

Internet Engineering

241-461

Robert Elz

kre@munnari.OZ.AU

kre@coe.psu.ac.th

<http://fivedots.coe.psu.ac.th/~kre>

ICMP Error Rules

- ◊ Never send ICMP error packet if:
 - packet with problem is not initial fragment
 - if Offset field in header != 0
 - packet with problem might have been delivered to many destinations
 - IP destination addr is broadcast or multicast
 - Link layer destination addr was broadcast or multicast
 - packet with the problem was an ICMP error packet
 - What about unknown ICMP types? (like type 6)
 - Only safe action is to treat as error packet
 - Sending ICMP is never required
 - Sending ICMP is sometimes prohibited

Do not send to you send!

Contents

- ◊ Why Fragmentation?
 - ◊ Why Reassembly?
 - ◊ When to Fragment
 - ◊ When to Reassemble
 - ◊ How to Fragment
 - ◊ How to Reassemble
 - ◊ Why not fragment
 - ◊
 - ◊ How to avoid fragmentation (PMTUD)
 - Send packets small enough
 - What is Small enough?
 - Path MTU Discovery
- ICMP

Avoiding Fragmentation

- ◊ If host sends small packets fragmentation will not be needed
- ◊ If packet is lost host can retransmit just that packet
- ◊ Other packets that reach destination can be retained and used
 - Transport Protocol issues (TCP)
- ◊ How small is small enough?
 - 68 bytes certainly
 - **Packets 68 bytes or less cannot be fragmented ... Why??**
 - Consider 60 byte IP header (which is possible)
 - Need at least 1 data byte (or infinite number of frags required)
 - Leading fragments must have multiple of 8 data bytes
- ◊ Network must be able to send 68 byte packets
 - on every link used for IP
- ◊ Using 68 byte packets
 - 20 byte IP header (minimum)
 - 20 byte TCP header (minimum)
 - 40 bytes of headers
 - 28 bytes remain for data
 - Almost 60% overhead - best case

Smaller multiple of 8.

Avoiding Fragmentation (2)

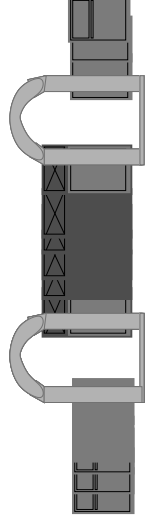
- ◊ 68 byte packets avoid fragmentation
 - 100% certainly
 - But too much overhead
- ◊ Can we send bigger packets
 - And still avoid fragmentation?
- ◊ Perhaps
 - It depends upon the network
- ◊ Usually assume 576 byte packets are OK
 - Certainly not guaranteed
 - Some slow links fragment around 200 bytes
 - so packets transmit quickly
- ◊ Almost always OK
 - Because of common mistake
 - Often believed that 576 is minimum required MTU
 - 576 is minimum packet reassembly size (host buffer)
 - RFC says hosts must be able to receive IP packets 576 bytes big
 - RFC does not say network must be able to transmit them

Avoiding Fragmentation (3)

- ◊ Can often send 1500 byte packets
 - No fragmentation required
 - < 3% header overhead (40 header bytes)
- ◊ Rarely possible to send bigger
 - While still avoiding fragmentation
 - Usually at least 1 ethernet in the path
- ◊ Is OFTEN good enough ?
 - No - not really
 - Too many links still < 1500 bytes MTU
 - One example
 - Put packet inside another packet
 - Send outside packet across net
 - Remove outer packet
 - Known as a tunnel

starting with packet size & control using the network

Tunnelling



- Red (outer) packet \leq 1500 bytes
 - to avoid fragmentation
- Thus green (inner) packet must be
 - \leq 1500 bytes - red packet header overhead
 - Or red packet will be fragmented
- If our packet might be green packet
 - And we have no way to know for sure
 - We must send $<$ 1500 bytes
 - or fragmentation will probably happen
- How big can we send?
 - And how do we find out?

Contents

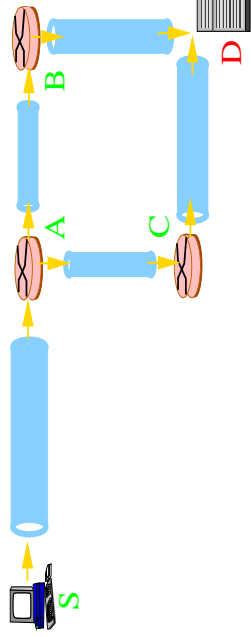
- ◇ Why Fragmentation?
- ◇ Why Reassembly?
- ◇ When to Fragment
- ◇ When to Reassemble
- ◇ How to Fragment
- ◇ How to Reassemble
- ◇ Why not fragment
- ◇ How to avoid fragmentation (PMTUD)
 - Send packets small enough
 - What is Small enough?
 - Path MTU Discovery

Path MTU Discovery

- ◇ We need a mechanism
 - To find the smallest MTU
 - On the current path
 - From sender (me) to the destination
- ◇ The smallest MTU on the path
 - The Minimum Maximum Transmission Unit
 - that is the choke point for fragmentation
 - send packets \leq this Minimum MTU in size
 - fragmentation never required on this path
 - can send packets this big
 - no need to send smaller
 - lower header overheads
- ◇ How do we find this important number ?
 - This is Path MTU Discovery PMTUD
 - Finds the MTU of the current path
 - from sender to recipient

Path MTU Discovery (2)

- ◊ Recall this example

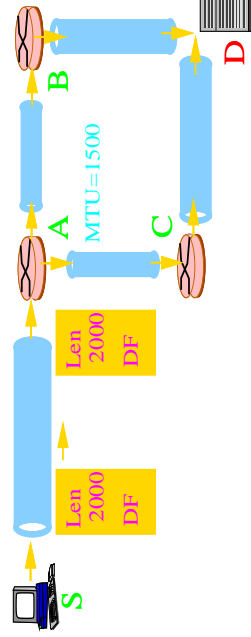


- We want to find the Path MTU (PMTU)
 - from sender S to destination D
- ◊ Start by assuming
 - PMTU is MTU of local link
 - certainly cannot be bigger than that
- ◊ Send biggest allowable packet
 - That will avoid fragmentation
- ◊ With the DF flag set in IP header

Path MTU Discovery (3)

- ◊ Sender sends biggest packet

- Local link will allow
- Sets DF in header



- Packet arrives at router
 - Where next link MTU
 - Smaller than packet size

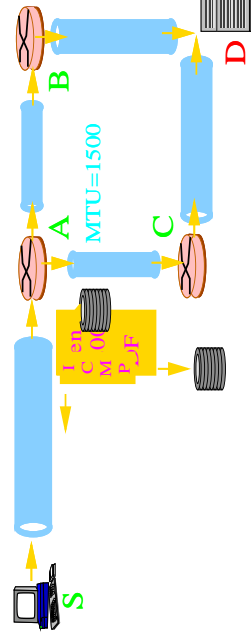
Path MTU Discovery (3)

- ◊ Packet bigger than MTU of link to C

- We assume path is S-A-C-D

- ◊ Router must discard packet

- Too big for link, and DF is set



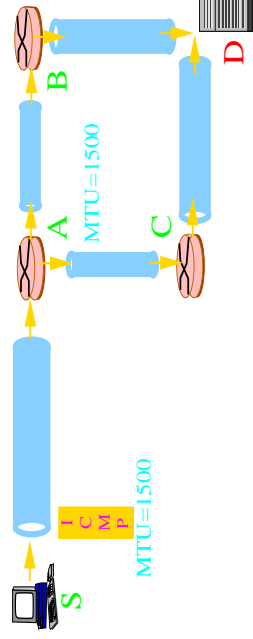
- ◊ Router also sends ICMP message to S

- Tells it that packet was discarded
- And why packet was discarded

Path MTU Discovery (4)

◊ S receives ICMP message

- Discovers that packet it sent was dropped
- Dropped because it was too big for link
 - And that 1500 is the available MTU

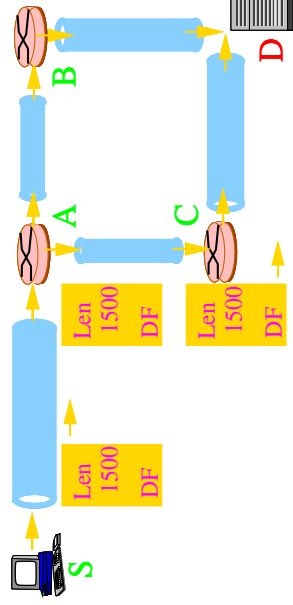


- S can now send a smaller packet
 - limited to 1500 bytes max
- And remember that 1500 is Path MTU to D
- This is Path MTU Discovery

More PMTUD

◊ S now knows 1500 as PMTU to D

- So sends packets that big

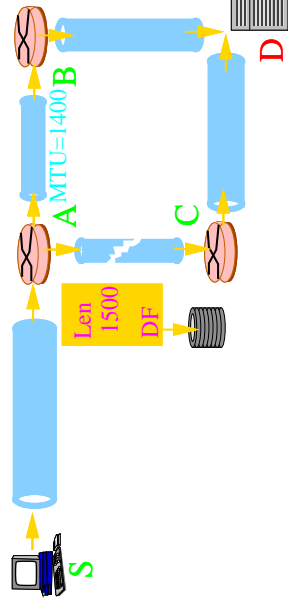


- These go via A and C, and arrive at D

More PMTUD (2)

◊ S keeps sending with DF flag set

- This is because...

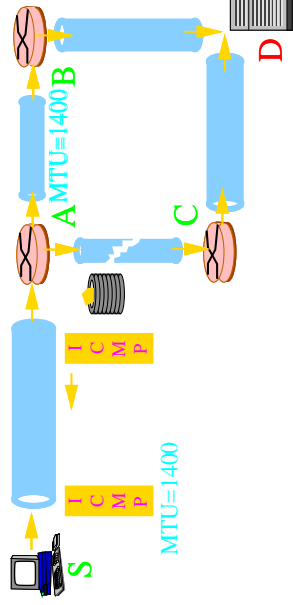


- The link from A to C might fail
- The alternate link (from A to B)
 - has an MTU of just 1400
- Because the DF flag remains set
 - A discards the packet

• 1500 bytes is too big for MTU=1400 link

More PMTUD (3)

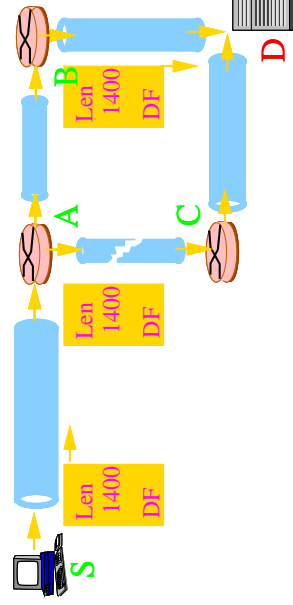
- ◇ A sends ICMP message
 - packet was discarded



- This tells S that MTU of Path is now 1400

More PMTUD (4)

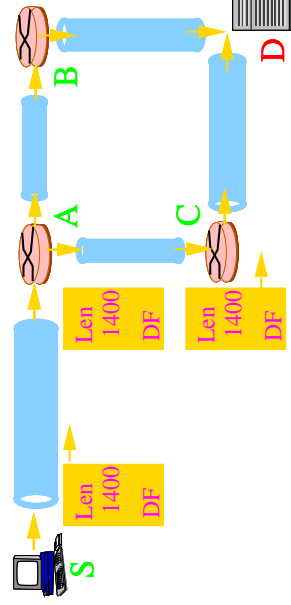
- ◇ S updates its record of the Path MTU to D
 - And ensures future packets
 - ▷ are no bigger than 1400 bytes



- This allows packets from S to reach D
 - ▷ without fragmentation

More PMTUD (5)

- ◇ Notice that when the A to C link is repaired
 - S still sends packets 1400 bytes



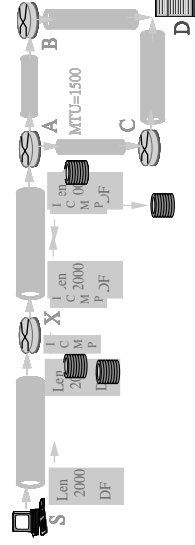
- Even though path now allows 1500 bytes
 - There is no Packet too small ICMP message!

Final PMTUD

- ◇ To allow for possibility
 - that PMTU has increased
- ◇ Sender occasionally
 - sends slightly bigger packet
- ◇ If that fails
 - ▷ ICMP Too Big is returned
 - Then nothing has changed
 - ▷ Determined PMTU remains as found earlier
 - Or as rediscovered now
- ◇ If bigger packet reaches D
 - Then PMTU has increased
 - ▷ S remembers bigger packets work
 - ▷ Tries again with even bigger packet
 - Until local link MTU reached
 - Or until ICMP is returned
 - Now the new PMTU is known

PMTUD Operational Problem

- ◇ Some network operators dislike ICMP
 - ▷ Some ICMP packets used Denial of Service attacks
 - So they prohibit all incoming ICMP packets!
- ◇ Problem for PMTUD ...



PMTUD Operational Problem (2)

- ◇ No packets delivered
 - No packets larger than Path MTU
- ◇ No failure indication received
- ◇ Solutions
 - Don't set DF
 - ▷ Allow fragmentation
 - No PMTUD
 - Black Hole Detection
 - ▷ Observe lack of replies
 - No replies to large packets
 - ▷ Guess ICMP filtering occurring
 - ▷ Reduce calculated PMTU
 - Try again with smaller packet
 - Not perfect - but required
- ◇ Requires
 - ▷ Guessing suitable smaller size
 - Also required for older ICMP messages