

60th IETF, PMTUD WG:

**Path MTU Discovery Using Options
draft-welzl-pmtud-options-01.txt**

Michael Welzl

<http://www.welzl.at> , michael.welzl@uibk.ac.at

**Distributed and Parallel Systems Group
Institute of Computer Science
University of Innsbruck, Austria**

Motivation

- In the end, (any kind of) PMTUD always loses a packet
 - it would be nice to avoid this
- Also, PMTUD should converge fast
- I am in favor of performance related signaling like ECN and XCP...
 - no “ECN flag” for PMTUD up to now (to avoid loss)
 - no “XCP” for PMTUD up to now (to converge faster)
- Proposal: add such signaling

How it works

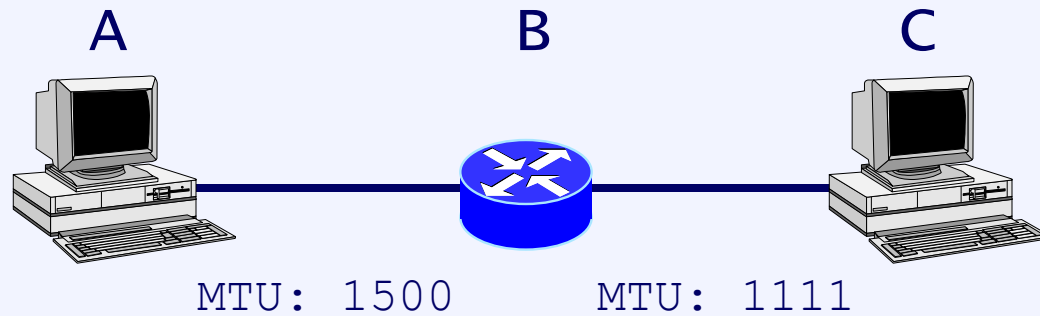
- Before doing (no matter which) PMTUD, include “Probe MTU” IP option
 - Initialized with MTU of outgoing link
 - Updated by routers if MTU of incoming or outgoing link is smaller
 - “TTL-Check” field decremented by each “Probe MTU” capable router: used to determine if all routers were involved
- Receiver feeds back result to source
 - either at IP layer (not recommended) or at packetization layer (specified for TCP, SCTP and DCCP, with IPv4 and IPv6)
- Sender reacts to feedback
 - Information complete (from TTL-Check): terminate immediately
 - Information incomplete: use as upper limit (i.e. starting point for RFC1191 PMTUD or to terminate PLPMTUD)

Potential benefits

- No loss, faster convergence
 - if lucky (result = PMTU)
- Less ICMP packets: less traffic, no risk of lost ICMP packet, reduced processing overhead for routers with small MTU
 - if lucky (result = PMTU)
- Less effort for tunnel endpoints (simply copy the option)
 - if supported by routers within a tunnel
- May circumvent Black Hole Detection
 - if upper limit from PMTU-Options < value that would cause troubles due to routers that don't send "Fragmentation needed"

Most beneficial for such routers, which are most beneficial for end points

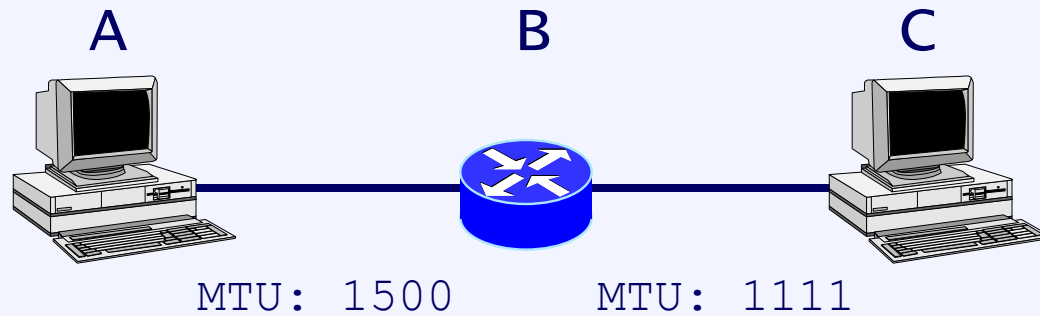
Example trace without PMTU-Options



3 normal PCs
 100 Mbit/s Ethernet
 Linux Kernel version 2.4.26
 RedHat 9.0 standard installation
 TCP file transfer with netcat

Nr.	Size	Sender	Receiver	Packet information
...				
6	1500	A	C	lost
7	576	B	A	ICMP Dest. unreachable
8	1500	A	C	lost
9	576	B	A	ICMP Dest. unreachable
10	1111	A	C	

Example trace with PMTU-Options



3 normal PCs
 100 Mbit/s Ethernet
 Linux Kernel version 2.4.26
 RedHat 9.0 standard installation
 TCP file transfer with netcat
kernel patch installed!

Nr.	IP Size	Sender	Receiver	Packet information
3	68	A	C	pmtu-ask 1500
8	60	C	A	pmtu-reply 1111
	...			
40	1111	A	C	
41	1111	A	C	

Problems with IP Options

- Slow Path processing
- Some routers drop these packets
- Series of measurement studies carried out with NOP IP Option... data from 2004 (100 pings of each type per host, 1 ping per second):
 - 12889 different hosts addresses, 14508 different router addresses
 - path lengths ranging from approx. 5 to 35 (majority around 15-25)
 - 29.48% of hosts did not respond when there was an IP option
 - average additional delay of 26.5% of a RTT
- Unknown problems
 - processing effort for routers
 - delay / drop results when a long series of packets carry options
 - Does Slow Path processing lead to reordering?

Deployment considerations

- Clearly not recommendable for all end-to-end TCP connections
- Also, security issues
 - lie about number of routers or send a MTU value that is too large: prevented by a Nonce
 - send a MTU value that is too small: cannot be prevented :-(
- Recommended for “special” scenarios only
 - detecting increased PMTU
 - tunnels
 - RTT-robust transport protocols
 - **Experimental** status envisioned

Patch, measurement results, future updates available from
<http://www.welzl.at/research/projects/ip-options/>