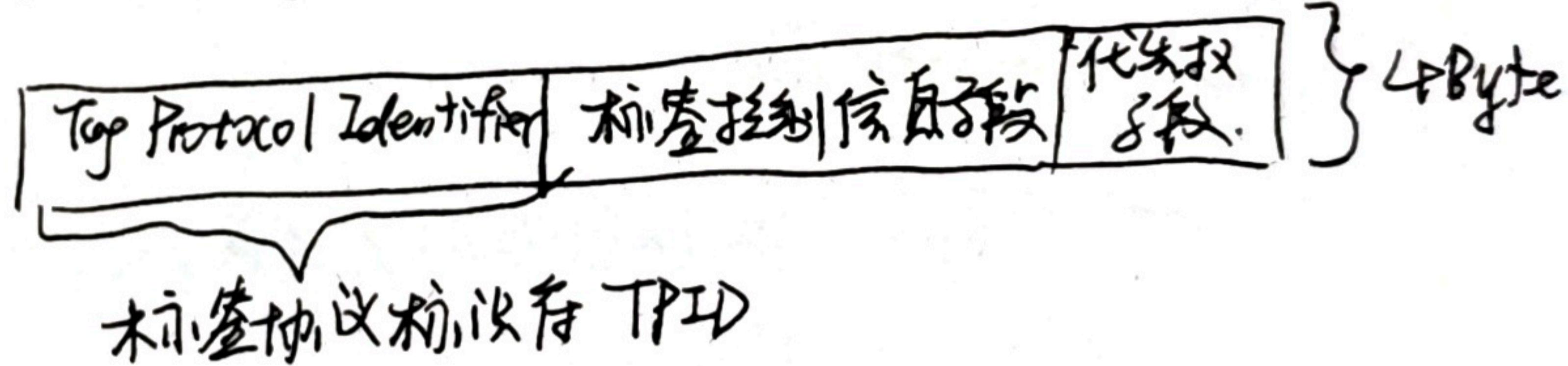


开辟一特殊端口, 属于所有的交换机所有的 VLAN。
如何判断该交换机的帧属于哪个 VLAN?

将以太帧头进行标记 → 802.1Q, 用于跨越 VLAN 的中继。

在标准以太帧中插入

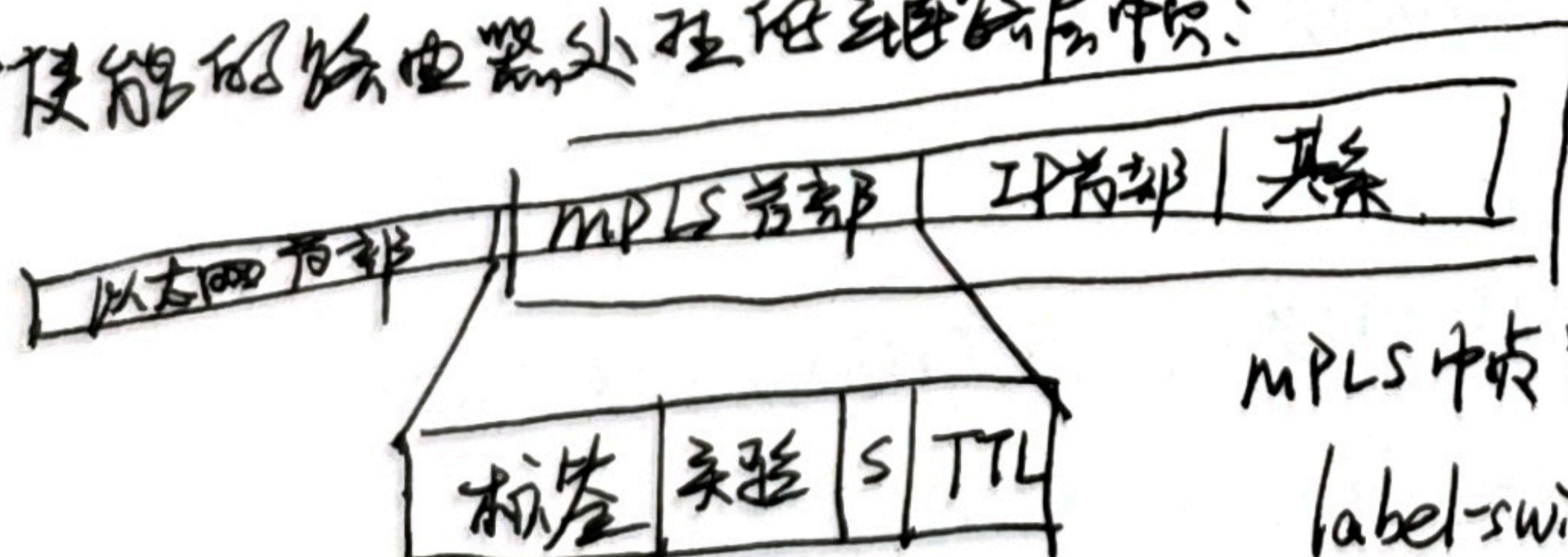
4 Bytes 的 VLAN tag.



三、链路虚拟化: 网络作为通路层。

多协议标签交换 (Multiprotocol Label Switching, MPLS), 为 IP 提供互联网的解决方案技术。

MPLS 使能的路由器处理链路层转发:



帧有效载荷。

MPLS 中只能在 MPLS 使能路由器 中使用。
label-switched router 标签交换路由器。

MPLS router 在转发中查找 MPLS tag, 然后转发, 不拆卸 IP 头, 不拆卸 MAC 地址。

基于标签执行交换, 不考虑 IP 头, 导致转发增加。

提供多条路径发布到流量工程

(traffic engineering)

Virtual Private Network, VPN.

将链路本身视为 MPLS 网络。

Internet

数据中的

数据包的 IP 网络 → 用户直接内部主机。

世界路由器 (border router)

主机 → TOF (blade), 托架顶部交换机 (Top of Rack, TOR)。

1. 负载均衡.

负载均衡器
load balancer

> 基于目的端口号和源的 IP 地址决策.

动作

向主机端口发请求，然后立即将负载.

类似于 NAT 功能，转换 IP，提供安全.

2. 等级集中架构.

数据集中应用 路由器和交换机等級结构(Hierarchy of router and switch)

并发现问题：不同机架内台式机间最大连接受限.

解决：①高速串接机架和路由器.

②相关服务和数据尽量靠近，以减少带宽.

③增强相邻交换机层级连接.

3. 发展趋势.

①成本降低. ④集中式 SDN 框架和模型. ③虚拟化

②物理分离. ⑤硬件模块化和定制化.

〈无线网络 和 有线网络〉.

一、概述.

wireless host, wireless communication link

base station 基站 → 接入点 (access point, AP), 并通过中继

建立与基站关联 → 基础设施模式 (infrastructure mode).

不关联 → 自组网无线 (ad hoc network)

主机直接和两个基站 → 切换 (handoff)

二、无线链路 和 网络特征.

1. 链路损耗 (path loss)

2. 共同干扰

3. 多径传播 (multipath propagation)

多径访问 } 隐藏终端问题 (hidden terminal problem) }

信号强度衰减 (fading)

⇒ CDMA 方，同时接收多路干扰，
由编码方式校正.

三、WiFi：802.11 无线局域网.

a) 体系结构.

基本构件 → 基本服务区 (Basic Service Set, BSS)



— Internet.

交换器或路由器

部署即插即用局域网为基础设施无线局域网
(infrastructure wireless LAN.)

(b) 信道关联.

每个无线站点在能发送/接收的服务层区域之前, 必须有一个时相关联.

AP > 服务集标识符 Service Set Identifier, SSID
信道号(共11个).

每个时间间隔发送(广播)信标帧(包含 SSID 和 MAC), 面对扫描信道周期. 时相关联.

passive scanning = AP发送信标帧.

站点, 而 AP 为关联请求帧

AP 向该站点发送关联响应帧.

active scanning = 站点主动搜索, 广播请求帧

AP 监听探测响应帧

站点, 而 AP 发送关联请求帧

AP 向站点发送关联响应帧.

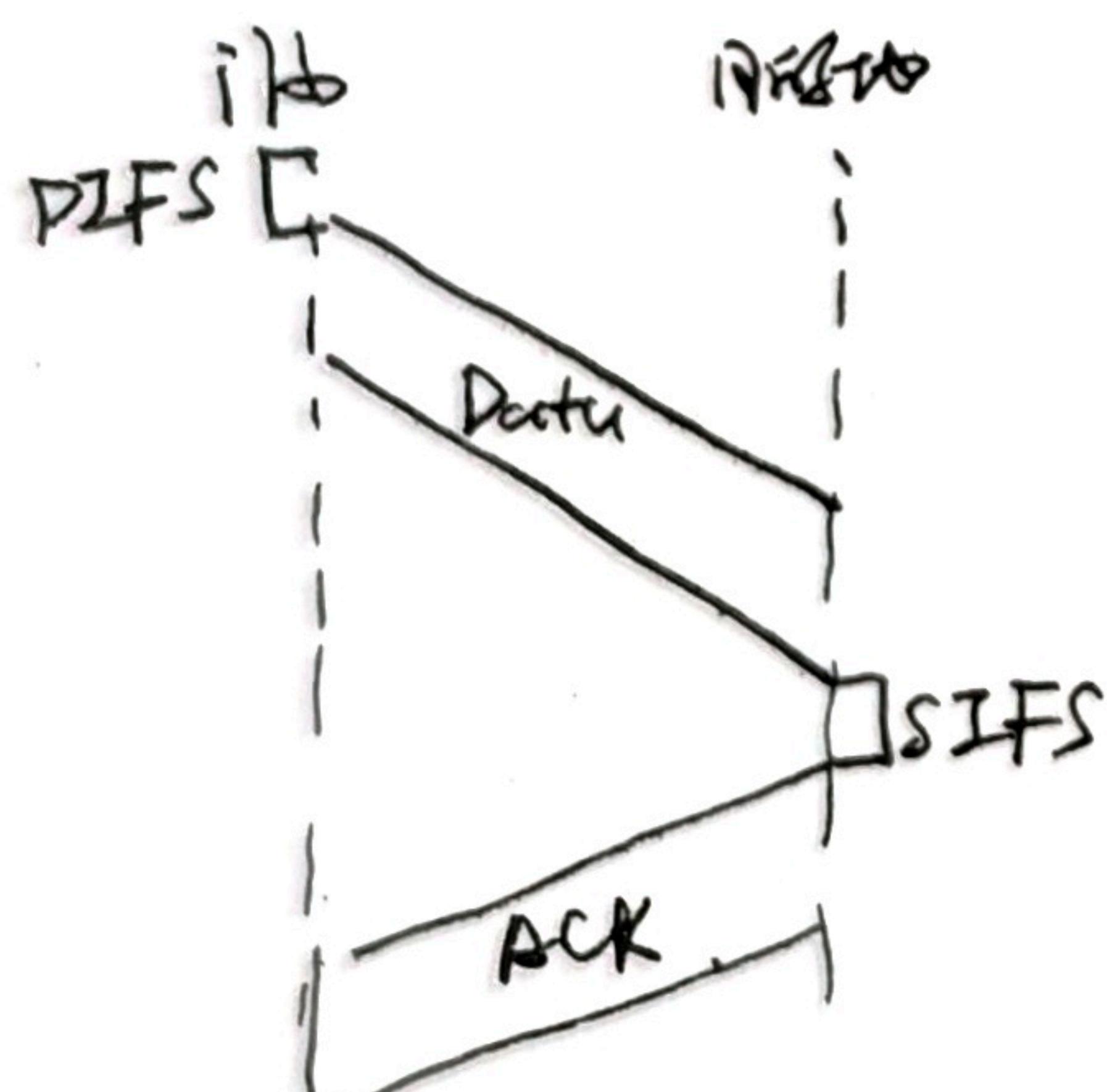
(c) 802.11 MAC 协议.

随机接入协议, 带碰撞避免的 CSMA (CSMA with collision avoidance, CSMA/CA)
比特差错率较高, 误用率较高, 不确认/重传 ARQ (Automatic Repeat reTransmit)
并未实现碰撞检测的原因:

1. 检测硬件代价大(具有同时发送和检测能力).

2. 室内隐藏终端和衰落问题无法检测所有碰撞.

① 碰撞避免.



link-layer acknowledgement

给定时间未收到 ACK 帧,
一直重传一直到 ACK, 放弃.

2) 处理隐藏终端问题: RTS/CTS.

分布式帧间间隔 (Distributed Inter-frame space)

短帧间间隔 (short Inter-frame space).

① 站点监听到信道空闲, 等待 DIFS, 发送.

② 监听不空闲, 丢弃并计数值, 在空闲时通过流程,
信道忙时计数值不变.

③ 计数值 > 0, 丢弃并等待确认.

④ 收到 ACK, 发送下一帧, 再计数值.

未 ACK, 选择更大范围计数值.

若站点相互隐藏, 无法检测.

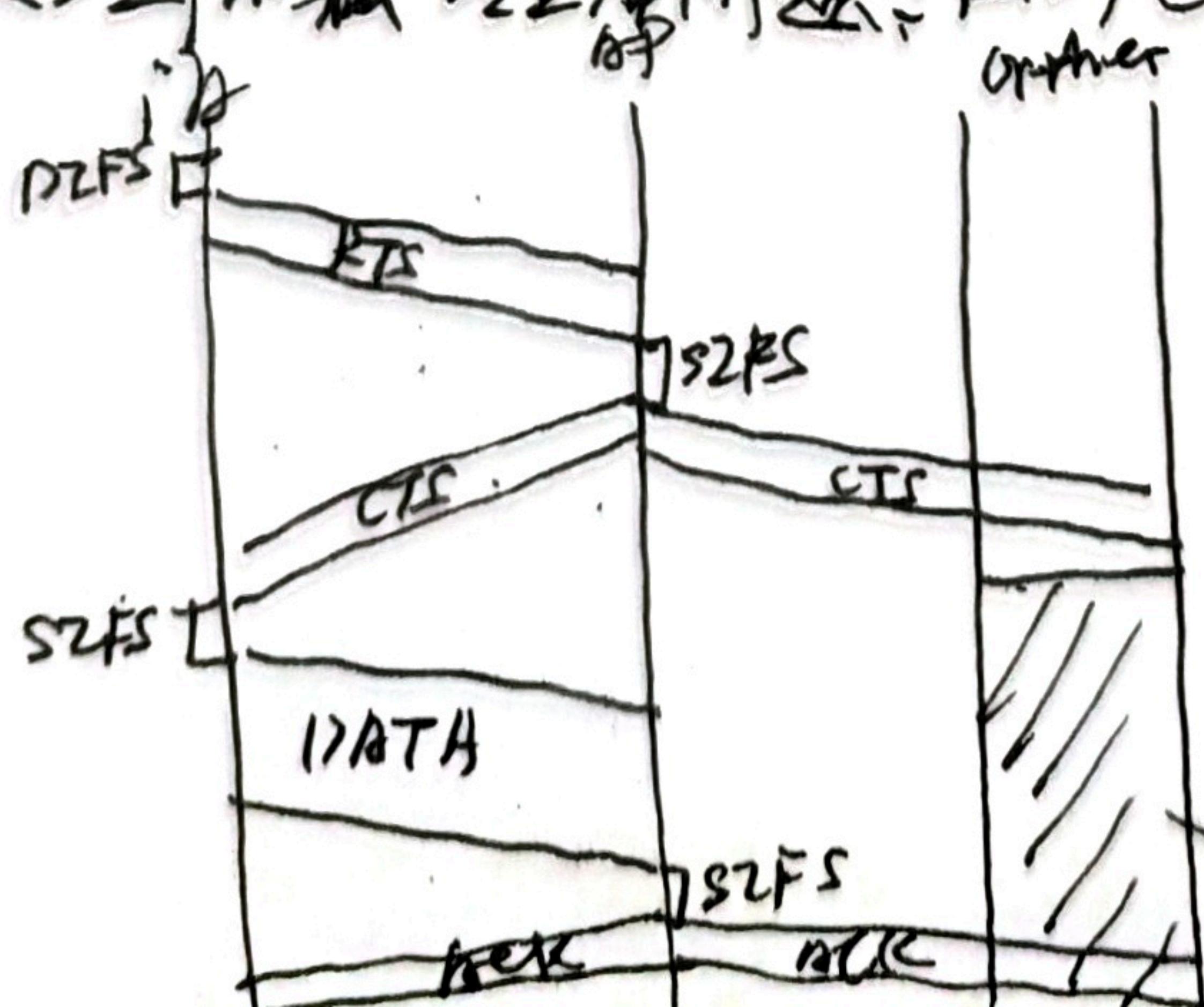
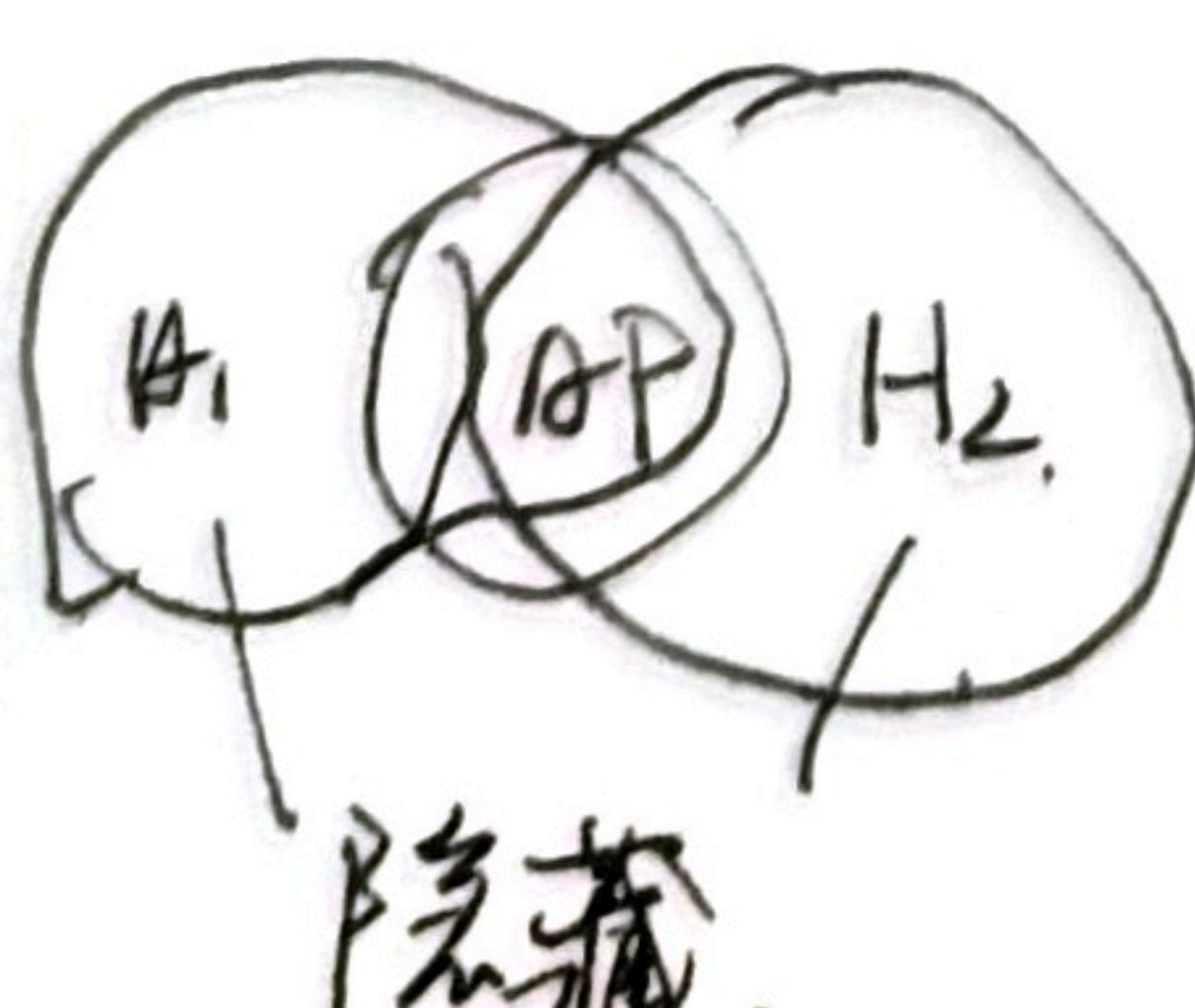
非带宽限制计数值, 则无碰撞.

Request to Send, Clear to Send 机制下,

RTS, CTS 都是广播发送

可以设置门限, 只将 RTS/CTS 用于长距离
数据的发送.

碰撞检测



d) IEEE 802.11 协议.

① 具有4个 MAC 地址字段。(自组织时才使用第4个)

地址1是接收帧该帧的源站点 MAC, CFP 或站点

地址2是目的 MAC 地址。(同一个 MAC LAN 内, 一般不使用虚地址)

地址3是源虚地址 MAC 地址, 可以是固定的也可以是动态的

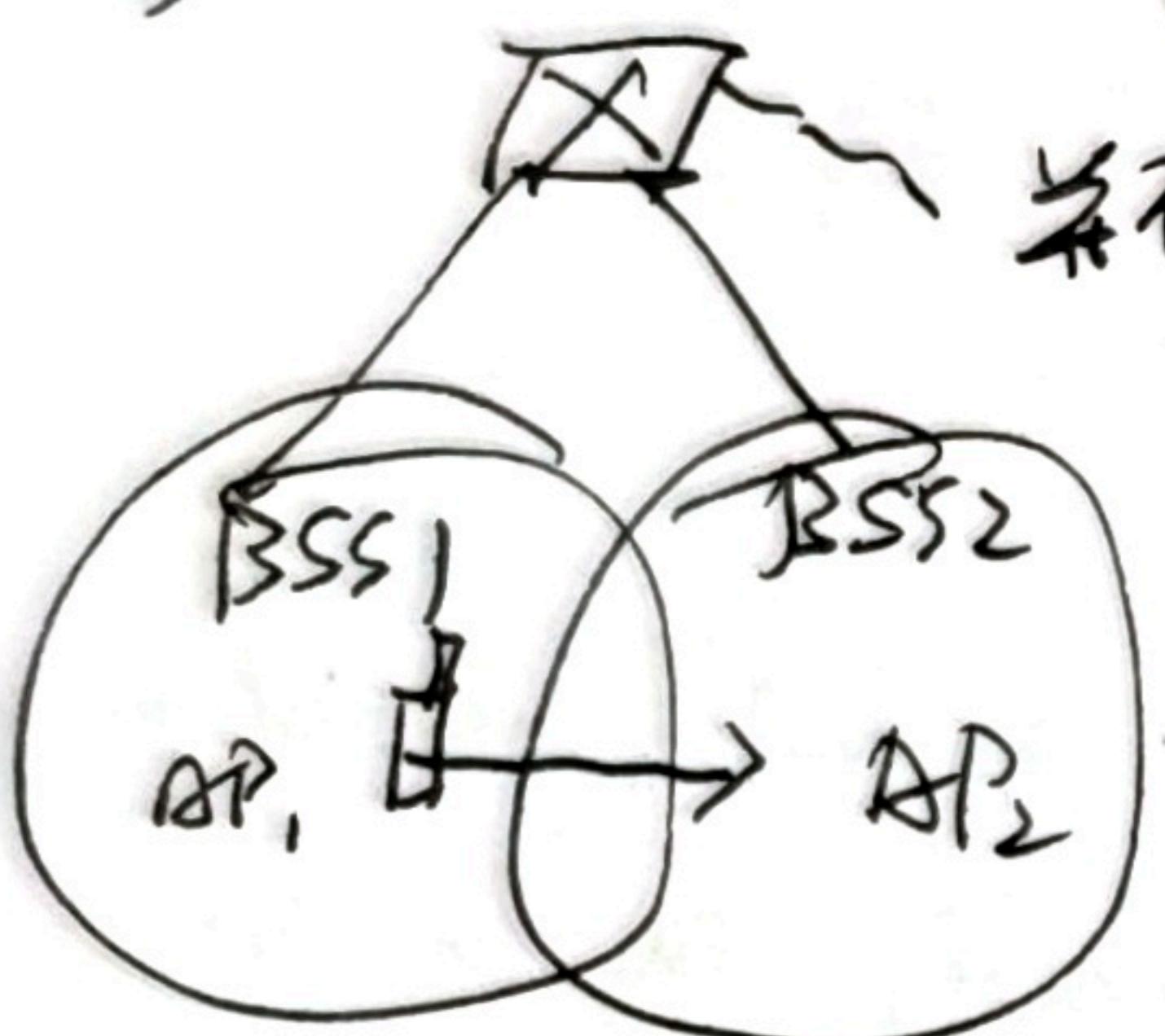
不仅可以接以太网中转和 WiFi 中转并转换。

② 具有信道子段, 及多帧传输的帧和虚连接。

③ 挂起期子段(重新时间和确认帧响应时间)

④ 帧控制子段(类型和子类型用于区分子关联, RTS, CTS, ACK 和数据帧)

⑤ 在相同子网中的移动性。



并不直接由累, 所以 BSS1 和 BSS2 属于同一子网, IP 和 TCP 不会变换。

信号衰弱, 收到其它信号帧 → 直接建立新关联(先解除)

对于切换, 接口改变, 利用直连, 让主机换 BSS2 之后换 AP2
而直接换 AP 会建立新的 MAC, 需换帧头重新交换。

f) 特色。

速率匹配, 自适应调制技术。

功率管理, 睡眠和唤醒。