

(2) Third Dup Ack:

Retransmit ² lost packet.

Assume (for time being) all ~~at~~
later packets made it
(continued "from end").

"Fast Recovery".

(decrease congestion window).

(3) SACK.

Retransmit lost packet. (not?).

Further? (problem!).

Decrease congestion window.

TCP is reliable.

Lost, damaged packets are re-transmitted.

How do we find out (detect, guess)
a packet is lost or damaged?

- (1) Time-Out.
 - (2) Third Duplicate Acknowledgment.
 - (3) (later) SACK.
-

Current situation (more or less, this is
implementation dependent):

(1) Time-Out:

Retransmit lost packet,
and all later packets.

(decrease congestion window)

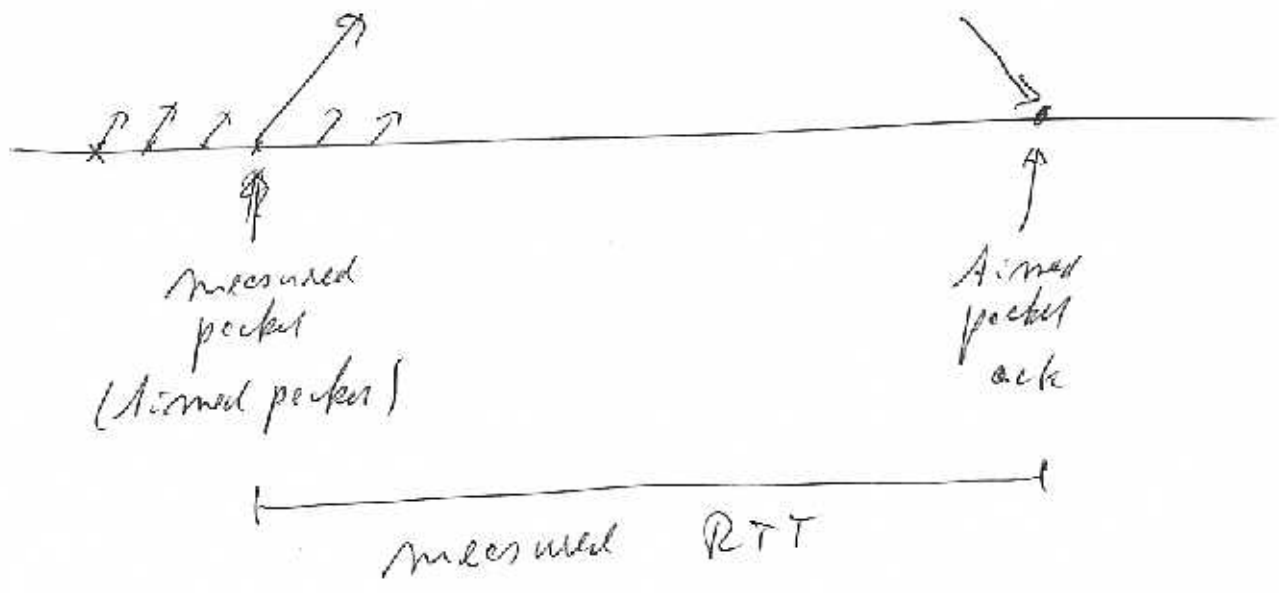
(Overload Control)

Round Trip Time Estimation.

Certain (not necessarily all) packets are

Aimed : RTT measured.

- (1) measure sending time
- (2) measure time of ACK.
- (3) difference is measured RTT.



RTT_{meas.}

In BSD Unix:



Only one outstanding Aimed packet at a time.

If "Aimed ack" : next detected becomes Aimed packet.

RTT_{meas}

measured RTT

RTT_{est}

estimated RTT.

RTT_{SD}

(estimated) standard deviation of RTT.

Algorithm:

When "Aimed ack" returns:

$$\Delta = |RTT_{meas} - RTT_{est}|$$

$$RTT_{est} = \frac{7}{8} RTT_{est} + \frac{1}{8} RTT_{meas}$$

$$RTT_{SD} = \frac{3}{4} RTT_{SD} + \frac{1}{4} \Delta$$

Re-Transmission Timer
 (Re-Transmission Time-Out)
 RTO

Standard:

$$RTO = 2 * RTT_{est}$$

BSD Unix:

$$RTO = RTT_{est} + \frac{1}{2} * RTT_{std}$$

(new values)

One way or another:

"we have an RTO".

Next: how do we use this?

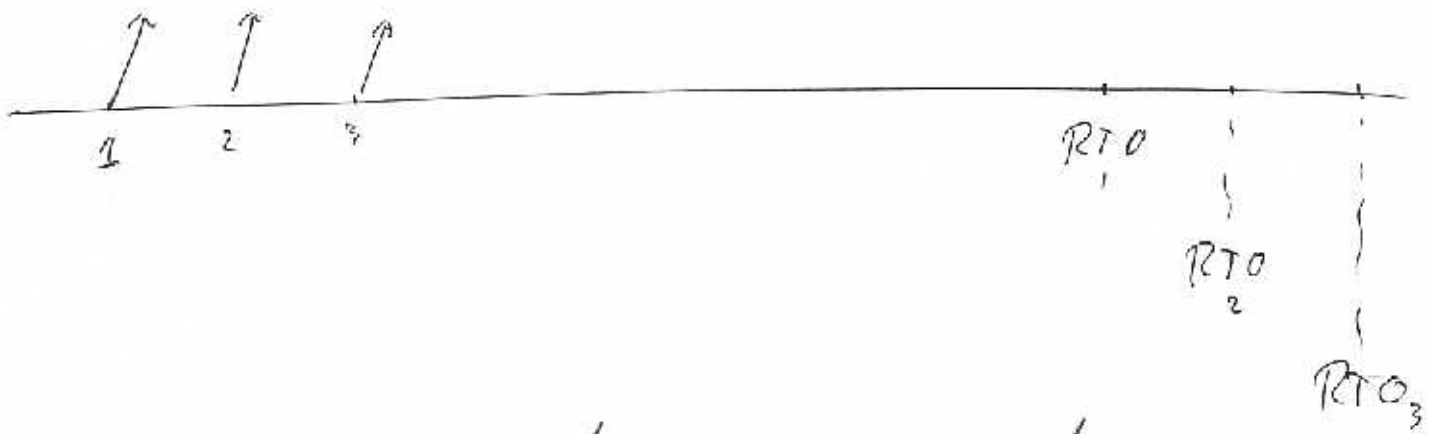
Slam here

03/30/2004

163

Naive:

~~Back to ...~~
RTO, and setting timer.



A "time-out" for every packet.

Heavy on the book-keeping!



BSD Unix:

If a time t a "good ack" arrives:

(1) kill RTO timer.

(2) Any more unacked bytes?

No: No new RTO timer.

Yes: set new RTO timer at $t + RTO$



This increases effective RTO interval
even more!

If we reach RTO timer:



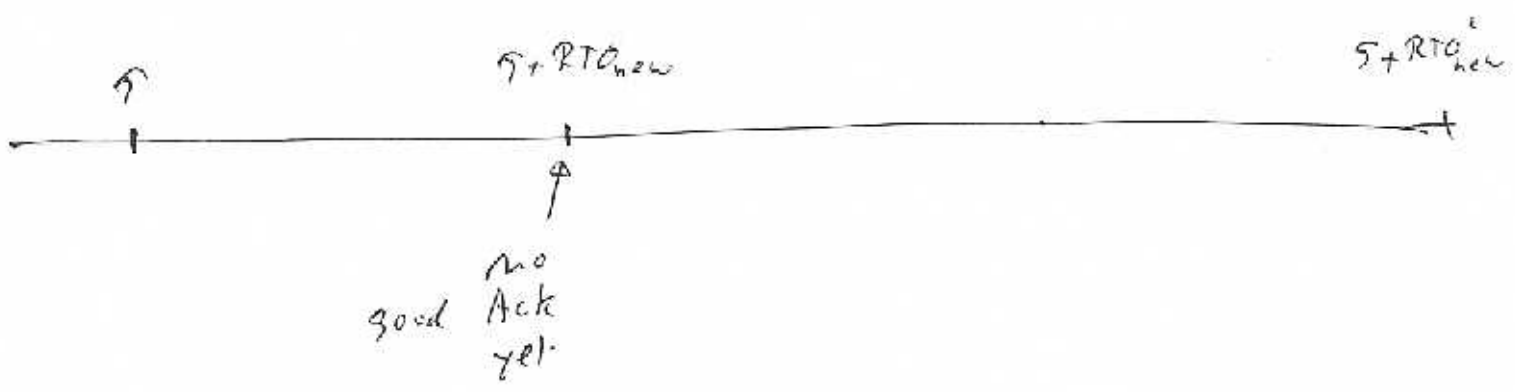
(1) Re-transmit "apparently missing" packet -

~~(Most implementations)~~

(2) Double RTO

$$RTO_{new} = 2 * RTO_{old}$$

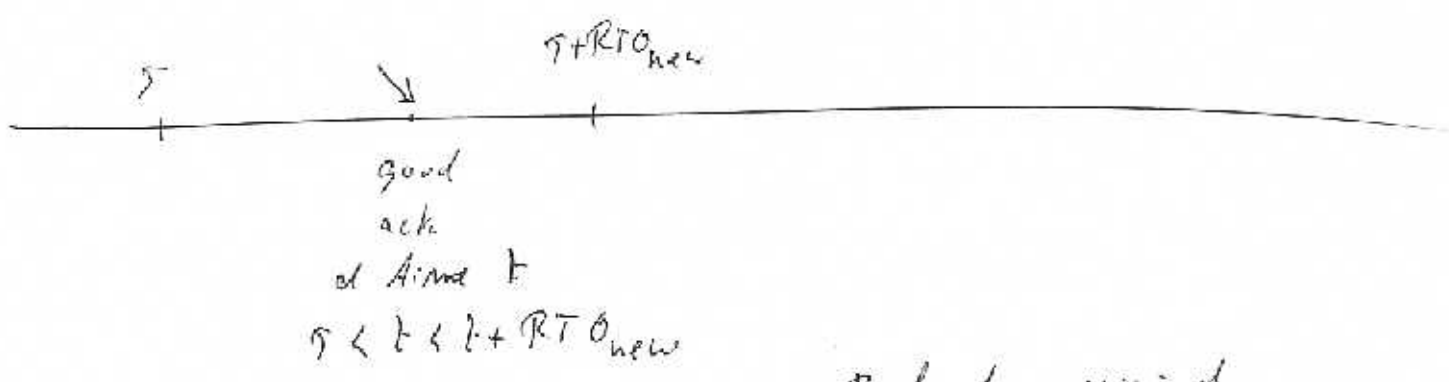
Next:



$RTO'_{new} = 2 \times RTO_{new}$
 (up to ~ 60 sec.)

etc.
 keep doubling

OR : hopefully :



Compute $RTO^* = ??$

Back to original
 from RTT_{est} etc.

More general:

$$RTT_{est, new} = \alpha * RTT_{est, old} + (1-\alpha) RTT_{meas.}$$

BSD Unix:

$$\alpha = \frac{7}{8}$$

If we time more packets per RTT:
 make α "somewhat" larger.
 ($(1-\alpha)$ "somewhat" smaller).

"Exponential Smoothing"

* Time series $(X_t)_t$ (say stock market)

$$\hat{X}_t = (1-\alpha) X_t + \alpha \hat{X}_{t-1} = (1-\alpha) X_t + \alpha(1-\alpha) X_{t-1} + \alpha^2 \hat{X}_{t-2} = \dots$$

$$\text{(previous estimated mean)} = (1-\alpha) \sum_{k=0}^{\infty} \alpha^k X_{t-k}$$

Mathematically: $|\alpha| < 1$

Common Sense: $0 < \alpha < 1$

Exponential smoothing determines "length of memory".
 α close to zero: short memory.
 α close to one: long memory.

Sharma's Algorithm

Simplification:

Stop estimating RTT during recovery.

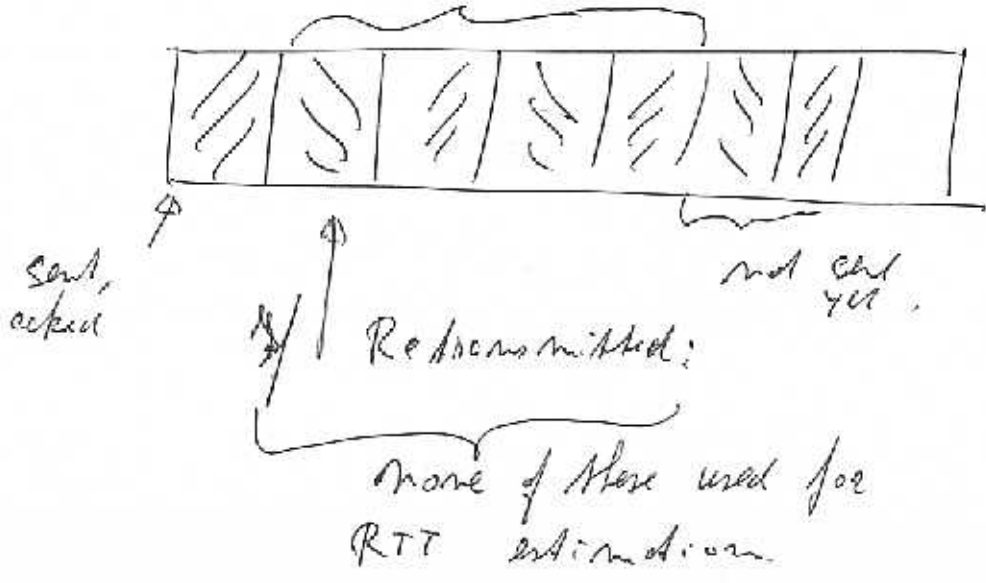
(Do not use data sent above for or more for RTT est update)

Complication: protect against sudden increase in β actual RTT.

Corner is risky - nearby.

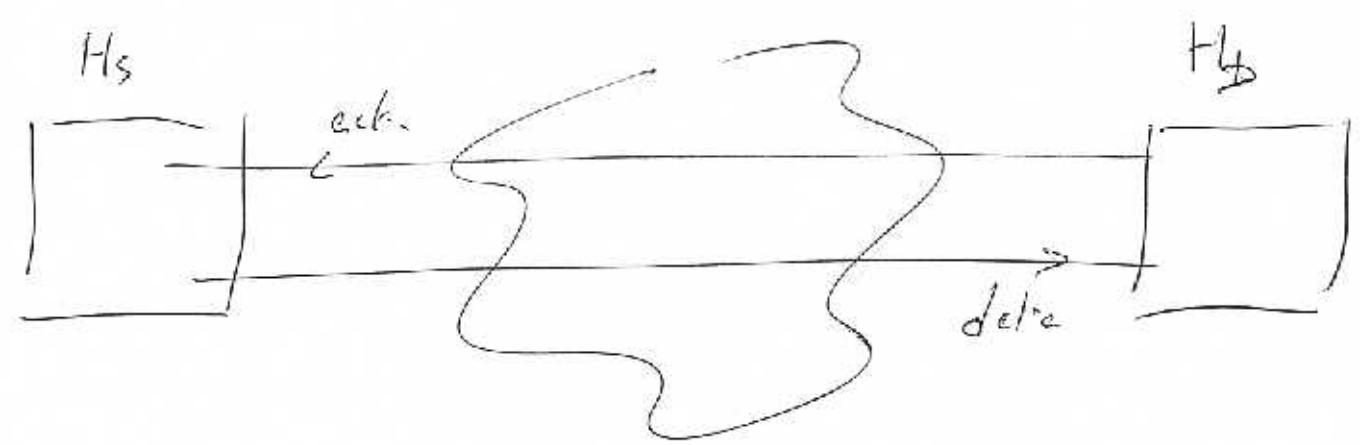
Send Buffer:

sent, not acked yet.

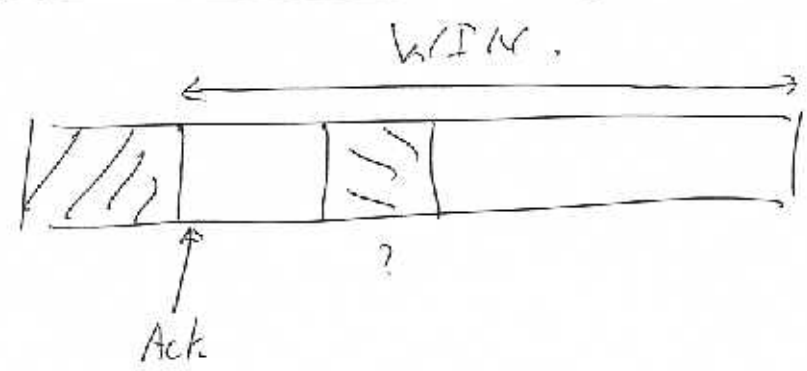


Plus Bells and Whistles!

Effective window.

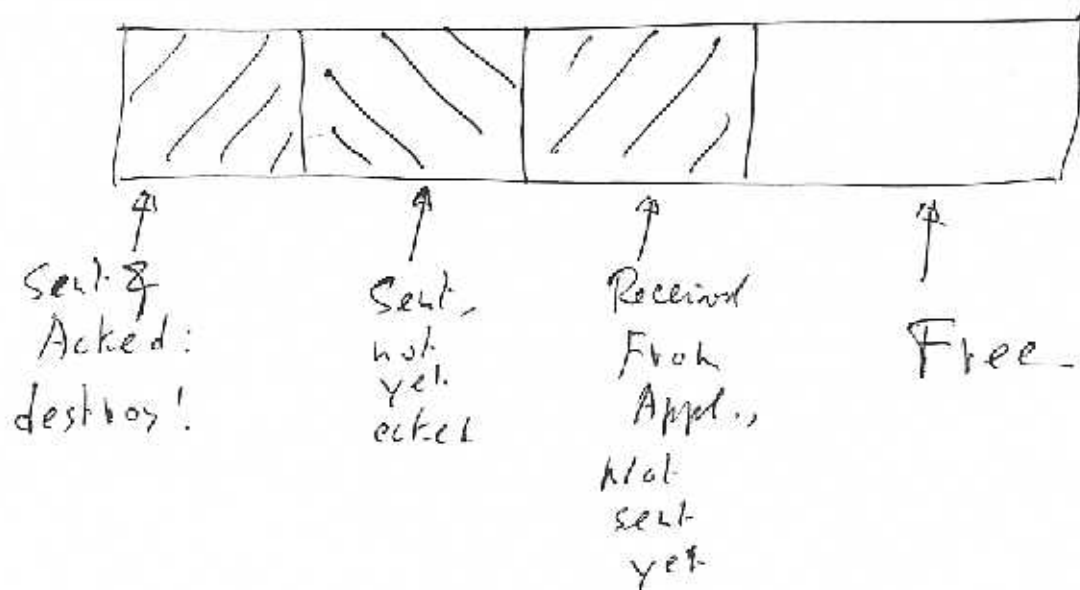


Defn : Receive Buffer



WIN = Advertised Window.

Source Buffer (Send Buffer)



Amount of stuff sent, not yet acked: \leq Send Buffer.

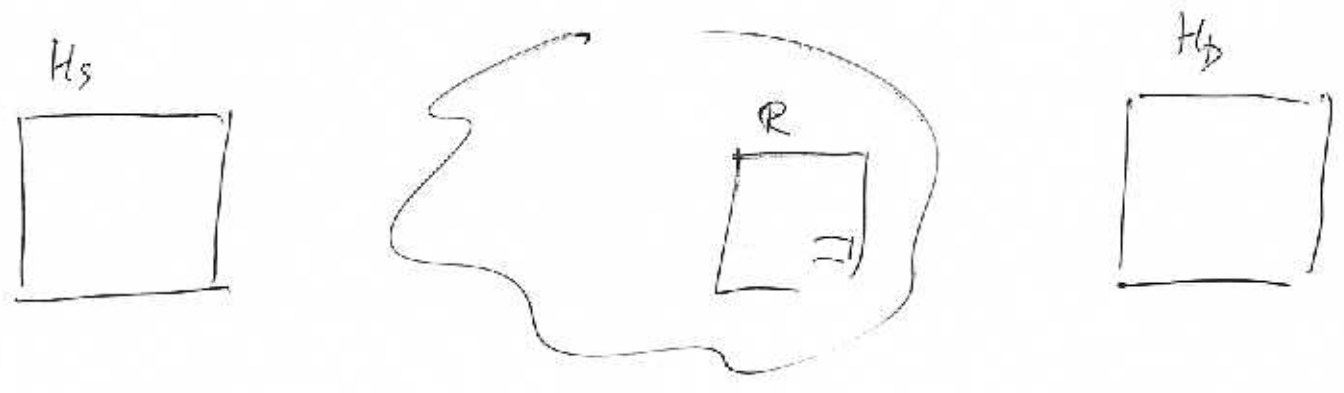
Combine:

Amount of stuff sent, not yet acked

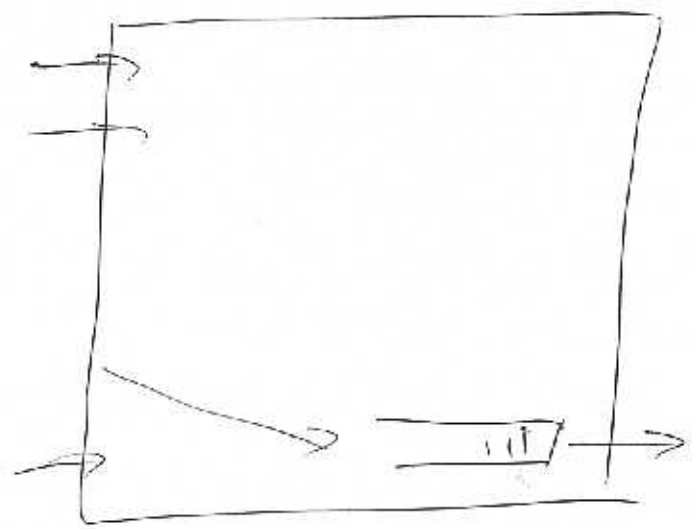
\leq \min (Advertised WIND, Send Buffer).

\leq \min (Send Buffer, Recv Buffer, $2^k \times \text{buffer}$).

Also : Congestion Window.



Router :



Too many packets to some output buffer in short time : Congestion!

Buffer overflow!

Those flows must decrease!

$$cwnd(h) = \begin{cases} cwnd(h) + \frac{1}{cwnd(h)} & \text{if Ok.} \\ \frac{cwnd(h)}{2} & \text{if not Ok.} \end{cases}$$

171 -

BSD Unix: $cwnd$ is expressed in Bytes,
not MSS's.

$$\frac{cwnd}{MSS} = \frac{cwnd}{MSS} + \frac{1}{\left(\frac{cwnd}{MSS}\right)} \quad \text{not MSS}$$

$$cwnd = cwnd + \frac{(MSS)^2}{cwnd}$$

Funny Story: See Stevens II p 977
(Since changed!)

Effective window size:

$\min(\text{cwnd}, \text{send_window}, \text{advertised window})$.

receiveWind, $2^k \times 65535$

(make sure to get units right).

If the WIN_{eff} is "constant" and the source has plenty of stuff to send:

$$\text{Thput} = \frac{WIN_{\text{eff}}}{RTT}.$$

1 T Byte sec in 4 hours.

10^{12} Bytes in $4 \times 60 \times 60$ sec:

$69,444,444.44$ Bytes/sec. = $555,555,555.55$ b./sec.

in 50 msec = $\frac{1}{20}$ sec.

$3,472,222.22$ bytes in 50 msec.

Receive Buffer and Send Buffer must be at
at least 3.5 M Bytes!

What kind of window scaling do we need? 133

$$\frac{3.5 \text{ MBytes}}{65535 \text{ Bytes}} \sim 52.98$$

$32 < 53 < 64$

We need $k=6$.

Choose $k=6$.

$$\frac{6}{2} =$$

$$2^6 * 65535 = 4,194,240 \text{ Bytes.}$$

Choose $k=6$,

Receive Buffer,

Send Buffer.

Both at least 4.2 M Bytes.

~~Don't forget~~ also modify TCP behavior!

Start here 04/02/04.

Silly window syndrome.

(e.g. Telnet).

"Many Small packets"

Can be due to slow source appl
(Telnet)

Slow dest. Appl (~~at~~ adv. window)

Solution:

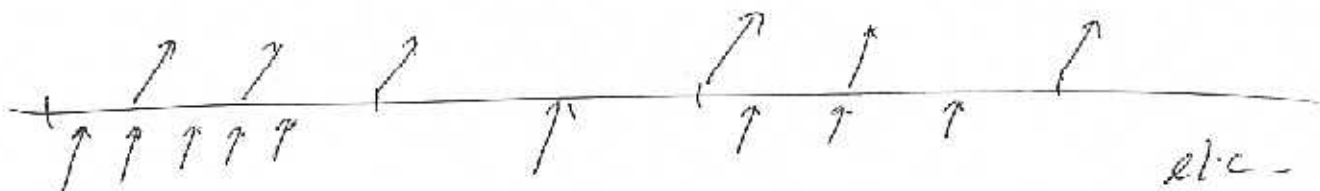
Two seemingly unrelated entities.

(1) At Source Side:
Nagle's algorithm.

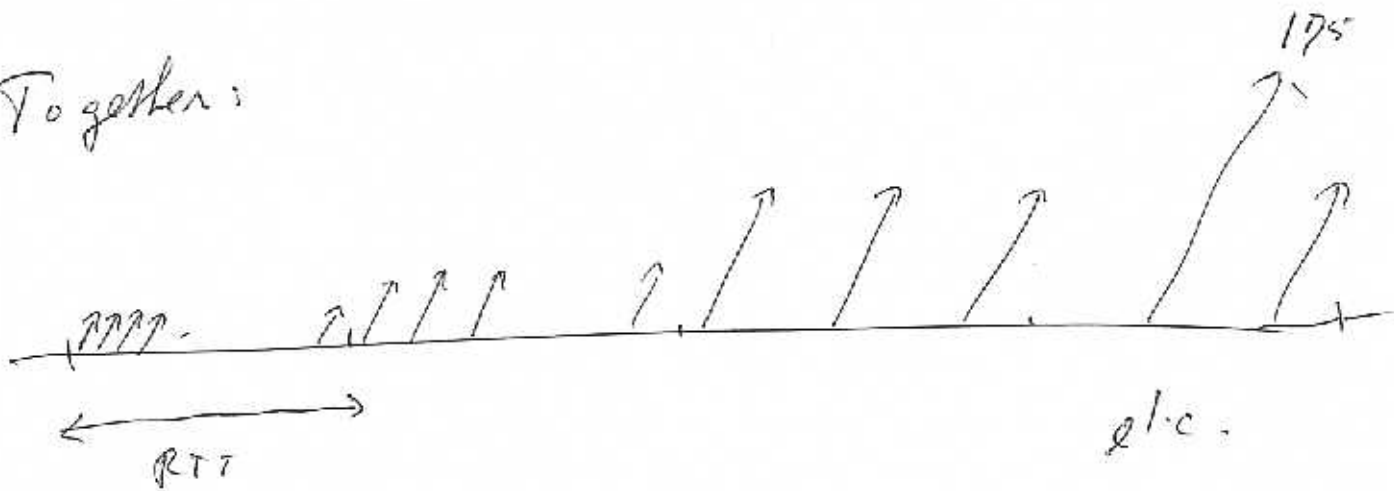
Start: just send.

Later: only if (1) whole 1955,
or (2) Good Ack.

(2) Dest Side. Delayed Acknowledgments.



Together:



Soon:
either 1 data packet per clock tick or other
or whole mess.

Under push



Routing.

Cormer p 253 etc.

Routers exchange information about the "state of the network", so each can construct a forwarding table.

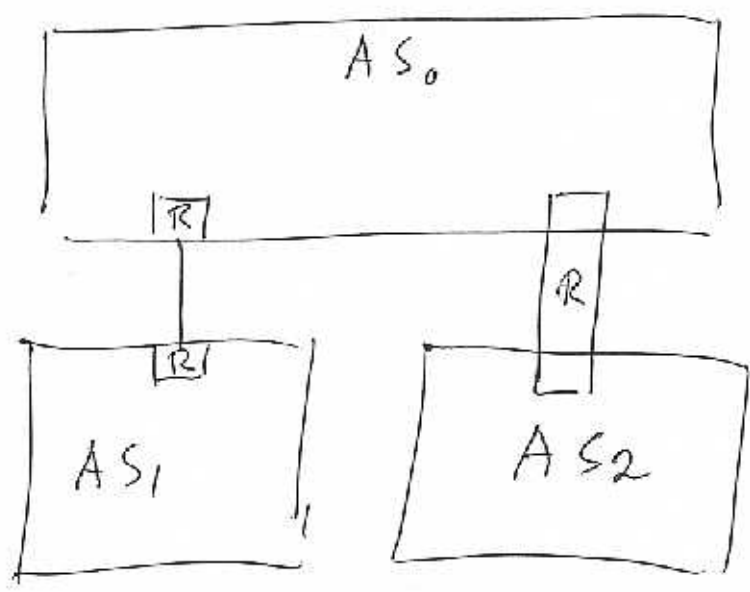
Question: do routers trust each other?

Autonomous Systems ASs

"Administratively Autonomous".

(Autonomous domains
Administrative domains).

Within AS: Routers "trust" each other.



Routing "between" ASs:

Exterior Gateway Protocol.

EGP.

Routing "inside" ASs:

Interior Gateway Protocol.

IGP.

Examples:

EGP: only one: BGP. Ch. 15
 Border Gateway Protocol.
 (IETF).

IGP:

RIP Routing Information Protocol.
 v1 and v2. Ch. 16

OSPF Open (Shortest Path First) Ch. 16

IS-IS Intermediate System -
 Intermediate System.

Hello.

RIP: IETF

IS-IS: originated in ISO
 (International Standards Org).

OSPF: IETF, grew out of IS-IS

BGP: IETF

BGP allows "policy routing".

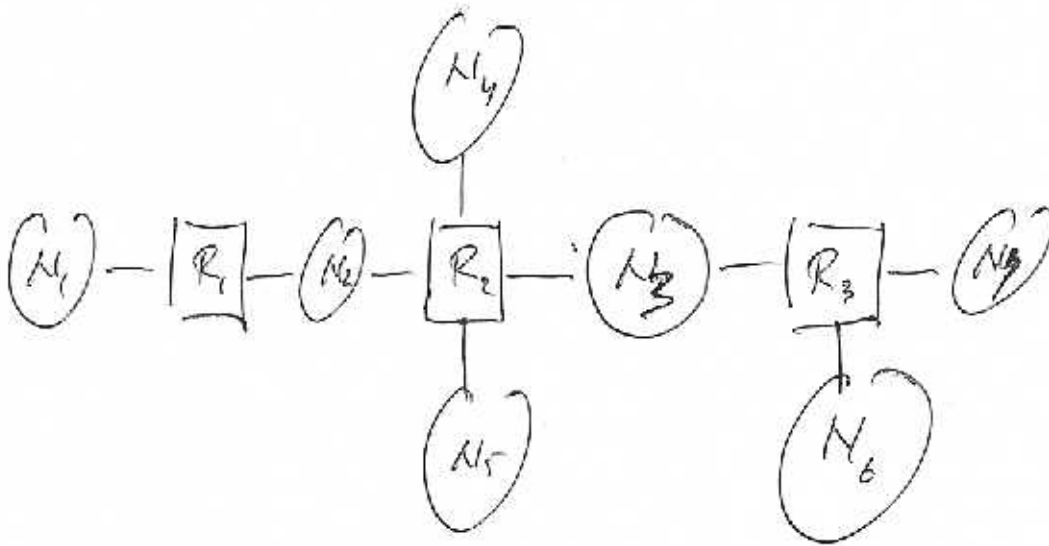
(e.g. traffic of customer A is not allowed on the network of long haul company Z).

Distance Vector Protocols

Basic idea: every router has a routing table and sends advertisements to its neighbors:

"I know how to reach network N, and this is how good the route is. If you send me stuff for N, I will forward it."

Start here 04/06/2004



Suppose R_1, R_2, R_3 boot at the same time. (Unrealistic!)

R_1

R_2

R_3

R_1

netN	dist	next
N_1	0	DD
N_2	0	DD

N_2	0	DD
N_3	0	DD
N_4	0	DD
N_5	0	DD

N_3	0	DD
N_6	0	DD
N_7	0	DD

advertisers :

RIP v1
RIP v2

advertises only network & distance (hopcount).
also advertises "next hop" and metric.
(OSPF)

N_1	0
N_2	0

N_2	0
N_3	0
N_4	0
N_5	0



Let's assume R_1 advertisements from R_2 .

R_2

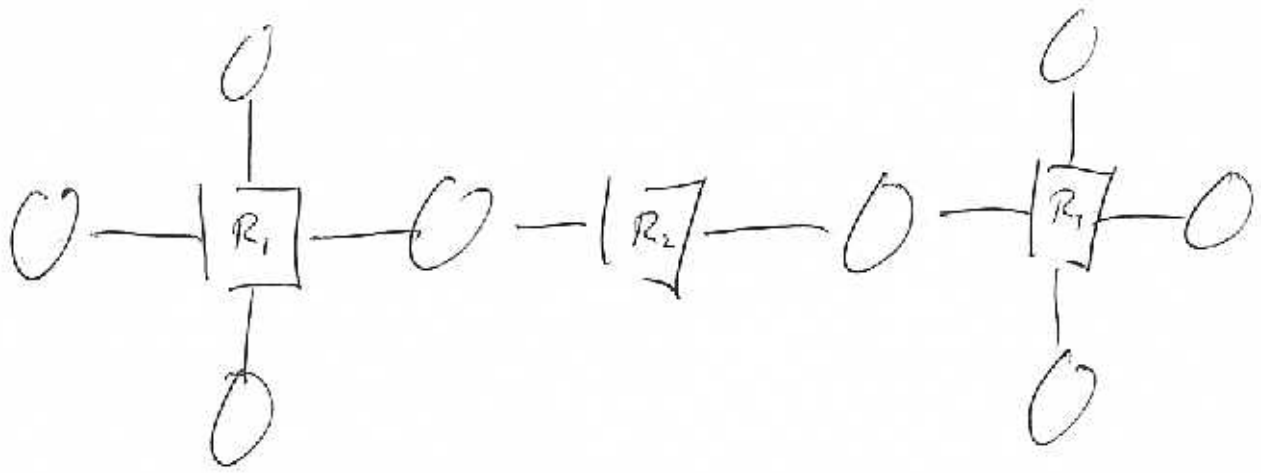
N_1	1	R_1
N_2	0	$\$$
N_3	0	$\$$
N_4	0	$\$$
N_5	0	$\$$

Let's assume next R_2 advertisements to R_3

N_1	1
N_2	0
N_3	0
N_4	0
N_5	0

R_3 :

N_1	2	R_2
N_2	1	R_2
N_3	0	$\$$
N_4	1	R_2
N_5	1	R_2
N_6	0	$\$$
N_7	0	$\$$



Router R_x has currently has hopcount $d_{x,N}$ to net work N .

Gets advertisement from router R_y : $d_{y,N}$.

- (1) If $d_{x,N} > d_{y,N} + 1$:
change route $d_{x,N} = d_{y,N} + 1$,
next hop Y .
- (2) If $d_{x,N} < d_{y,N} + 1$ and next hop is Y :
 $d_{x,N} = d_{y,N} + 1$, next hop is Y .
- (3) Otherwise: no change.

Problem with RIP:

"slow convergence"

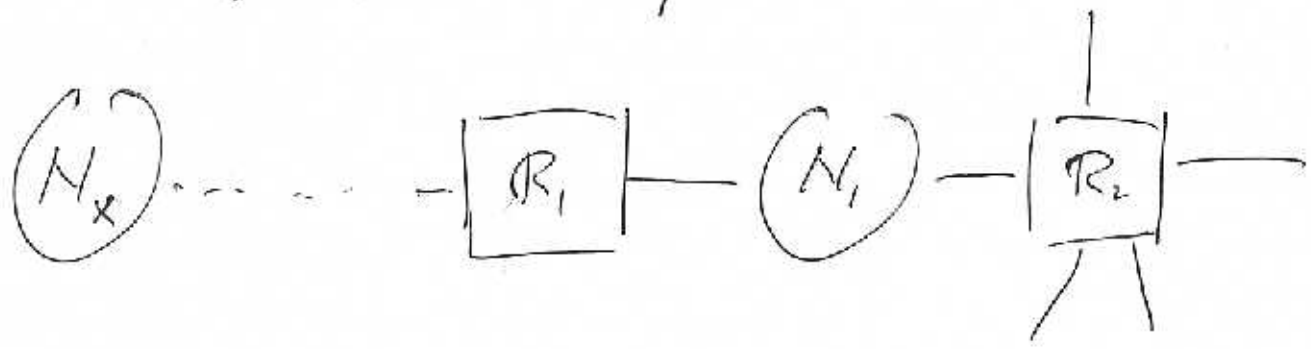
advertise one every 15 sec.

distance Φ : takes $\Phi * 15$ sec

Problem with RIP:

"Walk to infinity."

"Count to infinity".



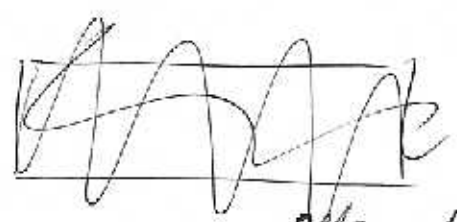
$$d_{R_1, N_x}$$

$$d_{R_2, N_x} = d_{R_1, N_x} + 1$$

Now the connection between R_1 and N_x breaks.

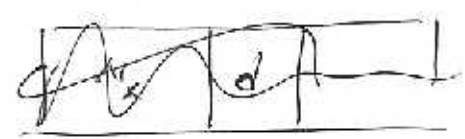
And, unfortunately, R_2 advertises to R_1 right after.

R_1



gets from R_2

N_x	"infinity"	none
-------	------------	------



R_x	d_{R_2, N_x}	
-------	----------------	--

R_1

N_x	d_{R_2, N_x}^{+1}	R_2
-------	---------------------	-------

advertises to R_2 :

R_2 :

N_x	d_{R_2, N_x}^{+2}	R_1
-------	---------------------	-------

etc.

"walk to infinity".

"Solution: make infinity small.
infinity = 16.

N	16	?
-----	----	---

"I know about N , but you can not get there from here".

Solutions:

Split Horizon.

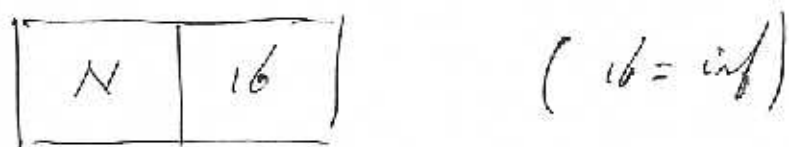


then R_x does not advertise network N
to router R_y

Poison Reverse. (Poisoned Reverse?)



then R_x advertises to R_y :



Triggered updates:

~~Send~~ Immediately send advertisements when receiving "Bad news".

Hold-down: ~~The opposite of~~ After "bad news", do not use "good news" for a while.

RIP uses UDP.

port. 520.

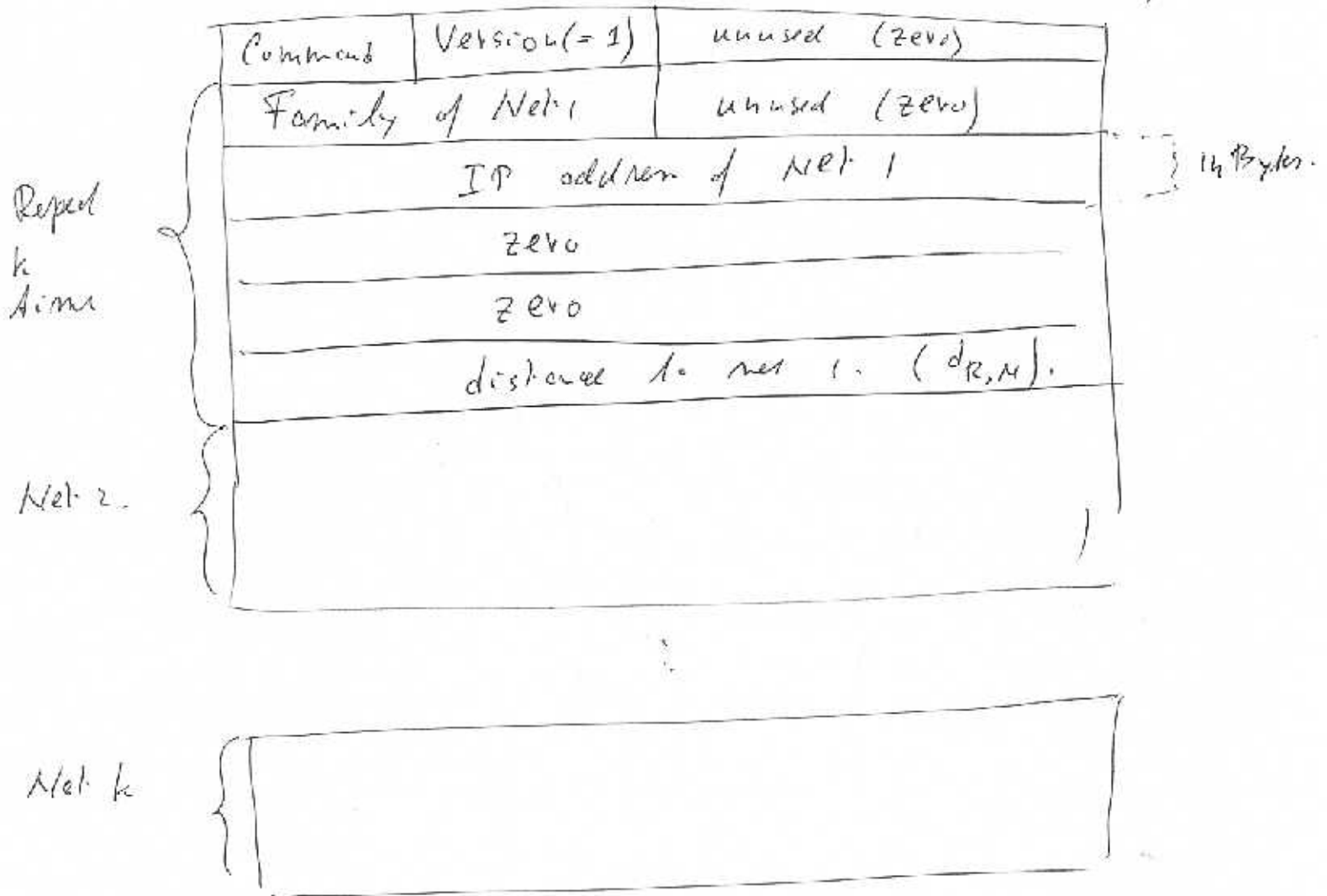
Is ok: neighbors only..

Start here 04/13/2004

RIP 1 format.

Cornes p. 301

once every 30 sec.
(or one every 15 sec.?)



Comments: See Cornes p 301

Version: 1 for RIP 1 (!)

Family: 16 bits = 2 for TCP/IP

Address: 14 Bytes (!)

{ 1 Request
 2 Response
 ..
 "

Route Tag

Used to identify "owner identifier" of route.

E.g. "Autonomous System Number" from BGP. (both are 16 bits).

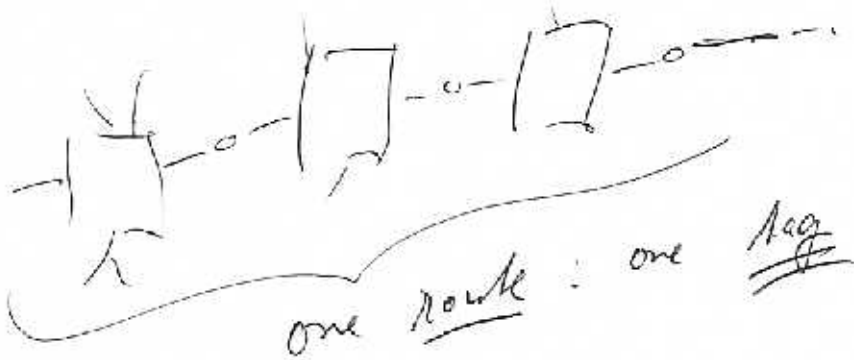
Also: Used to identify "Route"

RIP: no message length (# net works).

RIP uses UDP port 520.

UDP message length!

UDP: Ok. neighbors only.



Hello Protocol:

There are two of them!

① Hello protocol and securities.
"whole protocol".

② Hello Protocol, part of OSPF.
"know thy neighbor".

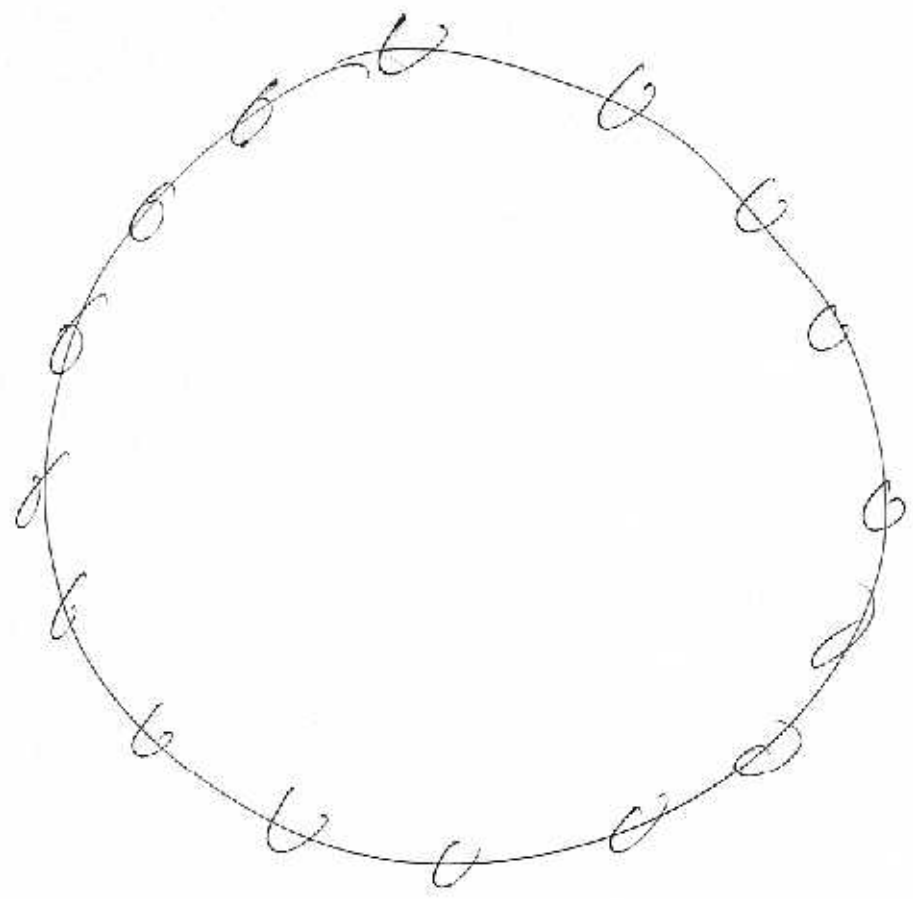
Hello Protocol.

Based on Delay Estimation.

In stead of a hop count :
estimated delay.

Find the "minimal delay" route.

Caused Oscillation.



Oscillatory Behavior

Can be countermeasured by changing the delay estimate "very slowly".

Still, a problem with
"all or nothing" routing.

OSPF.

OSPF.

192 ~~198~~

(1) This is a link state protocol.

Every router sends advertisements describing the local network.

Includes Type of Service routing.

Delay

~~Throughput~~ Throughput

Reliability.

Each router develops a map of the whole network (whole area) and in that area finds the "best route" to all networks.

• These paths had better be consistent.

Cost of a link.

Try to find the route that

minimizes the cost.

$$\text{minimize } \sum_{l \in \text{Path}} c_l$$

Just like in RIP:

a mathematical optimization.

How are the costs used?

To distribute load

Cornet p 309

ToS Routing

Load Balancing: Equal Cost Multipath

Areas

Authentication!

All kinds of routers.

→ Link State

Virtual network ← Transistors

External link, External link

OSPF has an option

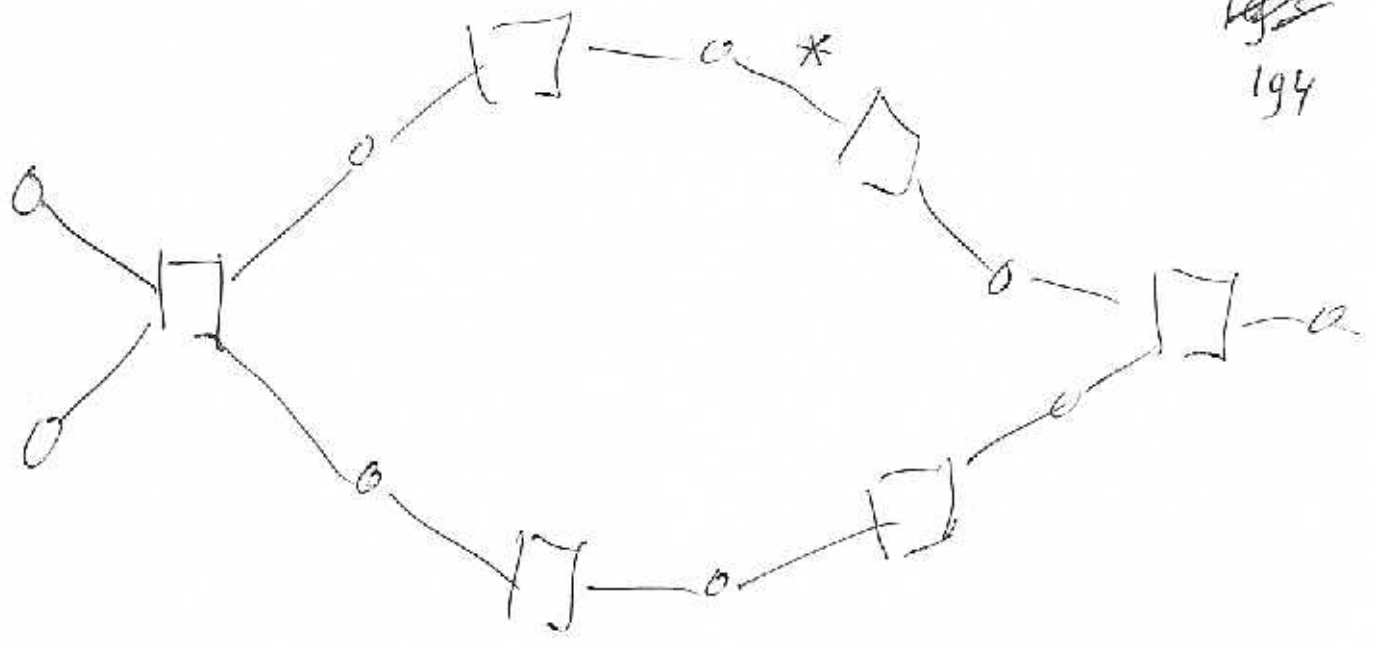
"~~Multi~~" "Equal Cost Multipath"

193.

B

OSPF: RFC 2328

OSPF equal cost multipath
RFC 2138. (?)



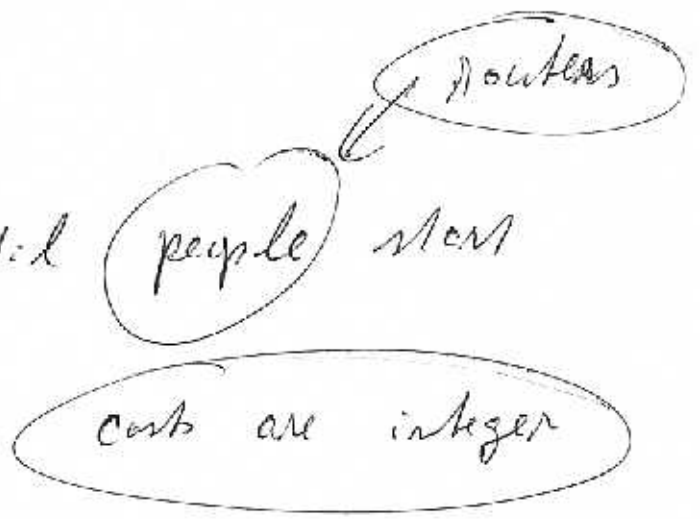
6. (Actually, many of these links probably are point to point)

* is "too busy".

increase its cost, until people start re-routing.

What happens?

"Bang".



OSPF: RFC 2328

OSPF { Multipath 2178

Equal Cost

(Check?)
(Wrong order).

RFC numbers need checking!

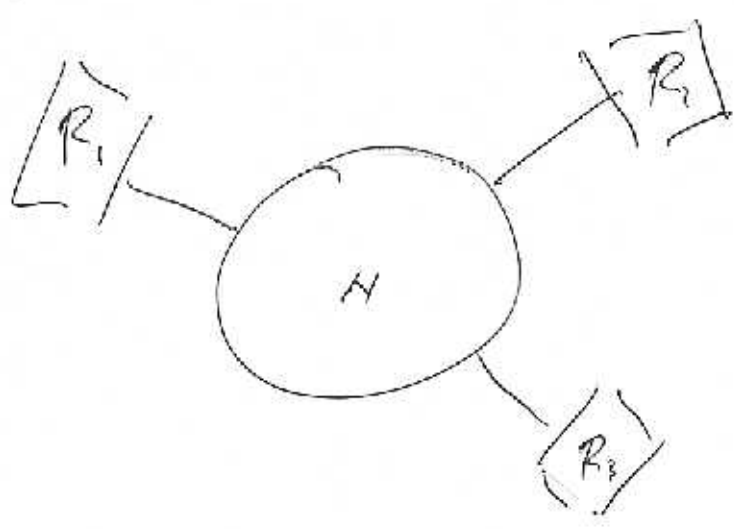
OSPF.

Types of "Links"

(1) Point to Point.



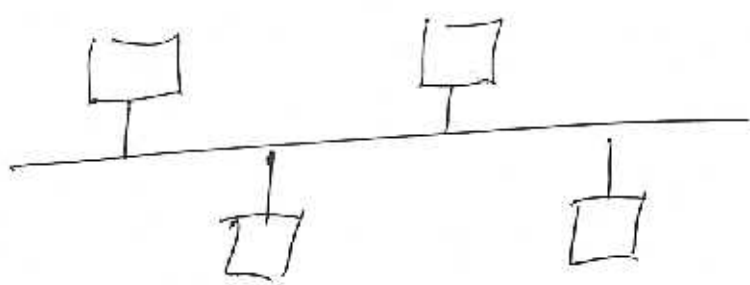
(2)



Transmits Link
Transmit Network.

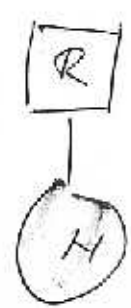
(3)

or :



(3)

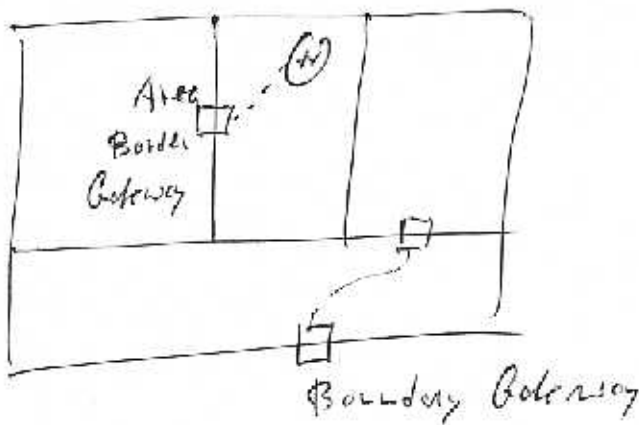
Stub



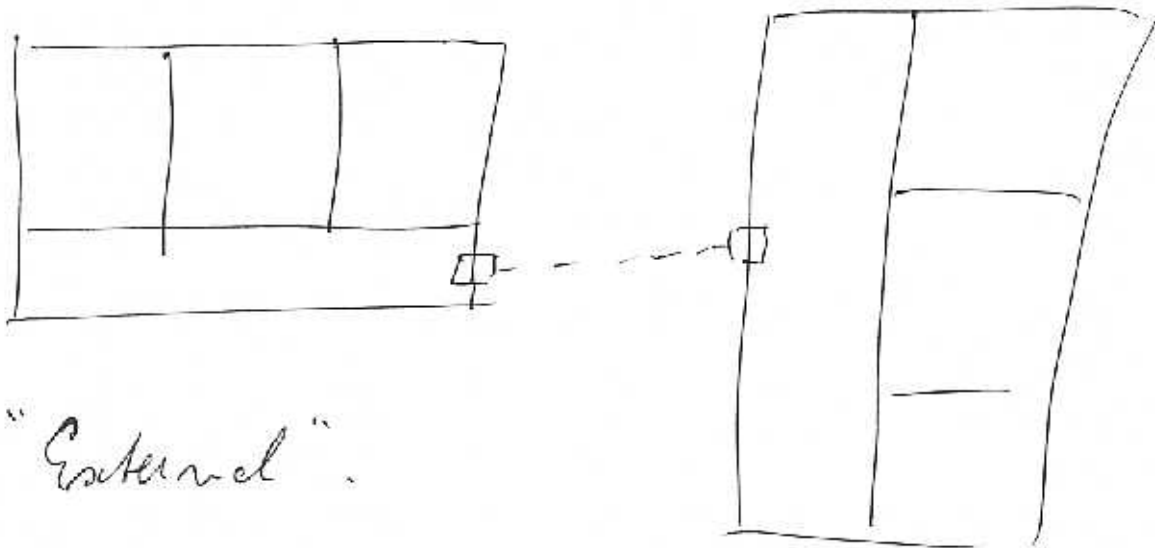
one router
on network.

PPP, Transient, Stub: are real.
(physical).

There are more that are "virtual".

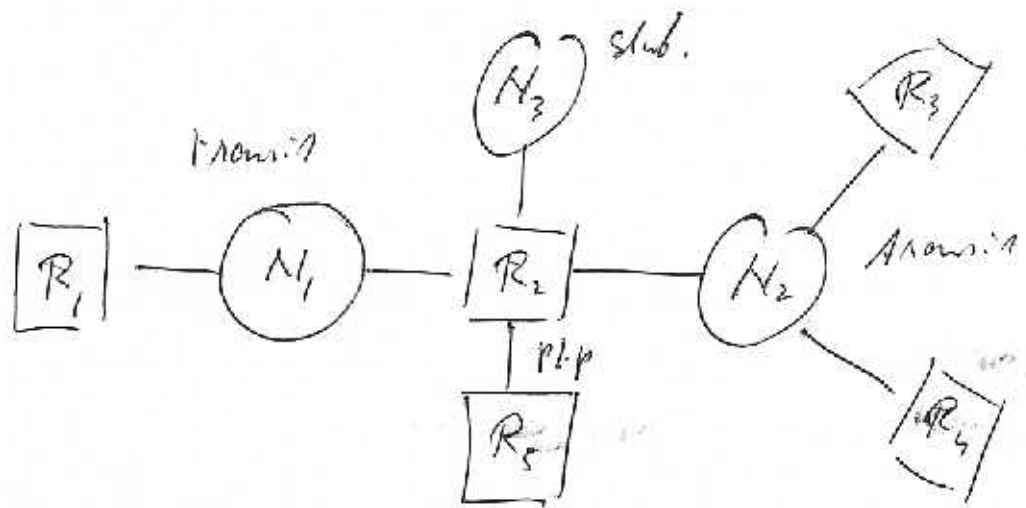


Summary Link



"External".

Such links are advertised
"kind of like a stub".



RFC 2328

From.

	R_1	R_2	R_3	R_4	R_5	N_1	N_2	N_3
R_1						0		
R_2					3	0	0	0 !
R_3							0	
R_4							0	
R_5	0	2						
N_1	3	2						
N_2		3	4	11				
N_3		7						

Blank is "inf".
not zero.

Each column is an advertisement.

R_1/c R_2 advertises

R_5	2
N_1	2
N_2	3
N_3	7

Who advertises for N_2 ?

one from R_2, R_3, R_4 .

The "designated router" for N_2 .

For example R_2 in that case:

R_2 advertises

R_5	2
N_1	2
N_2	3
N_3	7

for itself

and it advertises

R_2	0
R_3	0
R_4	0

for N_2 .

Recursion:

Every Router gets picture of whole network.
(Hopefully: Same!)

Peer-to-peer:

Distance from to (next) router
is always zero!

Distances asymmetric.

Who advertises for N_3 ?

Nobody.

Routing with partial information